

# CSCI567 Machine Learning (Fall 2018)

Prof. Haipeng Luo

U of Southern California

Sep 19, 2018

September 19, 2018 1 / 50

## Outline

- 1 Review of last lecture
- 2 Convolutional neural networks
- 3 Kernel methods

September 19, 2018 3 / 50

## Administration

HW 1 grading will be completed next week

HW 2 to be released very soon

- deadline is 10/07
- but solutions for written part will be released on 09/27

Midterm (10/03)

- location is finalized: THH 101 and THH 201
- coverage on Lecture 6 will be minimum
- you can start preparing after this lecture

September 19, 2018 2 / 50

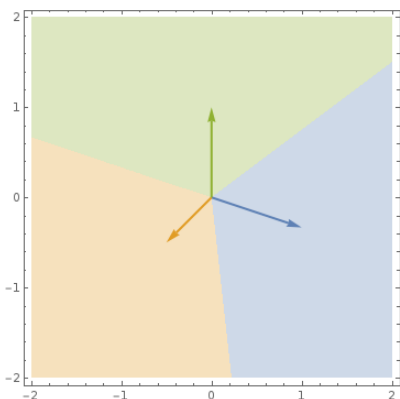
Review of last lecture

## Outline

- 1 Review of last lecture
- 2 Convolutional neural networks
- 3 Kernel methods

September 19, 2018 4 / 50

## Linear models: from binary to multiclass



$$\begin{aligned} \mathbf{w}_1 &= (1, -\frac{1}{3}) \\ \mathbf{w}_2 &= (-\frac{1}{2}, -\frac{1}{2}) \\ \mathbf{w}_3 &= (0, 1) \end{aligned}$$

- Blue class:  
 $\{\mathbf{x} : 1 = \operatorname{argmax}_k \mathbf{w}_k^T \mathbf{x}\}$
- Orange class:  
 $\{\mathbf{x} : 2 = \operatorname{argmax}_k \mathbf{w}_k^T \mathbf{x}\}$
- Green class:  
 $\{\mathbf{x} : 3 = \operatorname{argmax}_k \mathbf{w}_k^T \mathbf{x}\}$

$$\mathcal{F} = \left\{ f(\mathbf{x}) = \operatorname{argmax}_{k \in [C]} \mathbf{w}_k^T \mathbf{x} \mid \mathbf{w}_1, \dots, \mathbf{w}_C \in \mathbb{R}^D \right\}$$

## MLE = minimizing cross-entropy loss

Maximize probability of see labels  $y_1, \dots, y_N$  given  $\mathbf{x}_1, \dots, \mathbf{x}_N$

$$P(\mathbf{W}) = \prod_{n=1}^N \mathbb{P}(y_n \mid \mathbf{x}_n; \mathbf{W}) = \prod_{n=1}^N \frac{e^{\mathbf{w}_{y_n}^T \mathbf{x}_n}}{\sum_{k \in [C]} e^{\mathbf{w}_k^T \mathbf{x}_n}}$$

By taking **negative log**, this is equivalent to minimizing

$$F(\mathbf{W}) = \sum_{n=1}^N \ln \left( \frac{\sum_{k \in [C]} e^{\mathbf{w}_k^T \mathbf{x}_n}}{e^{\mathbf{w}_{y_n}^T \mathbf{x}_n}} \right) = \sum_{n=1}^N \ln \left( 1 + \sum_{k \neq y_n} e^{(\mathbf{w}_k - \mathbf{w}_{y_n})^T \mathbf{x}_n} \right)$$

This is the *multiclass logistic loss*, a.k.a *cross-entropy loss*.

## Comparisons of multiclass-to-binary reductions

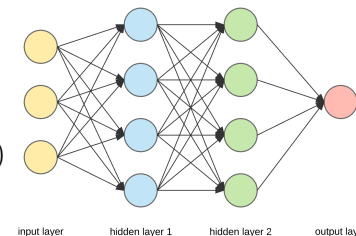
In big O notation,

Reduction	#training points	test time	Idea
OvA	CN	C	is class $k$ or not?
OvO	CN	$C^2$	is class $k$ or class $k'$ ?
ECOC	LN	L	is bit $b$ on or off?
Tree	$(\log_2 C)N$	$\log_2 C$	belong to which half of the label set?

## Math formulation of neural nets

An L-layer neural net can be written as

$$\mathbf{f}(\mathbf{x}) = \mathbf{h}_L(\mathbf{W}_L \mathbf{h}_{L-1}(\mathbf{W}_{L-1} \cdots \mathbf{h}_1(\mathbf{W}_1 \mathbf{x})))$$



To ease notation, for a given input  $\mathbf{x}$ , define recursively

$$\mathbf{o}_0 = \mathbf{x}, \quad \mathbf{a}_\ell = \mathbf{W}_\ell \mathbf{o}_{\ell-1}, \quad \mathbf{o}_\ell = \mathbf{h}_\ell(\mathbf{a}_\ell) \quad (\ell = 1, \dots, L)$$

where

- $\mathbf{W}_\ell \in \mathbb{R}^{D_\ell \times D_{\ell-1}}$  is the weights for layer  $\ell$
- $D_0 = D, D_1, \dots, D_L$  are numbers of neurons at each layer
- $\mathbf{a}_\ell \in \mathbb{R}^{D_\ell}$  is input to layer  $\ell$
- $\mathbf{o}_\ell \in \mathbb{R}^{D_\ell}$  is output to layer  $\ell$
- $\mathbf{h} : \mathbb{R}^{D_\ell} \rightarrow \mathbb{R}^{D_\ell}$  is activation functions at layer  $\ell$

## Backprop = SGD for neural nets

The **backpropagation** algorithm (**Backprop**)

Initialize  $\mathbf{W}_1, \dots, \mathbf{W}_L$  (all  $\mathbf{0}$  or randomly). Repeat:

- 1 randomly pick one data point  $n \in [N]$
- 2 **forward propagation**: for each layer  $\ell = 1, \dots, L$ 
  - compute  $\mathbf{a}_\ell = \mathbf{W}_\ell \mathbf{o}_{\ell-1}$  and  $\mathbf{o}_\ell = \mathbf{h}_\ell(\mathbf{a}_\ell)$  ( $\mathbf{o}_0 = \mathbf{x}_n$ )
- 3 **backward propagation**: for each  $\ell = L, \dots, 1$ 
  - compute

$$\frac{\partial \mathcal{E}_n}{\partial \mathbf{a}_\ell} = \begin{cases} \left( \mathbf{W}_{\ell+1}^T \frac{\partial \mathcal{E}_n}{\partial \mathbf{a}_{\ell+1}} \right) \circ \mathbf{h}'_\ell(\mathbf{a}_\ell) & \text{if } \ell < L \\ 2(\mathbf{h}_L(\mathbf{a}_L) - \mathbf{y}_n) \circ \mathbf{h}'_L(\mathbf{a}_L) & \text{else} \end{cases}$$

- update weights

$$\mathbf{W}_\ell \leftarrow \mathbf{W}_\ell - \eta \frac{\partial \mathcal{E}_n}{\partial \mathbf{W}_\ell} = \mathbf{W}_\ell - \eta \frac{\partial \mathcal{E}_n}{\partial \mathbf{a}_\ell} \mathbf{o}_{\ell-1}^T$$

*Think about how to do the last two steps properly!*

## Outline

- 1 Review of last lecture
- 2 Convolutional neural networks
  - Motivation
  - Architecture
- 3 Kernel methods

## 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.

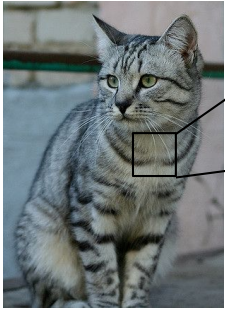
## Image Classification: A core task in Computer Vision



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

→ cat

## The Problem: Semantic Gap



This image by Nikita is licensed under CC-BY 2.0

```

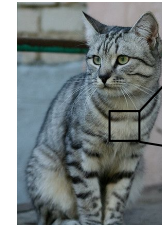
[[185 112 188 111 184 99 186 99 96 183 112 110 184 97 93 87]
 [ 91 98 182 186 184 79 98 183 99 185 113 116 118 185 94 85]
 [ 78 85 98 185 188 180 87 98 99 113 112 180 183 99 85]
 [ 89 81 81 93 128 131 127 188 95 98 182 99 96 93 181 84]
 [ 88 81 81 84 80 81 88 85 181 189 98 75 84 96 95]
 [114 188 85 55 55 85 84 84 84 87 112 129 98 74 84 91]
 [122 127 147 183 85 81 88 82 54 74 84 182 93 85 82]
 [128 127 144 180 85 88 82 65 63 63 68 73 86 83]
 [125 128 148 137 110 111 84 85 79 88 65 64 64 72 88]
 [127 125 147 137 110 111 86 89 75 64 72 88]
 [115 114 188 123 158 148 111 118 113 189 188 92 74 65 72 78]
 [ 89 93 98 87 188 147 111 118 113 114 113 189 186 95 77 88]
 [ 83 77 86 81 77 79 182 123 117 115 117 125 125 138 115 87]
 [ 82 85 82 89 78 71 88 181 124 128 118 181 114 113 113]
 [ 83 85 75 88 89 71 82 81 128 138 135 185 81 88 118 118]
 [ 87 85 71 87 186 95 69 65 78 138 138 187 82 84 182 113]
 [118 97 82 86 117 123 116 86 41 51 95 93 89 95 182 187]
 [184 148 132 88 82 128 124 78 68 65 68 81 182 189]
 [157 178 127 128 93 86 114 112 112 97 69 55 78 82 99 84]
 [138 128 134 141 139 188 180 118 121 114 114 87 65 69 86]
 [128 112 96 117 158 144 128 115 184 187 182 93 87 81 72 78]
 [123 187 86 86 113 123 149 122 188 184 78 88 112 95]
 [122 121 182 88 82 86 94 117 145 148 153 182 58 78 92 187]
 [122 144 148 183 71 56 78 83 93 183 118 128 82 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



```

[[185 112 188 111 184 99 186 99 96 183 112 110 184 97 93 87]
 [ 91 98 182 186 184 79 98 183 99 185 113 116 118 185 94 85]
 [ 78 85 98 185 188 180 87 98 99 113 112 180 183 99 85]
 [ 89 81 81 93 128 131 127 188 95 98 182 99 96 93 181 84]
 [ 88 81 81 84 80 81 88 85 181 189 98 75 84 96 95]
 [114 188 85 55 55 85 84 84 84 87 112 129 98 74 84 91]
 [122 127 147 183 85 81 88 82 54 74 84 182 93 85 82]
 [128 127 144 180 85 88 82 65 63 63 68 73 86 83]
 [125 128 148 137 110 111 84 85 79 88 65 64 64 72 88]
 [127 125 147 137 110 111 86 89 75 64 72 88]
 [115 114 188 123 158 148 111 118 113 189 188 92 74 65 72 78]
 [ 89 93 98 87 188 147 111 118 113 114 113 189 186 95 77 88]
 [ 83 77 86 81 77 79 182 123 117 115 117 125 125 138 115 87]
 [ 82 85 82 89 78 71 88 181 124 128 118 181 114 113 113]
 [ 83 85 75 88 89 71 82 81 128 138 135 185 81 88 118 118]
 [ 87 85 71 87 186 95 69 65 78 138 138 187 82 84 182 113]
 [118 97 82 86 117 123 116 86 41 51 95 93 89 95 182 187]
 [184 148 132 88 82 128 124 78 68 65 68 81 182 189]
 [157 178 127 128 93 86 114 112 112 97 69 55 78 82 99 84]
 [138 128 134 141 139 188 180 118 121 114 114 87 65 69 86]
 [128 112 96 117 158 144 128 115 184 187 182 93 87 81 72 78]
 [123 187 86 86 113 123 149 122 188 184 78 88 112 95]
 [122 121 182 88 82 86 94 117 145 148 153 182 58 78 92 187]
 [122 144 148 183 71 56 78 83 93 183 118 128 82 61 69 84]]
    
```

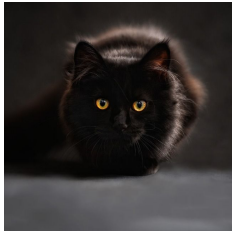
All pixels change when the camera moves!

This image by Nikita is licensed under CC-BY 2.0

## Challenges: Illumination



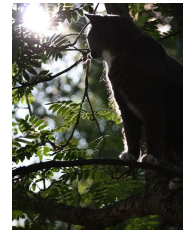
This image is CC0 1.0 public domain



This image is CC0 1.0 public domain

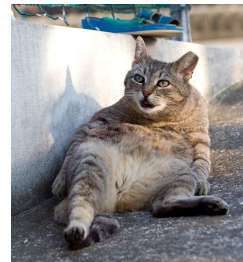


This image is CC0 1.0 public domain



This image is CC0 1.0 public domain

## Challenges: Deformation



This image by Umberto Salvagnin is licensed under CC-BY 2.0



This image by Umberto Salvagnin is licensed under CC-BY 2.0



This image by sara boax is licensed under CC-BY 2.0



This image by Tom Thair is licensed under CC-BY 2.0

## Challenges: Occlusion



This image is CC0.1.0 public domain



This image is CC0.1.0 public domain



This image by jomston is licensed under CC-BY 2.0

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 2 - 11

April 6, 2017

## Challenges: Background Clutter



This image is CC0.1.0 public domain



This image is CC0.1.0 public domain

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 2 - 12

April 6, 2017

## Challenges: Intra-class variation



This image is CC0.1.0 public domain

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 2 - 13

April 6, 2017

Convolutional neural networks Motivation

## Fundamental problems in vision

### The key challenge

How to train a model that can tolerate all those variations?

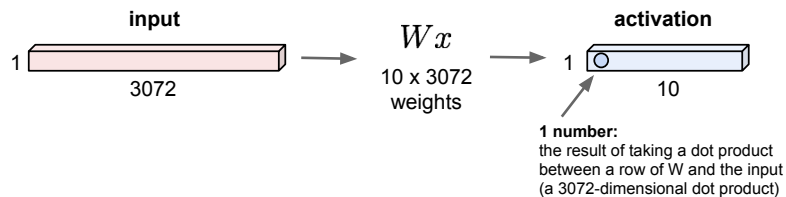
### Main ideas

- 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



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

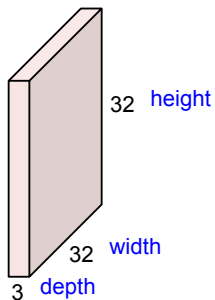
*Spatial structure is lost!*

## Convolution layer

Arrange neurons as a 3D volume naturally

### Convolution Layer

32x32x3 image -> preserve spatial structure

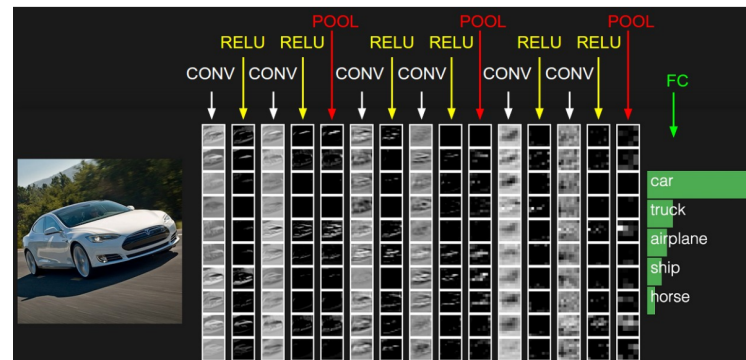


Fei-Fei Li & Justin Johnson & Serena Yeung Lecture 5 - 28 April 18, 2017

## Solution: Convolutional Neural Net (ConvNet/CNN)

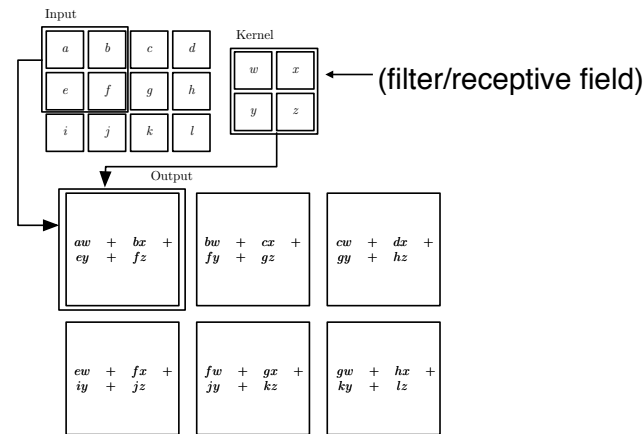
A special case of fully connected neural nets

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



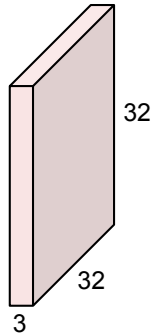
## Convolution

### 2D Convolution



# Convolution Layer

32x32x3 image



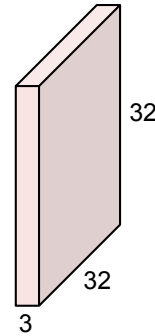
5x5x3 filter



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

# Convolution Layer

32x32x3 image



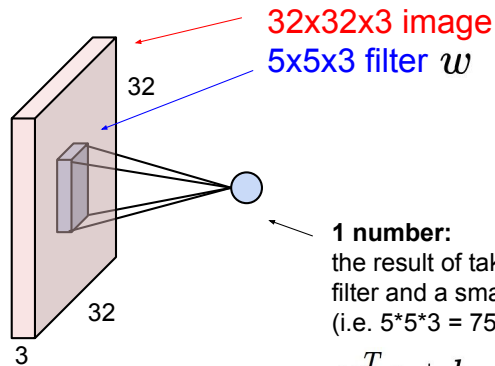
5x5x3 filter



Filters always extend the full  
depth of the input volume

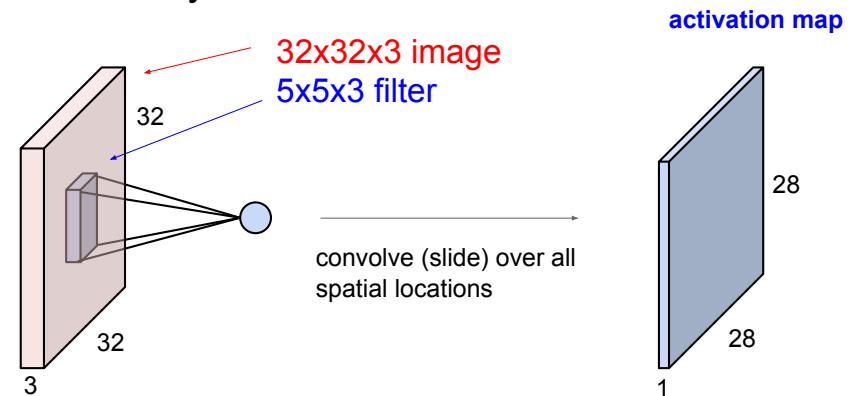
**Convolve** the filter with the image  
i.e. "slide over the image spatially,  
computing dot products"

# Convolution Layer



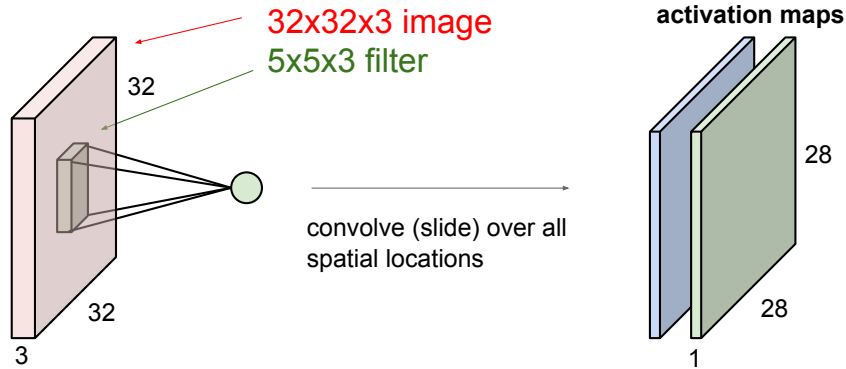
**1 number:**  
the result of taking a dot product between the  
filter and a small 5x5x3 chunk of the image  
(i.e.  $5*5*3 = 75$ -dimensional dot product + bias)  
 $w^T x + b$

# Convolution Layer

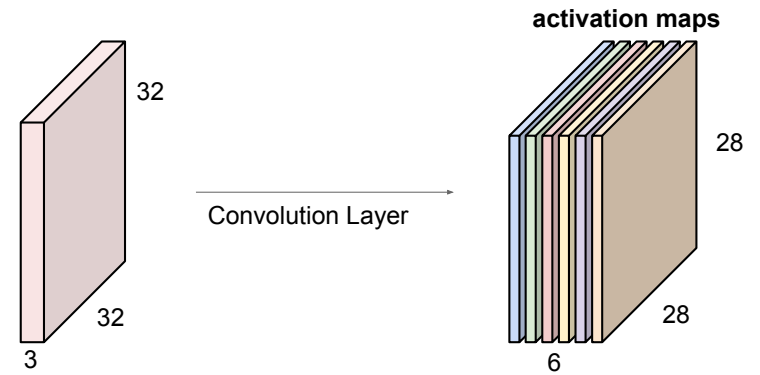


# Convolution Layer

consider a second, **green** filter

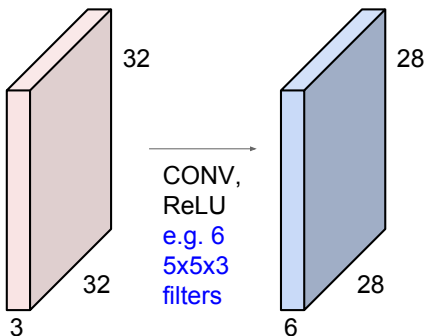


For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

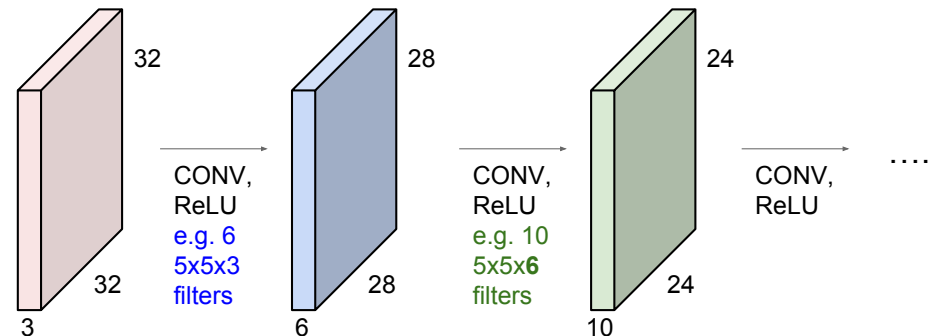


We stack these up to get a "new image" of size 28x28x6!

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



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

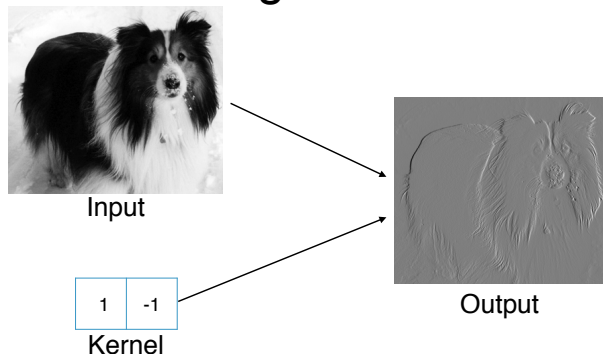




## Why convolution makes sense?

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

### A simple example why filtering is useful



September 19, 2018 17 / 50

## Connection to fully connected NNs

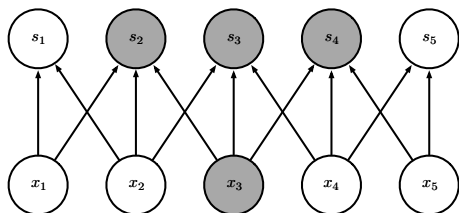
A convolution layer is a special case of a fully connected layer:

- filter = weights with **sparse connection**

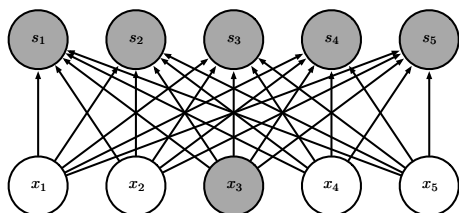
September 19, 2018 18 / 50

## Local Receptive Field Leads to Sparse Connectivity (affects less)

Sparse connections due to small convolution kernel



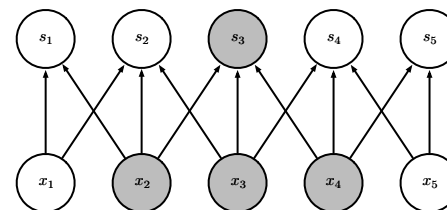
Dense connections



(Goodfellow 2016)

## Sparse connectivity: being affected by less

Sparse connections due to small convolution kernel



Dense connections

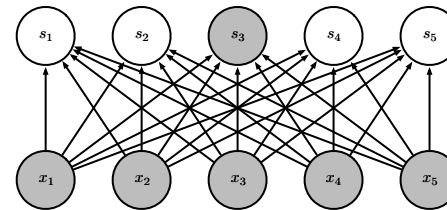


Figure 9.3

(Goodfellow 2016)

## 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

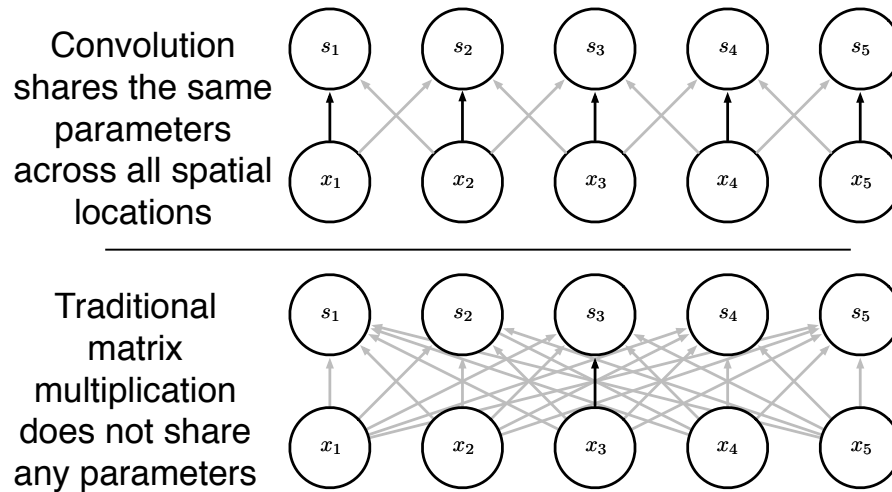


Figure 9.5

(Goodfellow 2016)

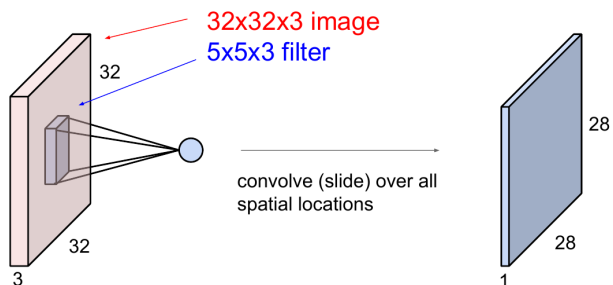
## Connection to fully connected NNs

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

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

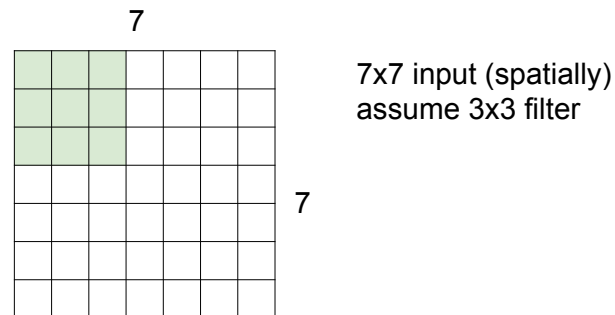
*Much less parameters!* Example (ignore bias terms):

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

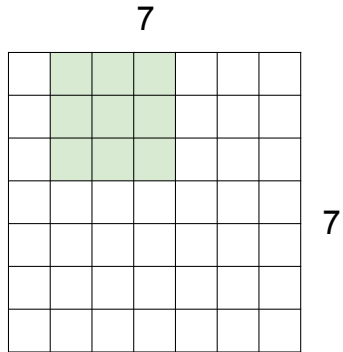


## Spatial arrangement: stride and padding

A closer look at spatial dimensions:

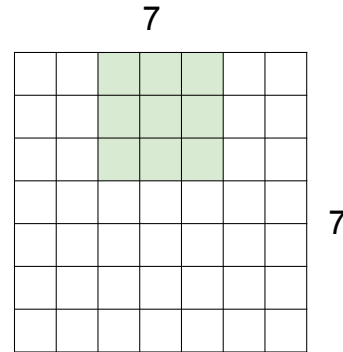


A closer look at spatial dimensions:



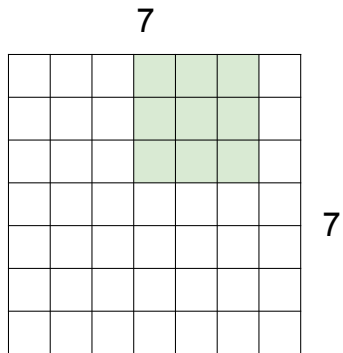
7x7 input (spatially)  
assume 3x3 filter

A closer look at spatial dimensions:



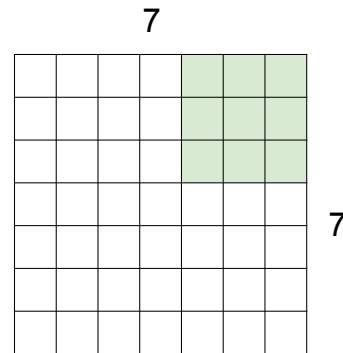
7x7 input (spatially)  
assume 3x3 filter

A closer look at spatial dimensions:



7x7 input (spatially)  
assume 3x3 filter

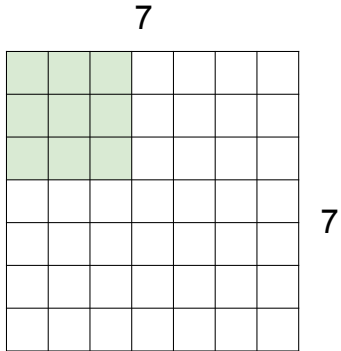
A closer look at spatial dimensions:



7x7 input (spatially)  
assume 3x3 filter

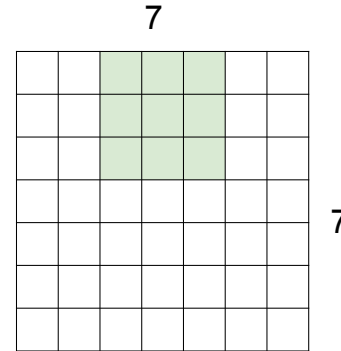
**=> 5x5 output**

A closer look at spatial dimensions:



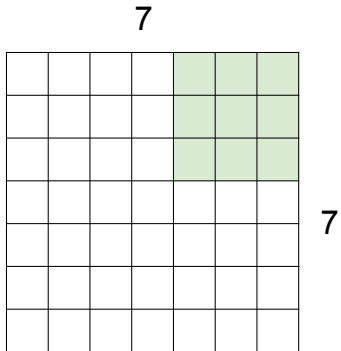
7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 2**

A closer look at spatial dimensions:



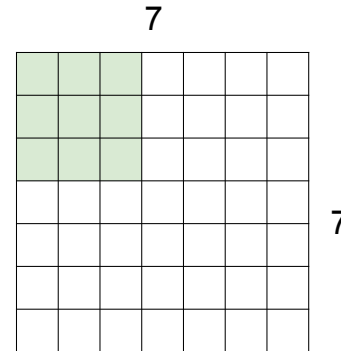
7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 2**

A closer look at spatial dimensions:



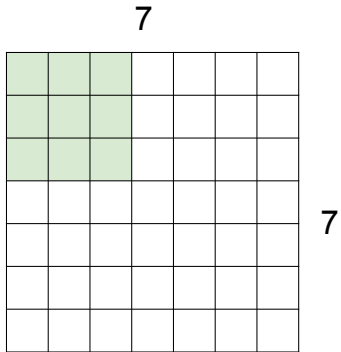
7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 2**  
**=> 3x3 output!**

A closer look at spatial dimensions:



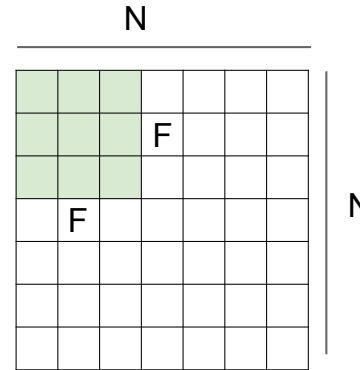
7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 3?**

A closer look at spatial dimensions:



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

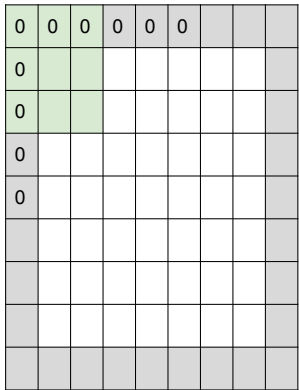
**doesn't fit!**  
cannot apply 3x3 filter on  
7x7 input with stride 3.



Output size:  
 **$(N - F) / \text{stride} + 1$**

e.g.  $N = 7, F = 3$ :  
stride 1  $\Rightarrow (7 - 3) / 1 + 1 = 5$   
stride 2  $\Rightarrow (7 - 3) / 2 + 1 = 3$   
stride 3  $\Rightarrow (7 - 3) / 3 + 1 = 2.33 \text{ :|}$

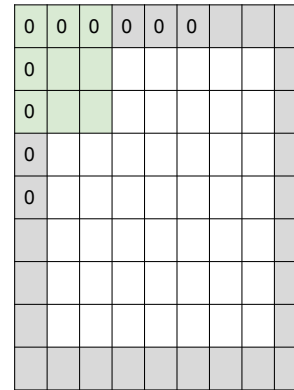
In practice: Common to zero pad the border



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

(recall:)  
 **$(N - F) / \text{stride} + 1$**

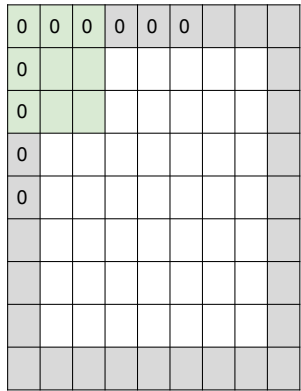
In practice: Common to zero pad the border



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

**7x7 output!**

## In practice: Common to zero pad the border

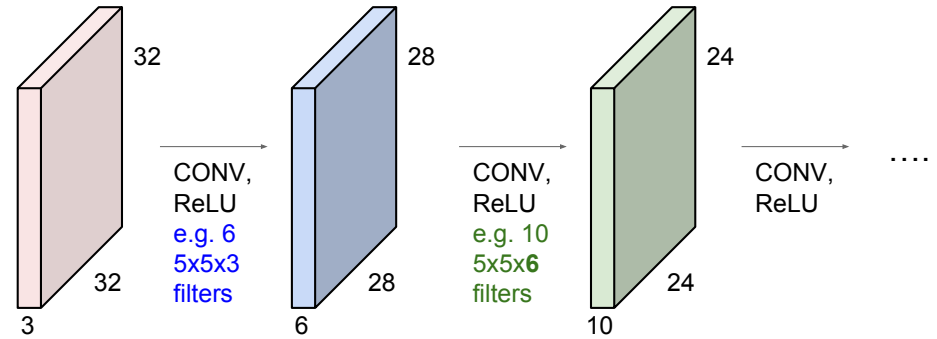


e.g. input 7x7  
**3x3** filter, applied with **stride 1**  
**pad with 1 pixel** border => what is the output?

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



## 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

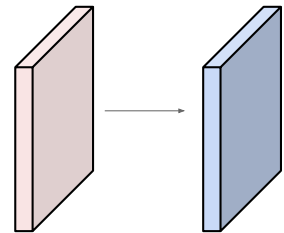
- $W_2 = (W_1 + 2P - F)/S + 1$
- $H_2 = (H_1 + 2P - F)/S + 1$
- $D_2 = K$

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

**Common setting:**  $F = 3, S = P = 1$

Examples time:

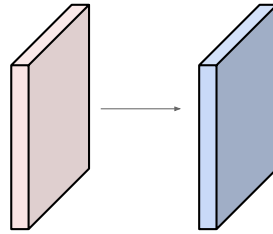
Input volume: **32x32x3**  
 10 5x5 filters with stride 1, pad 2



Output volume size: ?

Examples time:

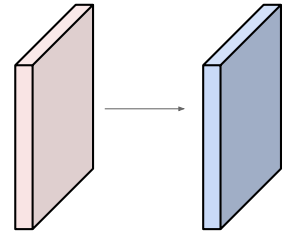
Input volume: **32x32x3**  
10 5x5 filters with stride 1, pad 2



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

Examples time:

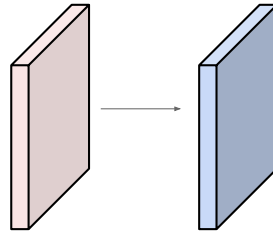
Input volume: **32x32x3**  
10 5x5 filters with stride 1, pad 2



Number of parameters in this layer?

Examples time:

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

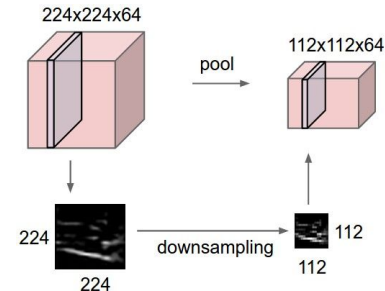


Number of parameters in this layer?  
each filter has  $5*5*3 + 1 = 76$  params (+1 for bias)  
 $\Rightarrow 76*10 = 760$

### Another element: pooling

#### Pooling layer

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

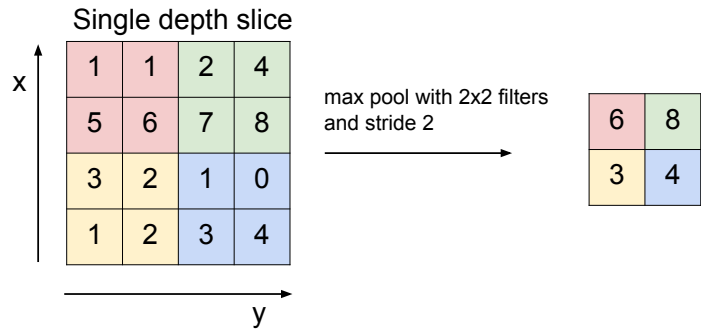


Types of pooling: average, L2-norm, max

## Max pooling

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

### MAX POOLING



Fei-Fei Li & Justin Johnson & Serena Yeung Lecture 5 - 73 April 18, 2017

September 19, 2018 24 / 50

## How to train a CNN?

*How do we learn the filters/weights?*

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

September 19, 2018 26 / 50

## 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.

September 19, 2018 25 / 50

## Outline

- 1 Review of last lecture
- 2 Convolutional neural networks
- 3 Kernel methods
  - Motivation
  - Dual formulation of linear regression
  - Kernel Trick

September 19, 2018 27 / 50



## Motivation

Recall the question: *how to choose nonlinear basis*  $\phi : \mathbb{R}^D \rightarrow \mathbb{R}^M$ ?

$$\mathbf{w}^T \phi(\mathbf{x})$$

- neural network is one approach: learn  $\phi$  from data
- **kernel method** is another one: sidestep the issue of choosing  $\phi$  by using *kernel functions*

## A closer look at the least square solution

By setting the gradient of  $F(\mathbf{w}) = \|\Phi\mathbf{w} - \mathbf{y}\|_2^2 + \lambda\|\mathbf{w}\|_2^2$  to be  $\mathbf{0}$ :

$$\Phi^T(\Phi\mathbf{w}^* - \mathbf{y}) + \lambda\mathbf{w}^* = \mathbf{0}$$

we know

$$\mathbf{w}^* = \frac{1}{\lambda} \Phi^T(\mathbf{y} - \Phi\mathbf{w}^*) = \Phi^T \boldsymbol{\alpha} = \sum_{n=1}^N \alpha_n \phi(\mathbf{x}_n)$$

Thus the least square solution is **a linear combination of features!**

Note this is true for perceptron and many other problems.

*Of course, the above calculation does not show what  $\boldsymbol{\alpha}$  is.*

## Case study: regularized linear regression

Kernel methods work for *many problems* and we take **regularized linear regression** as an example.

Recall the regularized least square solution:

$$\begin{aligned} \mathbf{w}^* &= \underset{\mathbf{w}}{\operatorname{argmin}} F(\mathbf{w}) \\ &= \underset{\mathbf{w}}{\operatorname{argmin}} (\|\Phi\mathbf{w} - \mathbf{y}\|_2^2 + \lambda\|\mathbf{w}\|_2^2) \\ &= (\Phi^T\Phi + \lambda\mathbf{I})^{-1} \Phi^T\mathbf{y} \end{aligned} \quad \left| \quad \Phi = \begin{pmatrix} \phi(\mathbf{x}_1)^T \\ \phi(\mathbf{x}_2)^T \\ \vdots \\ \phi(\mathbf{x}_N)^T \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}$$

Issue: *operate in space  $\mathbb{R}^M$  and  $M$  could be huge or even infinity!*

## Why is this helpful?

Assuming we know  $\boldsymbol{\alpha}$ , the prediction of  $\mathbf{w}^*$  on a new example  $\mathbf{x}$  is

$$\mathbf{w}^{*T} \phi(\mathbf{x}) = \sum_{n=1}^N \alpha_n \phi(\mathbf{x}_n)^T \phi(\mathbf{x})$$

Therefore we do not really need to know  $\mathbf{w}^*$ . *Only inner products in the new feature space matter!*

Kernel methods are exactly about computing inner products *without knowing  $\phi$* .

But we need to figure out what  $\boldsymbol{\alpha}$  is first!

## How to find $\alpha$ ?

Plugging in  $w = \Phi^T \alpha$  into  $F(w)$  gives

$$\begin{aligned}
 G(\alpha) &= F(\Phi^T \alpha) \\
 &= \|\Phi \Phi^T \alpha - y\|_2^2 + \lambda \|\Phi^T \alpha\|_2^2 \\
 &= \|\mathbf{K} \alpha - y\|_2^2 + \lambda \alpha^T \mathbf{K} \alpha && (\mathbf{K} = \Phi \Phi^T) \\
 &= \alpha^T \mathbf{K}^T \mathbf{K} \alpha - 2y^T \mathbf{K} \alpha + \lambda \alpha^T \mathbf{K} \alpha + \text{cnt.} \\
 &= \alpha^T (\mathbf{K}^2 + \lambda \mathbf{K}) \alpha - 2y^T \mathbf{K} \alpha + \text{cnt.} && (\mathbf{K}^T = \mathbf{K})
 \end{aligned}$$

This is sometime called the *dual formulation* of linear regression.

$\mathbf{K} = \Phi \Phi^T \in \mathbb{R}^{N \times N}$  is called **Gram matrix** or **kernel matrix** where the  $(i, j)$  entry is

$$\phi(x_i)^T \phi(x_j)$$

## Examples of kernel matrix

3 data points in  $\mathbb{R}$

$$x_1 = -1, x_2 = 0, x_3 = 1$$

$\phi$  is polynomial basis with degree 4:

$$\phi(x) = \begin{pmatrix} 1 \\ x \\ x^2 \\ x^3 \end{pmatrix}$$

$$\phi(x_1) = \begin{pmatrix} 1 \\ -1 \\ 1 \\ -1 \end{pmatrix} \quad \phi(x_2) = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} \quad \phi(x_3) = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$$

## Calculation of the Gram matrix

$$\phi(x_1) = \begin{pmatrix} 1 \\ -1 \\ 1 \\ -1 \end{pmatrix} \quad \phi(x_2) = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} \quad \phi(x_3) = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$$

### Gram/Kernel matrix

$$\begin{aligned}
 \mathbf{K} &= \begin{pmatrix} \phi(x_1)^T \phi(x_1) & \phi(x_1)^T \phi(x_2) & \phi(x_1)^T \phi(x_3) \\ \phi(x_2)^T \phi(x_1) & \phi(x_2)^T \phi(x_2) & \phi(x_2)^T \phi(x_3) \\ \phi(x_3)^T \phi(x_1) & \phi(x_3)^T \phi(x_2) & \phi(x_3)^T \phi(x_3) \end{pmatrix} \\
 &= \begin{pmatrix} 4 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & 1 & 4 \end{pmatrix}
 \end{aligned}$$

## Gram matrix vs covariance matrix

	dimensions	entry $(i, j)$	property
$\Phi \Phi^T$	$N \times N$	$\phi(x_i)^T \phi(x_j)$	both are <b>symmetric</b> and <b>positive semidefinite</b>
$\Phi^T \Phi$	$M \times M$	$\sum_{n=1}^N \phi(x_n)_i \phi(x_n)_j$	

## How to find $\alpha$ ?

Minimize the dual formulation

$$G(\alpha) = \alpha^T(K^2 + \lambda K)\alpha - 2\mathbf{y}^T K\alpha + \text{cnt.}$$

Setting the derivative to  $\mathbf{0}$  we have

$$\mathbf{0} = (K^2 + \lambda K)\alpha - K\mathbf{y} = K((K + \lambda I)\alpha - \mathbf{y})$$

Thus  $\alpha = (K + \lambda I)^{-1}\mathbf{y}$  is a **minimizer** and we obtain

$$\mathbf{w}^* = \Phi^T \alpha = \Phi^T (K + \lambda I)^{-1} \mathbf{y}$$

Exercise: *are there other minimizers?*

## Then what is the difference?

First, computing  $(\Phi\Phi^T + \lambda I)^{-1}$  can be more efficient than computing  $(\Phi^T\Phi + \lambda I)^{-1}$  when  $N \leq M$ .

More importantly, computing  $\alpha = (\Phi\Phi^T + \lambda I)^{-1}$  also *only requires computing inner products in the new feature space!*

Now we can conclude that the exact form of  $\phi(\cdot)$  is not essential; *all we need is computing inner products  $\phi(\mathbf{x})^T \phi(\mathbf{x}')$ .*

For some  $\phi$  it is indeed possible to compute  $\phi(\mathbf{x})^T \phi(\mathbf{x}')$  without computing/knowing  $\phi$ . This is the *kernel trick*.

## Comparing two solutions

**Minimizing  $F(\mathbf{w})$  gives  $\mathbf{w}^* = (\Phi^T\Phi + \lambda I)^{-1}\Phi^T\mathbf{y}$**

**Minimizing  $G(\alpha)$  gives  $\mathbf{w}^* = \Phi^T(\Phi\Phi^T + \lambda I)^{-1}\mathbf{y}$**

Note  $I$  has different dimensions in two formulas.

Natural question: *are they the same or different?*

**They are the same**

$$\begin{aligned} & (\Phi^T\Phi + \lambda I)^{-1}\Phi^T\mathbf{y} \\ &= (\Phi^T\Phi + \lambda I)^{-1}\Phi^T(\Phi\Phi^T + \lambda I)(\Phi\Phi^T + \lambda I)^{-1}\mathbf{y} \\ &= (\Phi^T\Phi + \lambda I)^{-1}(\Phi^T\Phi\Phi^T + \lambda\Phi^T)(\Phi\Phi^T + \lambda I)^{-1}\mathbf{y} \\ &= (\Phi^T\Phi + \lambda I)^{-1}(\Phi^T\Phi + \lambda I)\Phi^T(\Phi\Phi^T + \lambda I)^{-1}\mathbf{y} \\ &= \Phi^T(\Phi\Phi^T + \lambda I)^{-1}\mathbf{y} \end{aligned}$$

## Example

Consider the following polynomial basis  $\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$ :

$$\phi(\mathbf{x}) = \begin{pmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{pmatrix}$$

What is the inner product between  $\phi(\mathbf{x})$  and  $\phi(\mathbf{x}')$ ?

$$\begin{aligned} \phi(\mathbf{x})^T \phi(\mathbf{x}') &= x_1^2 x_1'^2 + 2x_1x_2x_1'x_2' + x_2^2 x_2'^2 \\ &= (x_1x_1' + x_2x_2')^2 = (\mathbf{x}^T \mathbf{x}')^2 \end{aligned}$$

Therefore, *the inner product in the new space is simply a function of the inner product in the original space.*

## Another example

$\phi : \mathbb{R}^D \rightarrow \mathbb{R}^{2D}$  is parameterized by  $\theta$ :

$$\phi_\theta(\mathbf{x}) = \begin{pmatrix} \cos(\theta x_1) \\ \sin(\theta x_1) \\ \vdots \\ \cos(\theta x_D) \\ \sin(\theta x_D) \end{pmatrix}$$

What is the inner product between  $\phi_\theta(\mathbf{x})$  and  $\phi_\theta(\mathbf{x}')$ ?

$$\begin{aligned} \phi_\theta(\mathbf{x})^\top \phi_\theta(\mathbf{x}') &= \sum_{d=1}^D \cos(\theta x_d) \cos(\theta x'_d) + \sin(\theta x_d) \sin(\theta x'_d) \\ &= \sum_{d=1}^D \cos(\theta(x_d - x'_d)) \end{aligned}$$

Once again, *the inner product in the new space is a simple function of the features in the original space.*

## Infinite dimensional mapping

When  $L \rightarrow \infty$ , even if we cannot compute  $\phi(x)$ , a vector of *infinite dimension*, we can still compute inner product:

$$\begin{aligned} \phi_\infty(\mathbf{x})^\top \phi_\infty(\mathbf{x}') &= \int_0^{2\pi} \sum_{d=1}^D \cos(\theta(x_d - x'_d)) d\theta \\ &= \sum_{d=1}^D \frac{\sin(2\pi(x_d - x'_d))}{x_d - x'_d} \end{aligned}$$

Again, a simple function of the original features.

Note that using this mapping in linear regression, we are *learning a weight  $\mathbf{w}^*$  with infinite dimension!*

## More complicated example

Based on  $\phi_\theta$ , define  $\phi_L : \mathbb{R}^D \rightarrow \mathbb{R}^{2D(L+1)}$  for some integer  $L$ :

$$\phi_L(\mathbf{x}) = \begin{pmatrix} \phi_0(\mathbf{x}) \\ \phi_{\frac{2\pi}{L}}(\mathbf{x}) \\ \phi_{2\frac{2\pi}{L}}(\mathbf{x}) \\ \vdots \\ \phi_{L\frac{2\pi}{L}}(\mathbf{x}) \end{pmatrix}$$

What is the inner product between  $\phi_L(\mathbf{x})$  and  $\phi_L(\mathbf{x}')$ ?

$$\begin{aligned} \phi_L(\mathbf{x})^\top \phi_L(\mathbf{x}') &= \sum_{\ell=0}^L \phi_{\frac{2\pi\ell}{L}}(\mathbf{x})^\top \phi_{\frac{2\pi\ell}{L}}(\mathbf{x}') \\ &= \sum_{\ell=0}^L \sum_{d=1}^D \cos\left(\frac{2\pi\ell}{L}(x_d - x'_d)\right) \end{aligned}$$

## Kernel functions

**Definition:** a function  $k : \mathbb{R}^D \times \mathbb{R}^D \rightarrow \mathbb{R}$  is called a *(positive semidefinite) kernel function* if there exists a function  $\phi : \mathbb{R}^D \rightarrow \mathbb{R}^M$  so that for any  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^D$ ,

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^\top \phi(\mathbf{x}')$$

Examples we have seen

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}')^2$$

$$k(\mathbf{x}, \mathbf{x}') = \sum_{d=1}^D \frac{\sin(2\pi(x_d - x'_d))}{x_d - x'_d}$$

## Using kernel functions

Choosing a nonlinear basis  $\phi$  becomes choosing a kernel function.

As long as computing the kernel function is more efficient, we should apply the kernel trick.

**Gram/kernel matrix** becomes:

$$\mathbf{K} = \Phi\Phi^T = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_2) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix}$$

In fact,  $k$  is a kernel if and only if  $\mathbf{K}$  is positive semidefinite for *any*  $N$  and *any*  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$  (**Mercer theorem**).

## More examples of kernel functions

Two most commonly used kernel functions in practice:

### Polynomial kernel

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + c)^d$$

for  $c \geq 0$  and  $d$  is a positive integer.

### Gaussian kernel or Radial basis function (RBF) kernel

$$k(\mathbf{x}, \mathbf{x}') = e^{-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\sigma^2}}$$

for some  $\sigma > 0$ .

Think about *what the corresponding  $\phi$  is* for each kernel.

## Using kernel functions

The prediction on a new example  $\mathbf{x}$  is

$$\mathbf{w}^{*\top} \phi(\mathbf{x}) = \sum_{n=1}^N \alpha_n \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}) = \sum_{n=1}^N \alpha_n k(\mathbf{x}_n, \mathbf{x})$$

## Composing kernels

Creating more kernel functions using the following rules:

If  $k_1(\cdot, \cdot)$  and  $k_2(\cdot, \cdot)$  are kernels, the followings are kernels too

- **linear combination:**  $\alpha k_1(\cdot, \cdot) + \beta k_2(\cdot, \cdot)$  if  $\alpha, \beta \geq 0$
- **product:**  $k_1(\cdot, \cdot) k_2(\cdot, \cdot)$
- **exponential:**  $e^{k(\cdot, \cdot)}$
- ...

Verify using the definition of kernel!

## Examples that are not kernels

Function

$$k(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} - \mathbf{x}'\|_2^2$$

is *not a kernel*, why?

Example: if it is a kernel, the kernel matrix for two data points  $\mathbf{x}_1$  and  $\mathbf{x}_2$ :

$$\mathbf{K} = \begin{pmatrix} 0 & \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 \\ \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 & 0 \end{pmatrix}$$

must be positive semidefinite, *but is it?*

## Example: Kernelized NNC

For NNC with **L2 distance**, the key is to compute for any two points  $\mathbf{x}, \mathbf{x}'$

$$d(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} - \mathbf{x}'\|_2^2 = \mathbf{x}^T \mathbf{x} + \mathbf{x}'^T \mathbf{x}' - 2\mathbf{x}^T \mathbf{x}'$$

With a kernel function  $k$ , we simply compute

$$d^{\text{KERNEL}}(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}, \mathbf{x}) + k(\mathbf{x}', \mathbf{x}') - 2k(\mathbf{x}, \mathbf{x}')$$

which by definition is the **L2 distance in a new feature space**

$$d^{\text{KERNEL}}(\mathbf{x}, \mathbf{x}') = \|\phi(\mathbf{x}) - \phi(\mathbf{x}')\|_2^2$$

## Kernelizing other ML algorithms

Kernel trick is applicable to **many ML algorithms**:

- nearest neighbor classifier
- perceptron
- logistic regression
- SVM
- ...