# CSCI567 Machine Learning (Fall 2024)

Prof. Dani Yogatama

University of Southern California

September 20, 2024

# Outline

# Acknowledgements

Not much math, a lot of empirical intuitions

## Acknowledgements

Not much math, a lot of empirical intuitions

The materials borrow heavily from the following sources:

- Stanford Course CS231n: `http://cs231n.stanford.edu/`
- Dr. Ian Goodfellow's lectures on deep learning: `http://deeplearningbook.org`

Both website provides tons of useful resources: notes, demos, videos, etc.
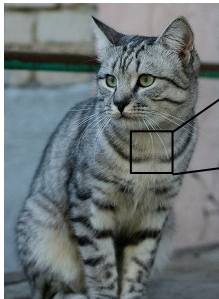
**Image Classification**: A core task in Computer Vision



This image by Nikita is
licensed under CC-BY 2.0

(assume given set of discrete labels)
{dog, cat, truck, plane, ...}

$\longrightarrow$  cat

**The Problem**: Semantic Gap



This image by Nikita is licensed under CC-BY 2.0

```
[[105 112 108 111 104  99 106  99  96 103 112 119 104  97  93  87]
 [ 91  98 102 106 104  79  98 103  99 105 123 136 110 105  94  85]
 [ 76  85  90 105 128 105  87  96  95  99 115 112 106 103  99  85]
 [ 99  81  81  93 120 131 127 100  95  98 102  99  96  93 101  94]
 [106  91  61  64  69  91  88  85 101 107 109  98  75  84  96  95]
 [114 108  85  55  55  69  64  54  64  87 112 129  98  74  84  91]
 [133 137 147 103  65  81  80  65  52  54  74  84 102  93  85  82]
 [128 137 144 140 109  95  86  70  62  65  63  63  60  73  86 101]
 [125 133 148 137 119 121 117  94  65  79  80  65  54  64  72  98]
 [127 125 131 147 133 127 126 131 111  96  89  75  61  64  72  84]
 [115 114 109 123 150 148 131 118 113 109 100  92  74  65  72  78]
 [ 89  93  90  97 108 147 131 118 113 114 113 109 106  95  77  80]
 [ 63  77  86  81  77  79 102 123 117 115 117 125 125 130 115  87]
 [ 62  65  82  89  78  71  80 101 124 126 119 101 107 114 131 119]
 [ 63  65  75  88  89  71  62  81 120 138 135 105  81  98 110 118]
 [ 87  65  71  87 106  95  69  45  76 130 126 107  92  94 105 112]
 [118  97  82  86 117 123 116  66  41  51  95  93  89  95 102 107]
 [164 146 112  80  82 120 124 104  76  48  45  66  88 101 102 109]
 [157 170 157 120  93  86 114 132 112  97  69  55  70  82  99  94]
 [130 128 134 161 139 100 109 118 121 134 114  87  65  53  69  86]
 [128 112  96 117 150 144 120 115 104 107 102  93  87  81  72  79]
 [123 107  96  86  83 112 153 149 122 109 104  75  80 107 112  99]
 [122 121 102  80  82  86  94 117 145 148 153 102  58  78  92 107]
 [122 164 148 103  71  56  78  83  93 103 119 139 102  61  69  84]]
```
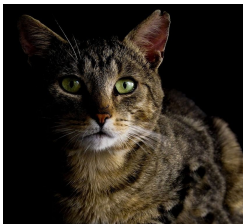
What the computer sees

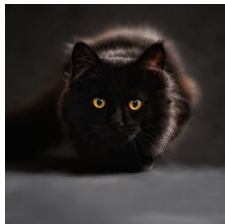An image is just a big grid of numbers between [0, 255]:
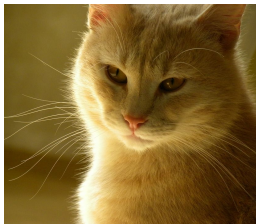
e.g. 800 x 600 x 3
(3 channels RGB)

**Challenges**: Viewpoint variation



All pixels change when
the camera moves!

**Challenges**: Illumination



This image is CC0 1.0 public domain   This image is CC0 1.0 public domain   This image is CC0 1.0 public domain   This image is CC0 1.0 public domain

**Challenges**: Deformation

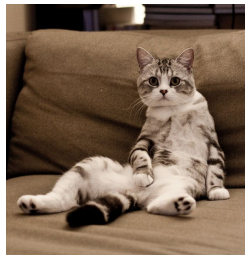**Challenges**: Occlusion



This image is CC0 1.0 public domain

This image is CC0 1.0 public domain

This image by jonsson is licensed under CC-BY 2.0

**Challenges**: Background Clutter

**Challenges**: Intraclass variation



This image is CC0 1.0 public domain

# Fundamental problems in vision

**The key challenge**
How to train a model that can tolerate all those variations?

# Fundamental problems in vision

**The key challenge**
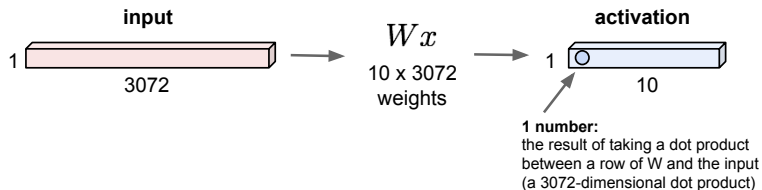How to train a model that can tolerate all those variations?

**Main ideas**

- need a lot of data that exhibits those variations

- need more specialized models to capture the invariance

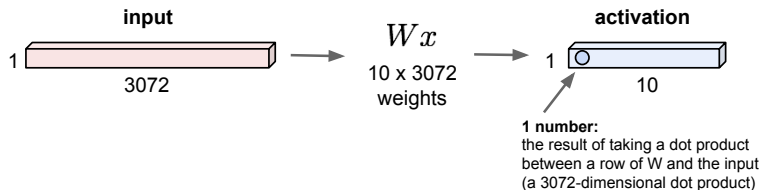# Issues of standard NN for image inputs

## Fully Connected Layer

32x32x3 image -> stretch to 3072 x 1



**input**

1 [                    ]
        3072

$Wx$
10 x 3072
weights

**activation**

1 [⊙                    ]
              10

**1 number:**
the result of taking a dot product
between a row of W and the input
(a 3072-dimensional dot product)

Fei-Fei Li & Justin Johnson & Serena Yeung          Lecture 5 - 27      April 18, 2017

# Issues of standard NN for image inputs

## Fully Connected Layer

32x32x3 image -> stretch to 3072 x 1



**input**

1 | 3072

$Wx$
10 x 3072
weights

**activation**

1 | 10

**1 number:**
the result of taking a dot product
between a row of W and the input
(a 3072-dimensional dot product)

Fei-Fei Li & Justin Johnson & Serena Yeung       Lecture 5 - 27     April 18, 2017

*Spatial structure is lost!*

# Solution: Convolutional Neural Net (ConvNet/CNN)

A special case of fully connected neural nets

# Solution: Convolutional Neural Net (ConvNet/CNN)

A special case of fully connected neural nets

- usually consist of **convolution layers**, ReLU layers, **pooling layers**, and regular fully connected layers

# Solution: Convolutional Neural Net (ConvNet/CNN)

A special case of fully connected neural nets

- usually consist of **convolution layers**, ReLU layers, **pooling layers**, and regular fully connected layers
- key idea: *learning from low-level to high-level features*

# Solution: Convolutional Neural Net (ConvNet/CNN)

A special case of fully connected neural nets

- usually consist of **convolution layers**, ReLU layers, **pooling layers**, and regular fully connected layers
- key idea: *learning from low-level to high-level features*

# Convolution layer

Arrange neurons as a **3D volume** naturally

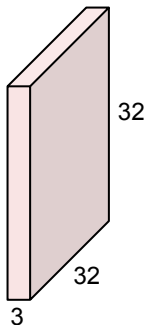## Convolution Layer

32x32x3 image -> preserve spatial structure



32 height

32 width

3 depth

Fei-Fei Li & Justin Johnson & Serena Yeung          Lecture 5 - 28       April 18, 2017
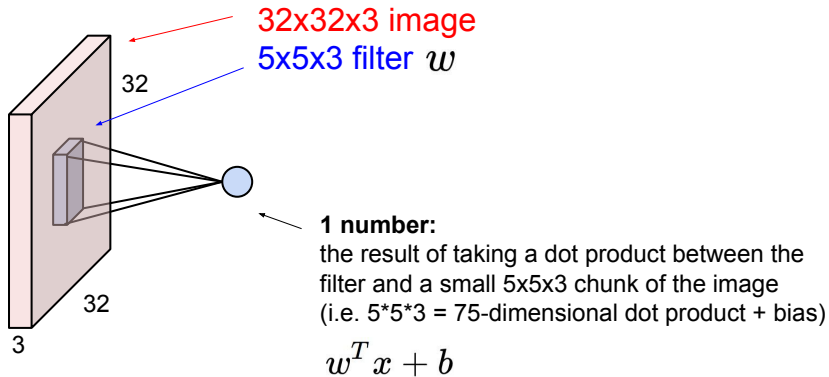
## Convolution

# Convolution Layer

32x32x3 image



32

32

3

5x5x3 filter

**Convolve** the filter with the image
i.e. "slide over the image spatially,
computing dot products"

# Convolution Layer

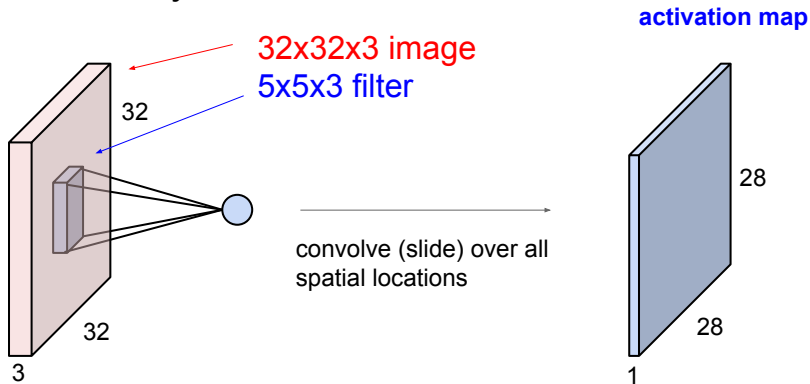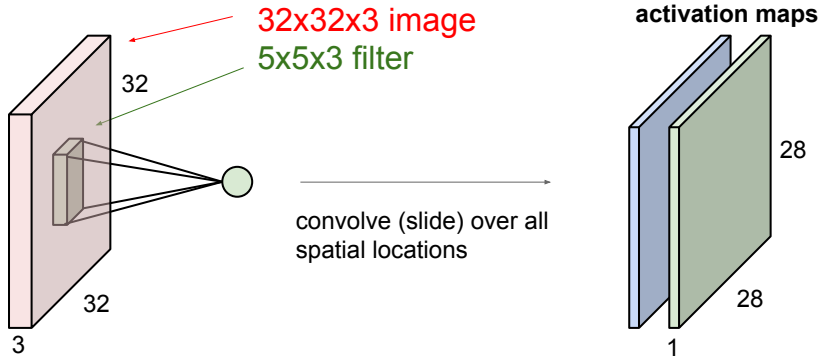**Filters always extend the full depth of the input volume**

32x32x3 image



5x5x3 filter

**Convolve** the filter with the image i.e. "slide over the image spatially, computing dot products"

32

32

3

# Convolution Layer



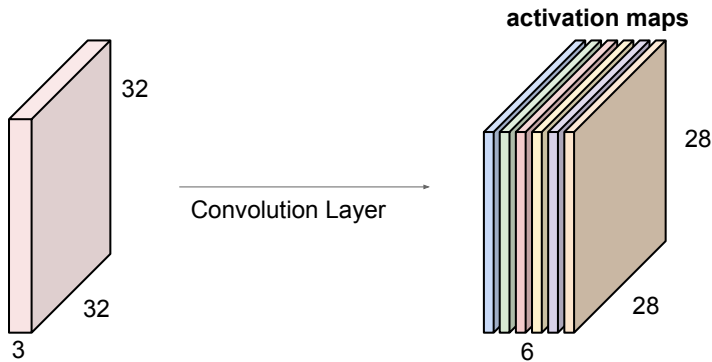32x32x3 image
5x5x3 filter $w$

32

32

3

**1 number:**
the result of taking a dot product between the filter and a small 5x5x3 chunk of the image (i.e. 5*5*3 = 75-dimensional dot product + bias)

$$w^T x + b$$

# Convolution Layer



32x32x3 image
5x5x3 filter

activation map

convolve (slide) over all spatial locations

Convolution Layer

consider a second, green filter

32x32x3 image
5x5x3 filter

32

32

3

convolve (slide) over all
spatial locations

**activation maps**
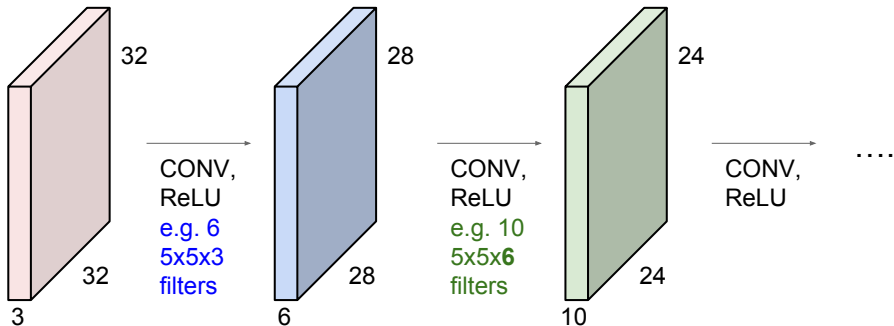
28

28

1

For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:



We stack these up to get a "new image" of size 28x28x6!

**Preview:** ConvNet is a sequence of Convolution Layers, interspersed with activation functions



32

32

3

CONV,
ReLU
e.g. 6
5x5x3
filters

28

28

6

**Preview:** ConvNet is a sequence of Convolutional Layers, interspersed with activation functions



CONV,
ReLU
e.g. 6
5x5x3
filters

CONV,
ReLU
e.g. 10
5x5x**6**
filters

CONV,
ReLU

....

# Why convolution makes sense?

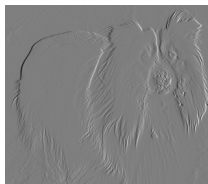Main idea: **if a filter is useful at one location, it should be useful at other locations.**

# Why convolution makes sense?

Main idea: **if a filter is useful at one location, it should be useful at other locations.**



A simple example why filtering is useful

Input

| 1 | -1 |

Kernel

Output

# Connection to fully connected NNs

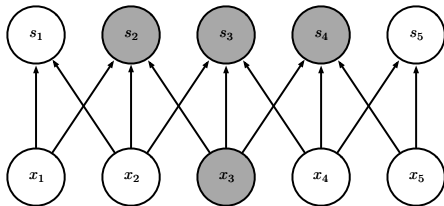A convolution layer is a special case of a fully connected layer:

# Connection to fully connected NNs

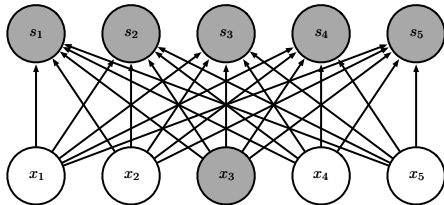A convolution layer is a special case of a fully connected layer:

- filter $=$ weights with **sparse connection**

# Local Receptive Field Leads to Sparse Connectivity (affects less)
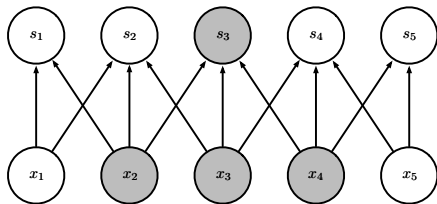
Sparse connections due to small convolution kernel

Dense connections

# Sparse connectivity: being affected by less

Sparse connections due to small convolution kernel
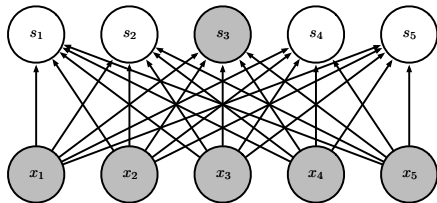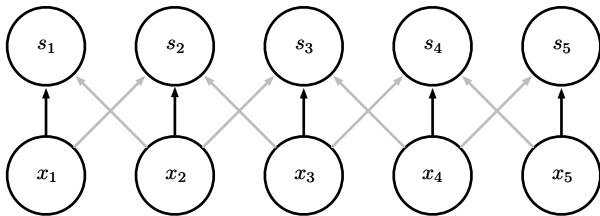


Dense connections

Figure 9.3

# Connection to fully connected NNs

A convolution layer is a special case of a fully connected layer:

- filter = weights with **sparse connection**
- **parameters sharing**

# Parameter Sharing

Convolution shares the same parameters across all spatial locations

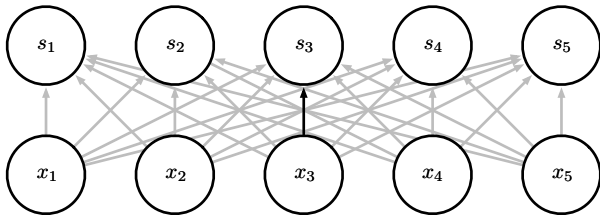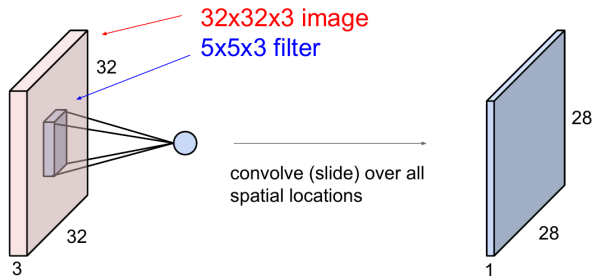Traditional matrix multiplication does not share any parameters



Figure 9.5

# Connection to fully connected NNs

A convolution layer is a special case of a fully connected layer:

- filter = weights with **sparse connection**
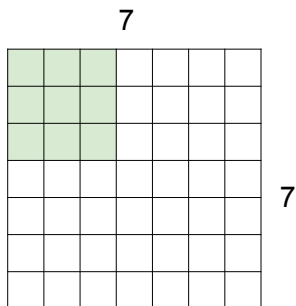- **parameters sharing**

*Much fewer parameters!* Example (ignore bias terms):

- FC: $(32 \times 32 \times 3) \times (28 \times 28) \approx 2.4M$
- CNN: $5 \times 5 \times 3 = 75$



32x32x3 image
5x5x3 filter

convolve (slide) over all spatial locations

# Spatial arrangement: stride and padding

A closer look at spatial dimensions:

7



7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:



7x7 input (spatially)
assume 3x3 filter
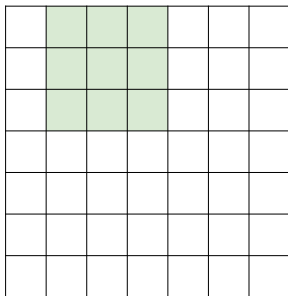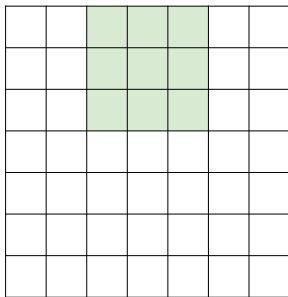
A closer look at spatial dimensions:



7x7 input (spatially)
assume 3x3 filter

A closer look at spatial dimensions:

7



7

7x7 input (spatially)
assume 3x3 filter
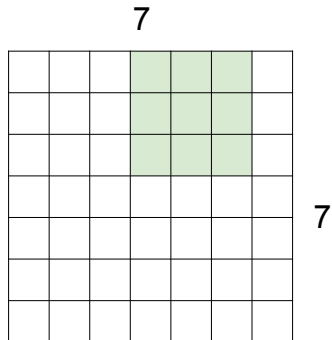
A closer look at spatial dimensions:
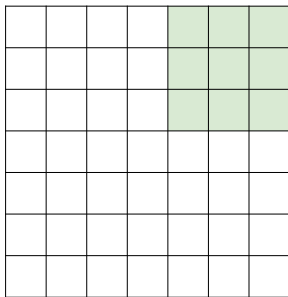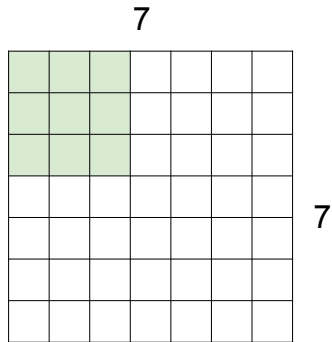
7



7x7 input (spatially)
assume 3x3 filter

**=> 5x5 output**

7

A closer look at spatial dimensions:

7



7

7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**

A closer look at spatial dimensions:



7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**

A closer look at spatial dimensions:



7x7 input (spatially)
assume 3x3 filter
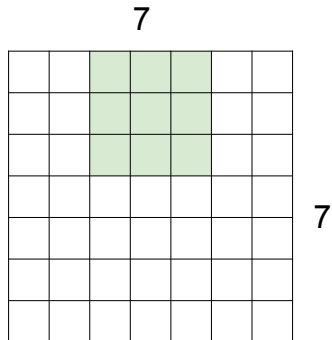applied **with stride 2
=> 3x3 output!**

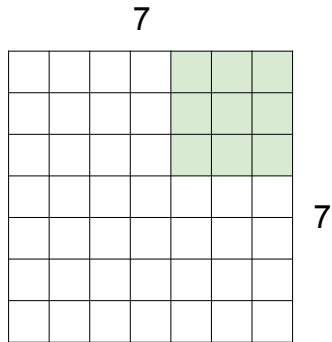A closer look at spatial dimensions:

7



7

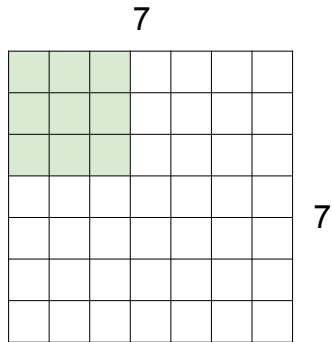7x7 input (spatially)
assume 3x3 filter
applied **with stride 3?**

A closer look at spatial dimensions:

7



7

7x7 input (spatially)
assume 3x3 filter
applied **with stride 3?**

**doesn't fit!**
cannot apply 3x3 filter on
7x7 input with stride 3.

Output size:
**(N - F) / stride + 1**

e.g. N = 7, F = 3:
stride 1 => (7 - 3)/1 + 1 = 5
stride 2 => (7 - 3)/2 + 1 = 3
stride 3 => (7 - 3)/3 + 1 = 2.33 :\

# In practice: Common to zero pad the border



e.g. input 7x7
**3x3** filter, applied with **stride 1**
**pad with 1 pixel** border => what is the output?

(recall:)
(N - F) / stride + 1

## In practice: Common to zero pad the border



| 0 | 0 | 0 | 0 | 0 | 0 | | | |
|---|---|---|---|---|---|---|---|---|
| 0 | | | | | | | | |
| 0 | | | | | | | | |
| 0 | | | | | | | | |
| 0 | | | | | | | | |
| | | | | | | | | |
| | | | | | | | | |
| | | | | | | | | |
| | | | | | | | | |

e.g. input 7x7
**3x3** filter, applied with **stride 1**
**pad with 1 pixel** border => what is the output?

**7x7 output!**

# In practice: Common to zero pad the border



e.g. input 7x7
**3x3** filter, applied with **stride 1**
**pad with 1 pixel** border => what is the output?

**7x7 output!**
in general, common to see CONV layers with stride 1, filters of size FxF, and zero-padding with (F-1)/2. (will preserve size spatially)
e.g. F = 3 => zero pad with 1
     F = 5 => zero pad with 2
     F = 7 => zero pad with 3

**Remember back to…**
E.g. 32x32 input convolved repeatedly with 5x5 filters shrinks volumes spatially!
(32 -> 28 -> 24 ...). Shrinking too fast is not good, doesn't work well.



32

28

24

CONV,
ReLU
e.g. 6
5x5x3
filters

CONV,
ReLU
e.g. 10
5x5x**6**
filters

CONV,
ReLU

….

32

28

24

3

6

10

# Summary for convolution layer

**Input**: a volume of size $W_1 \times H_1 \times D_1$

# Summary for convolution layer

**Input**: a volume of size $W_1 \times H_1 \times D_1$

**Hyperparameters**:

- $K$ filters of size $F \times F$
- stride $S$
- amount of zero padding $P$ (for one side)

# Summary for convolution layer

**Input**: a volume of size $W_1 \times H_1 \times D_1$

**Hyperparameters**:

- $K$ filters of size $F \times F$
- stride $S$
- amount of zero padding $P$ (for one side)

**Output**: a volume of size $W_2 \times H_2 \times D_2$ where

- $W_2 =$
- $H_2 =$
- $D_2 =$

# Summary for convolution layer

**Input**: a volume of size $W_1 \times H_1 \times D_1$

**Hyperparameters**:

- $K$ filters of size $F \times F$
- stride $S$
- amount of zero padding $P$ (for one side)

**Output**: a volume of size $W_2 \times H_2 \times D_2$ where

- $W_2 = (W_1 + 2P - F)/S + 1$
- $H_2 =$
- $D_2 =$

# Summary for convolution layer

**Input**: a volume of size $W_1 \times H_1 \times D_1$

**Hyperparameters**:

- $K$ filters of size $F \times F$
- stride $S$
- amount of zero padding $P$ (for one side)

**Output**: a volume of size $W_2 \times H_2 \times D_2$ where

- $W_2 = (W_1 + 2P - F)/S + 1$
- $H_2 = (H_1 + 2P - F)/S + 1$
- $D_2 =$

# Summary for convolution layer

**Input**: a volume of size $W_1 \times H_1 \times D_1$

**Hyperparameters**:

- $K$ filters of size $F \times F$
- stride $S$
- amount of zero padding $P$ (for one side)

**Output**: a volume of size $W_2 \times H_2 \times D_2$ where

- $W_2 = (W_1 + 2P - F)/S + 1$
- $H_2 = (H_1 + 2P - F)/S + 1$
- $D_2 = K$

# Summary for convolution layer

**Input**: a volume of size $W_1 \times H_1 \times D_1$

**Hyperparameters**:

- $K$ filters of size $F \times F$
- stride $S$
- amount of zero padding $P$ (for one side)

**Output**: a volume of size $W_2 \times H_2 \times D_2$ where

- $W_2 = (W_1 + 2P - F)/S + 1$
- $H_2 = (H_1 + 2P - F)/S + 1$
- $D_2 = K$

**#parameters**: $(F \times F \times D_1 + 1) \times K$ weights

# Summary for convolution layer

**Input**: a volume of size $W_1 \times H_1 \times D_1$

**Hyperparameters**:

- $K$ filters of size $F \times F$
- stride $S$
- amount of zero padding $P$ (for one side)

**Output**: a volume of size $W_2 \times H_2 \times D_2$ where

- $W_2 = (W_1 + 2P - F)/S + 1$
- $H_2 = (H_1 + 2P - F)/S + 1$
- $D_2 = K$

**#parameters**: $(F \times F \times D_1 + 1) \times K$ weights
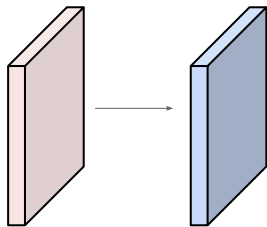
**Common setting**: $F = 3, S = P = 1$

Examples time:

Input volume: **32x32x3**
10 5x5 filters with stride 1, pad 2

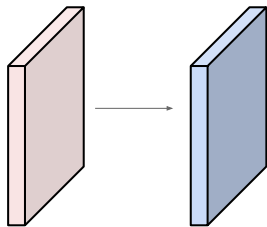Output volume size: ?

Examples time:

Input volume: **32x32x3**
10 5x5 filters with stride 1, pad 2

Output volume size:
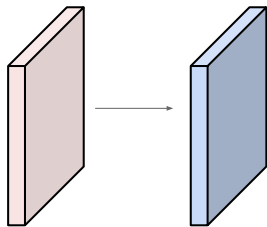(32+2*2-5)/1+1 = 32 spatially, so
**32x32x10**

Examples time:

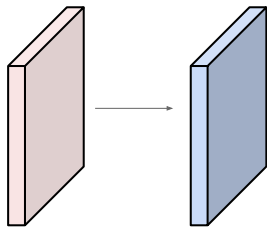Input volume: **32x32x3**
10 5x5 filters with stride 1, pad 2

Number of parameters in this layer?

Examples time:

Input volume: **32x32x3**
10 5x5 filters with stride 1, pad 2



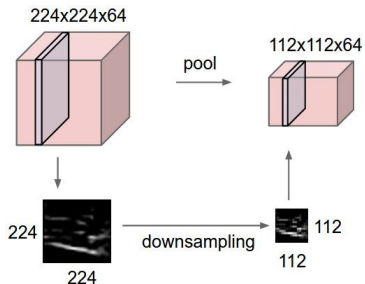Number of parameters in this layer?
each filter has 5*5*3 + 1 = 76 params  (+1 for bias)
=> 76*10 = **760**

# Another element: pooling

## Pooling layer
- makes the representations smaller and more manageable
- operates over each activation map independently:



Fei-Fei Li & Justin Johnson & Serena Yeung          Lecture 5 - 72      April 18, 2017

# Pooling

Similar to a filter, except

- depth is always 1

# Pooling

Similar to a filter, except

- depth is always 1
- different operations: average, L2-norm, max

# Pooling

Similar to a filter, except

- depth is always 1
- different operations: average, L2-norm, max
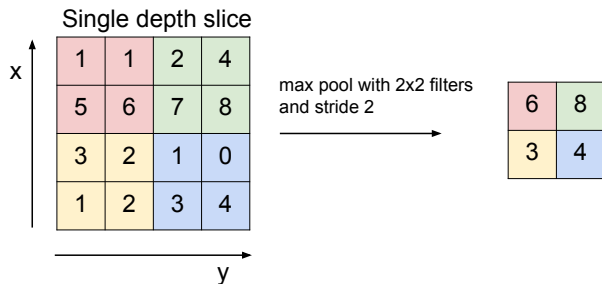- no parameters to be learned

# Pooling

Similar to a filter, except

- depth is always 1
- different operations: average, L2-norm, max
- no parameters to be learned

**Max pooling** with $2 \times 2$ filter and stride 2 is very common

## MAX POOLING

# Putting everything together

**Typical architecture for CNNs:**

Input → [[Conv → ReLU]*N → Pool?]*M → [FC → ReLU]*Q → FC

# Putting everything together

**Typical architecture for CNNs:**

Input → [[Conv → ReLU]*N → Pool?]*M → [FC → ReLU]*Q → FC

Common choices: $N \leq 5, Q \leq 2$, $M$ is large

# Putting everything together

**Typical architecture for CNNs:**

Input $\rightarrow$ [[Conv $\rightarrow$ ReLU]*N $\rightarrow$ Pool?]*M $\rightarrow$ [FC $\rightarrow$ ReLU]*Q $\rightarrow$ FC

Common choices: $N \leq 5, Q \leq 2$, $M$ is large

**Well-known CNNs**: LeNet, AlexNet, ZF Net, GoogLeNet, VGGNet, etc.

All achieve excellent performance on image classification tasks.

# How to train a CNN?

*How do we learn the filters/weights?*

# How to train a CNN?

*How do we learn the filters/weights?*

Essentially the same as FC NNs: apply **SGD/backpropagation**