

# CSCI-567: Machine Learning

Prof. Victor Adamchik

U of Southern California

July 13, 2020

Your model is only as good as your data.

## Midterm Exam

### Instructions:

- This is a take-home open-book midterm exam.
- The exams is paper-and-pencil, which you scan it and submit electronically.
- Questions should be answered concisely.
- Write legibly, avoid cursive writings.

## Outline

① Convolutional neural networks

② Kernel methods

- 1 Convolutional neural networks
- 2 Kernel methods

The materials I used to develop these notes include the following sources:

- Stanford Course Cs231n: <http://cs231n.stanford.edu/> (taught by Prof. Fei-Fei Li and her students).
- Dr. Ian Goodfellow's lectures on deep learning: <http://deeplearningbook.org>

Both website provides tons of useful resources: notes, demos, videos, etc.

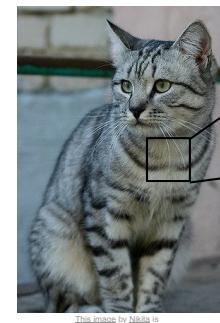
## Image Classification: A core task in Computer Vision



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

→ cat

## The Problem: Semantic Gap



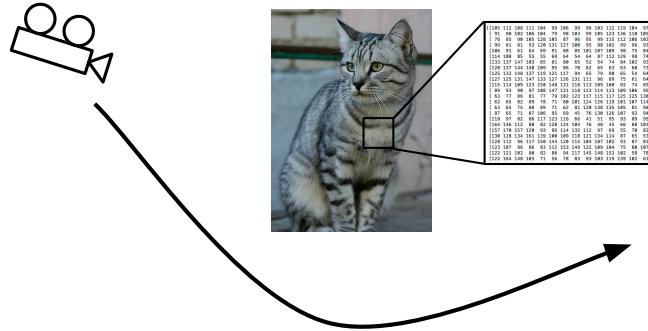
1185 112 180 111 184 99 186 99 186 183 112 110 184 97 93 871
1186 112 180 111 184 99 186 99 186 183 112 110 184 97 93 871
1187 85 98 185 128 185 87 96 99 99 115 112 186 183 99 851
1188 85 98 185 128 185 87 96 99 99 115 112 186 183 99 851
1189 81 81 93 128 131 127 188 95 98 182 99 99 93 181 841
1190 81 81 93 128 131 127 188 95 98 182 99 99 93 181 841
1191 148 85 95 55 55 69 54 64 87 112 120 96 74 84 911
1192 148 85 95 55 55 69 54 64 87 112 120 96 74 84 911
1193 137 144 148 189 95 89 78 62 65 63 63 68 73 86 1811
1194 137 144 148 189 95 89 78 62 65 63 63 68 73 86 1811
1195 133 148 137 119 121 117 94 65 79 88 65 54 64 72 981
1196 133 148 137 119 121 117 94 65 79 88 65 54 64 72 981
1197 134 149 137 119 121 117 94 65 79 88 65 54 64 72 981
1198 134 149 137 119 121 117 94 65 79 88 65 54 64 72 981
1199 134 149 137 119 121 117 94 65 79 88 65 54 64 72 981
1200 77 86 81 77 79 182 123 117 115 117 125 125 138 115 871
1201 77 86 81 77 79 182 123 117 115 117 125 125 138 115 871
1202 65 75 88 89 71 62 91 128 138 135 180 91 98 118 1181
1203 65 75 88 89 71 62 91 128 138 135 180 91 98 118 1181
1204 65 75 88 89 71 62 91 128 138 135 180 91 98 118 1181
1205 65 75 88 89 71 62 91 128 138 135 180 91 98 118 1181
1206 146 112 88 82 128 121 186 76 45 66 68 88 181 182 1891
1207 146 112 88 82 128 121 186 76 45 66 68 88 181 182 1891
1208 134 161 139 189 180 180 118 121 134 114 87 65 53 69 861
1209 112 96 117 158 144 128 115 184 187 182 93 81 72 791
1210 112 96 117 158 144 128 115 184 187 182 93 81 72 791
1211 125 125 182 88 82 86 94 117 145 148 153 182 58 78 92 1871
1212 125 125 182 88 82 86 94 117 145 148 153 182 58 78 92 1871
1213 125 125 182 88 82 86 94 117 145 148 153 182 58 78 92 1871

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



## Challenges: Illumination



This image by Nikita is licensed under CC-BY 2.0

## Challenges: Deformation



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

## Challenges: Occlusion



This image by Jonson is licensed under CC-BY 2.0

## Challenges: Background Clutter



This image is CC0 1.0 public domain



This image is CC0 1.0 public domain

## Challenges: Intraclass variation



This image is CC0 1.0 public domain

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 2 - 13

April 5, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 2 - 14

April 5, 2018

## Fundamental problems in vision

An image classifier is not like sorting a list of numbers: there no obvious way to develop an algorithm for recognizing a cat, or other classes.

### The key challenge

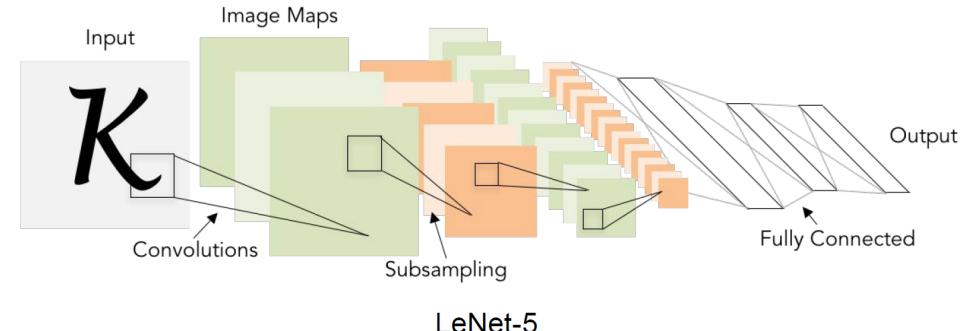
How to train a model that can tolerate all those variations?

### Main ideas

- need a lot of data that exhibits those variations
- need more specialized models to capture the invariance

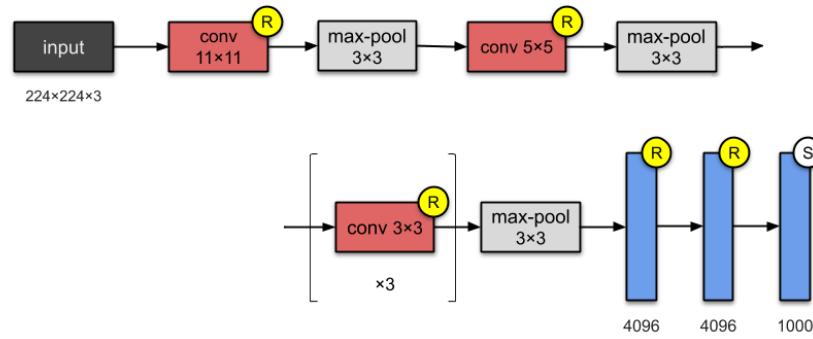
## A bit of history:

Gradient-based learning applied to document recognition [LeCun, Bottou, Bengio, Haffner 1998]



## A bit of history:

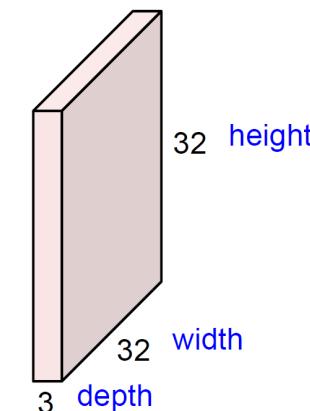
ImageNet Classification with Deep Convolutional Neural Networks  
[Krizhevsky, Sutskever, Hinton, 2012]



## Convolution layer

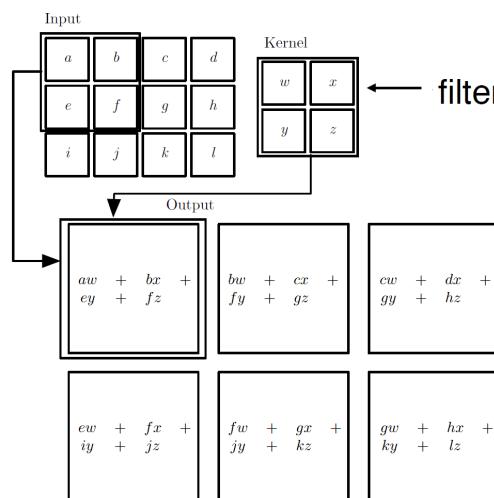
Arrange neurons as a **3D matrix/tensor** naturally.

$32 \times 32 \times 3$  image – this preserves spatial structure



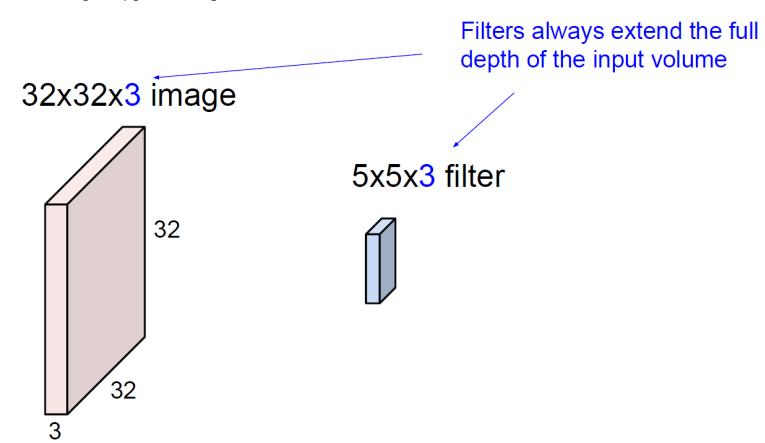
## 2D Convolution

This is a computation of a dot product.



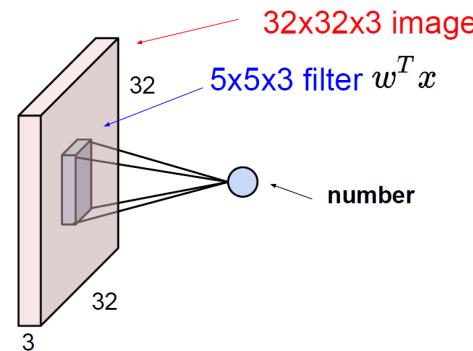
## Convolution layer

Convolve the filter with the image i.e. “slide over the image spatially, computing dot products”



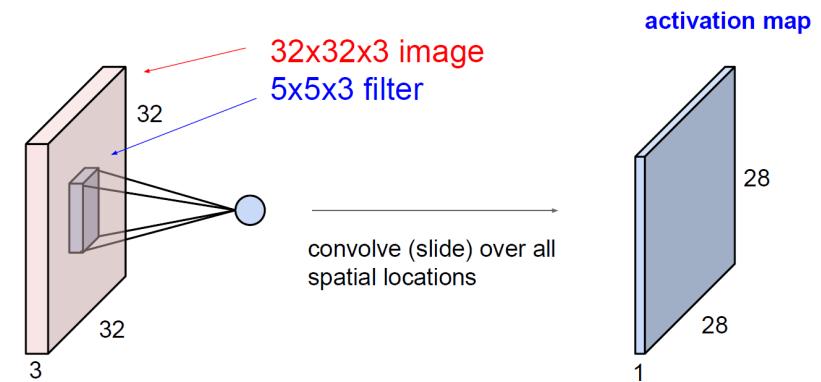
## Convolution layer

The result of taking a dot product between the filter  $w^T x$  and a small  $5 \times 5 \times 3$  chunk of the image.



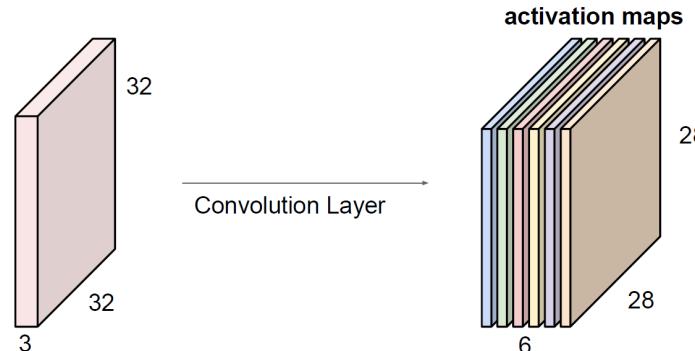
## Convolution layer

The result of taking a dot product between the filter  $w^T x$  and the image.



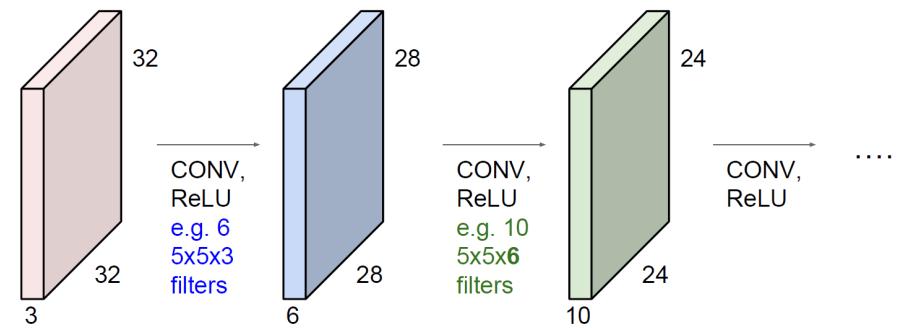
## Convolution layer

If we had six  $5 \times 5 \times 3$  filters, we'll get 6 separate activation maps. We stack them up to get a new image of size  $28 \times 28 \times 6$ .



## Convolution layer

Convolution Network is a sequence of Convolutional Layers, interspersed with activation functions



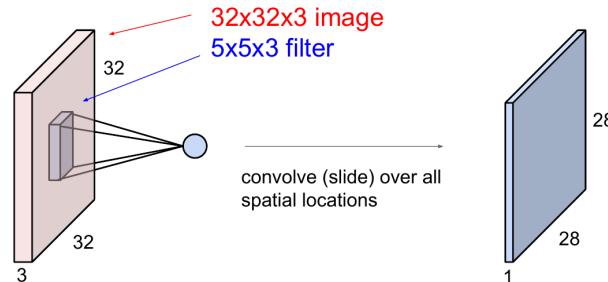
## Connection to fully connected NNs

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

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

*Much less parameters!* Example:

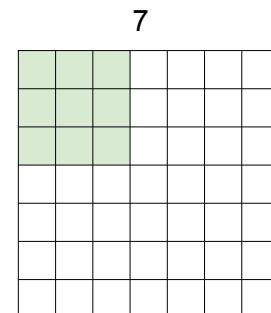
- FC:  $(32 \times 32 \times 3) \times (28 \times 28) \approx 2.4M$
- CNN:  $5 \times 5 \times 3 = 75$



July 13, 2020 17 / 56

## Spatial arrangement: stride and padding

A closer look at spatial dimensions:



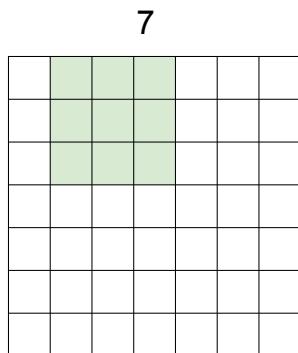
7x7 input (spatially)  
assume 3x3 filter

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 42 April 18, 2017

July 13, 2020 18 / 56

A closer look at spatial dimensions:

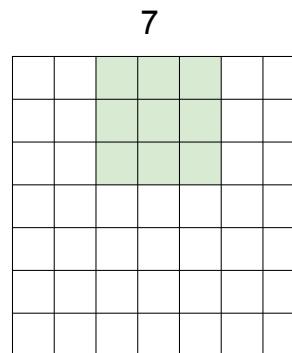


7x7 input (spatially)  
assume 3x3 filter

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 43 April 18, 2017

A closer look at spatial dimensions:

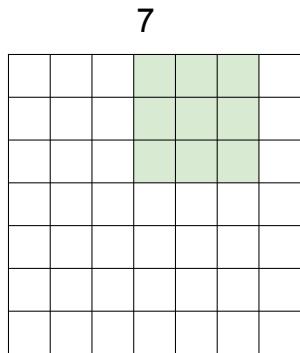


7x7 input (spatially)  
assume 3x3 filter

Fei-Fei Li & Justin Johnson & Serena Yeung

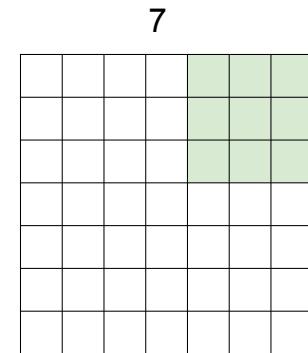
Lecture 5 - 44 April 18, 2017

A closer look at spatial dimensions:



7x7 input (spatially)  
assume 3x3 filter

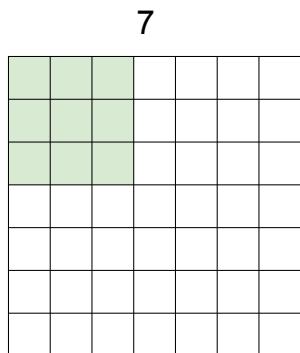
A closer look at spatial dimensions:



7x7 input (spatially)  
assume 3x3 filter

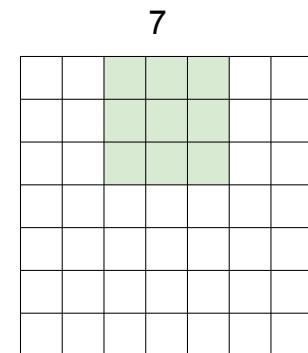
=> **5x5 output**

A closer look at spatial dimensions:



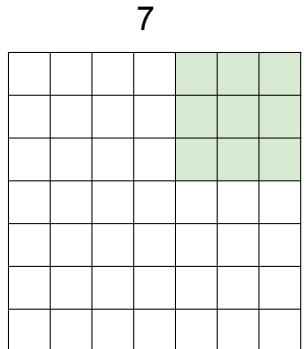
7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 2**

A closer look at spatial dimensions:



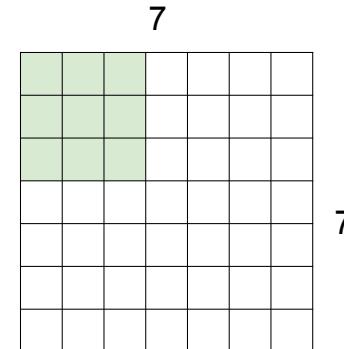
7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 2**

A closer look at spatial dimensions:



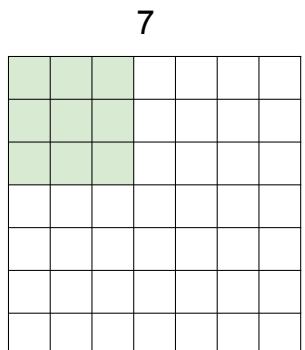
7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 2**  
**=> 3x3 output!**

A closer look at spatial dimensions:



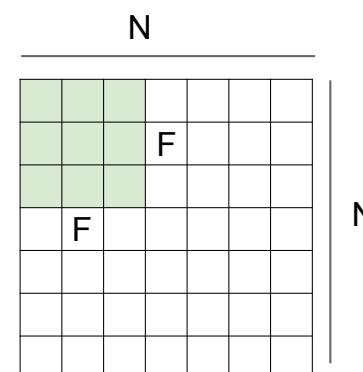
7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 3?**

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) / stride + 1**

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

## 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?

(recall:)

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

## 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

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 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)

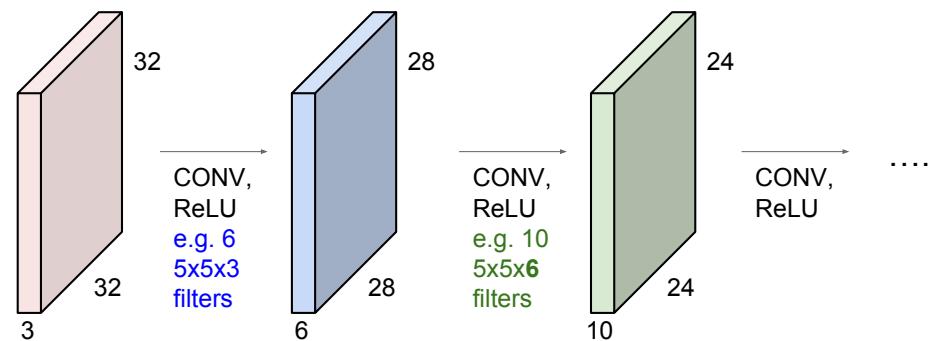
e.g.  $F = 3 \Rightarrow$  zero pad with 1

$F = 5 \Rightarrow$  zero pad with 2

$F = 7 \Rightarrow$  zero pad with 3

## Remember back to...

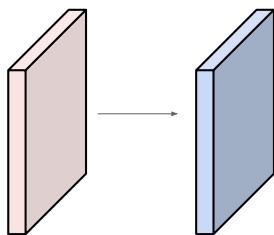
E.g. 32x32 input convolved repeatedly with 5x5 filters shrinks volumes spatially! (32  $\rightarrow$  28  $\rightarrow$  24 ...). Shrinking too fast is not good, doesn't work well.



Examples time:

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2

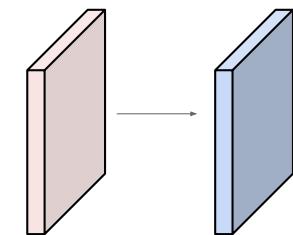


Output volume size: ?

Examples time:

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2



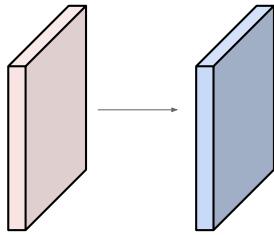
Output volume size:

$(32+2^2-5)/1+1 = 32$  spatially, so  
**32x32x10**

Examples time:

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2

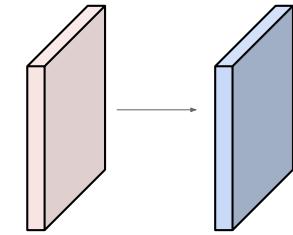


Number of parameters in this layer?

Examples time:

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2



Number of parameters in this layer?

each filter has  $5 \times 5 \times 3 + 1 = 76$  params  
=> **76 \* 10 = 760**

(+1 for bias)

## Summary for convolution layer

**Input:** a tensor 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 tensor of size  $W_2 \times H_2 \times D_2$  where

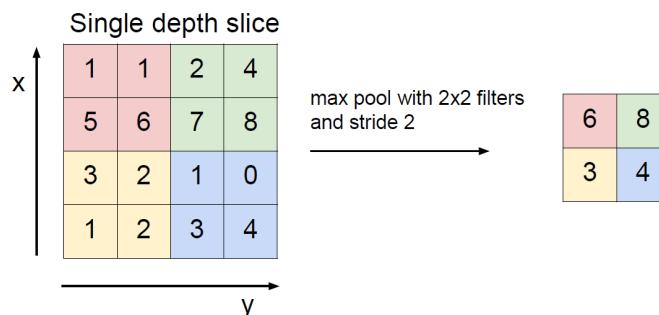
- $W_2 = (W_1 + 2P - F)/S + 1$
- $H_2 = (H_1 + 2P - F)/S + 1$
- $D_2 = K$

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

**Common setting:**  $F = 3, S = P = 1$  and  $K$  is a power of two.

## Max pooling

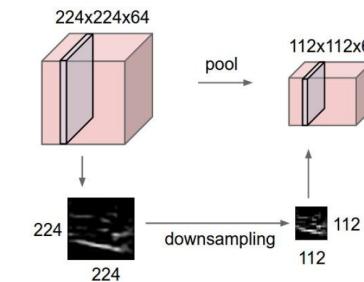
**Max pooling** with  $2 \times 2$  filter and stride 2 is very common.



## 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

## Putting everything together

**Typical architecture for CNNs:**

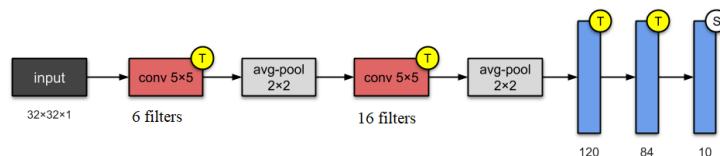
Input  $\rightarrow$  [[Conv  $\rightarrow$  ReLU] $^N$   $\rightarrow$  Pool] $^M$   $\rightarrow$  [FC  $\rightarrow$  ReLU] $^Q$   $\rightarrow$  Softmax

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

**Well-known CNNs:**

LeNet (1998), AlexNet(2012), ZF Net(2013), GoogLeNet(2014), VGGNet(2014), ResNet(2015), SqueezeNet (2016), SENet(2017)

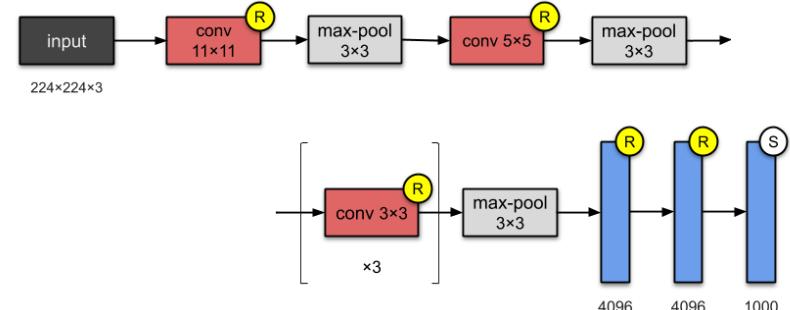
## LeNet - 1996



$T$  stands for the tanh function.

$S$  stands for the softmax function

## AlexNet - 2012

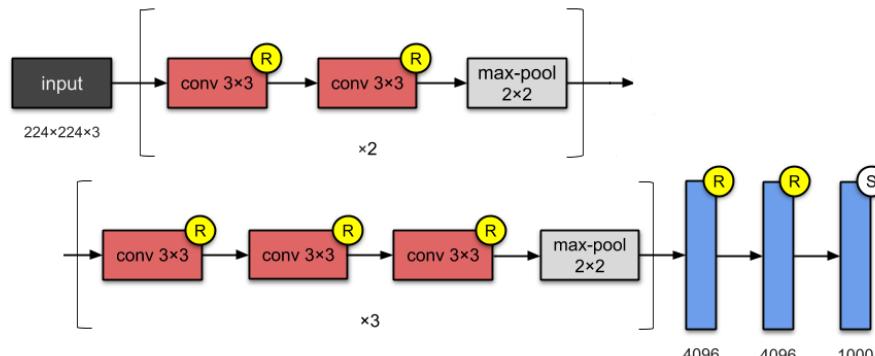


$R$  stands for the ReLU function.

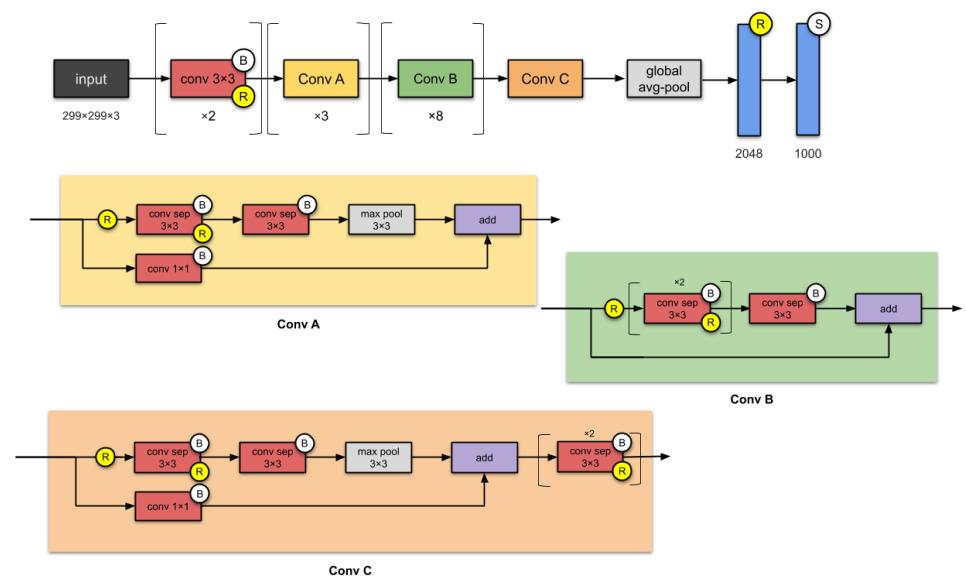
July 13, 2020 23 / 56

July 13, 2020 24 / 56

## Visual Geometry Group - 2014



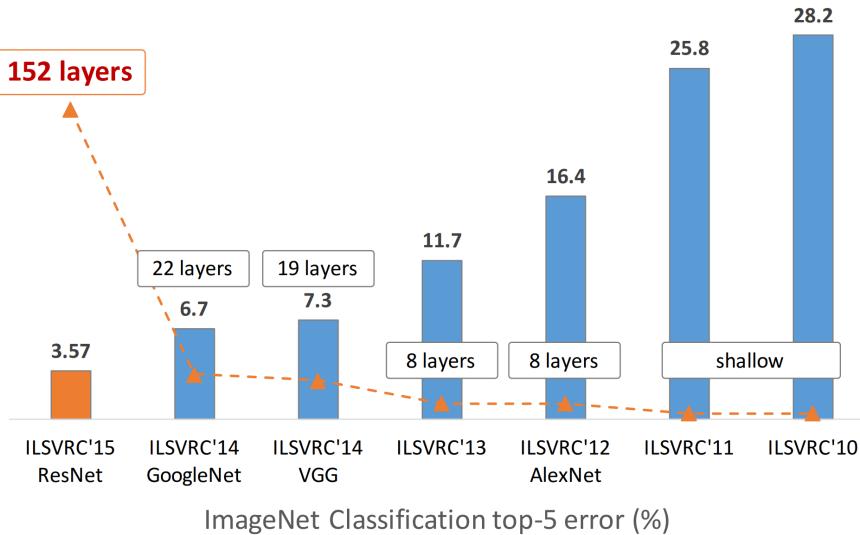
## ResNet - 2015



July 13, 2020 25 / 56

July 13, 2020 26 / 56

# Revolution of Depth



Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition".

July 13, 2020 27 / 56

## Outline

1 Convolutional neural networks

2 Kernel methods

- Dual formulation of linear regression
- Kernel Trick
- Kernelizing ML algorithms

## How to train a CNN?

*How do we learn the filters/weights?*

Run the image through 30+ layers of convolutions and train the filters with **SGD/backpropagation**

*What computer architecture to use?*

buy 10 or more GPUs (NVIDIA graphic card)

*What software package to use?*

State-of-the-art: TensorFlow plus Keras or Caffe on the top.

See demo at:

<http://cs.stanford.edu/people/karpathy/convnetjs/demo/cifar10.html>

July 13, 2020 28 / 56

## 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*

July 13, 2020 29 / 56

July 13, 2020 30 / 56

Recall the regularized least square solution:

$$\begin{aligned} \mathbf{w}^* &= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \\ \mathbf{w}^* &= (\Phi^T \Phi + \lambda \mathbf{I})^{-1} \Phi^T \mathbf{y} \end{aligned} \quad \left| \begin{array}{l} \mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_N^T \end{pmatrix}, \Phi = \begin{pmatrix} \phi(\mathbf{x}_1)^T \\ \phi(\mathbf{x}_2)^T \\ \vdots \\ \phi(\mathbf{x}_N)^T \end{pmatrix} \end{array} \right.$$

Here  $\mathbf{X}^T \mathbf{X}$  is  $D \times D$  matrix, and  $\Phi^T \Phi$  is  $M \times M$  matrix,

Issue:  *$M$  could be huge or even infinity!*

We will rewrite the solution in a different form.

## Gram matrix

We call  $\mathbf{K} = \Phi \Phi^T$  **Gram matrix** or **kernel matrix** where the  $(i, j)$  entry is

$$\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

Therefore,

$$\begin{aligned} \mathbf{K} &= \Phi \Phi^T \\ &= \begin{pmatrix} \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_1) & \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_2) & \cdots & \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_N) \\ \phi(\mathbf{x}_2)^T \phi(\mathbf{x}_1) & \phi(\mathbf{x}_2)^T \phi(\mathbf{x}_2) & \cdots & \phi(\mathbf{x}_2)^T \phi(\mathbf{x}_N) \\ \cdots & \cdots & \cdots & \cdots \\ \phi(\mathbf{x}_N)^T \phi(\mathbf{x}_1) & \phi(\mathbf{x}_N)^T \phi(\mathbf{x}_2) & \cdots & \phi(\mathbf{x}_N)^T \phi(\mathbf{x}_N) \end{pmatrix} \in \mathbb{R}^{N \times N} \end{aligned}$$

**Another minimizer is**

$$\begin{aligned} \mathbf{w}^* &= \Phi^T (\Phi \Phi^T + \lambda \mathbf{I})^{-1} \mathbf{y} \\ \mathbf{w}^* &= \Phi^T (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y} \\ \mathbf{w}^* &= \Phi^T \boldsymbol{\alpha} = \sum_{n=1}^N \alpha_n \phi(\mathbf{x}_n) \end{aligned}$$

where  $\mathbf{K} = \Phi \Phi^T \in \mathbb{R}^{N \times N}$  is the Gram/Kernel matrix and  $\boldsymbol{\alpha}$  is a new vector.

Solution  $\mathbf{w}^*$  is a linear combination of features!

## Another solution

Here we prove that two solutions

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

are the same.

$$\begin{aligned} &(\Phi^T \Phi + \lambda \mathbf{I}_1)^{-1} \Phi^T \mathbf{y} \\ &= (\Phi^T \Phi + \lambda \mathbf{I}_1)^{-1} \Phi^T (\Phi \Phi^T + \lambda \mathbf{I}_2) (\Phi \Phi^T + \lambda \mathbf{I}_2)^{-1} \mathbf{y} \\ &= (\Phi^T \Phi + \lambda \mathbf{I}_1)^{-1} (\Phi^T \Phi \Phi^T + \lambda \Phi^T) (\Phi \Phi^T + \lambda \mathbf{I}_2)^{-1} \mathbf{y} \\ &= (\Phi^T \Phi + \lambda \mathbf{I}_1)^{-1} (\Phi^T \Phi + \lambda \mathbf{I}_1) \Phi^T (\Phi \Phi^T + \lambda \mathbf{I}_2)^{-1} \mathbf{y} \\ &= \Phi^T (\Phi \Phi^T + \lambda \mathbf{I}_2)^{-1} \mathbf{y} \end{aligned}$$

## Then what is the difference?

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

More importantly, computing  $(K + \lambda I)^{-1}$  **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  $\phi(x)^T \phi(x')$ .**

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

## Why is this helpful?

The prediction of  $w^*$  on a new example  $x$  is

$$w^{*T} \phi(x) = \left( \sum_{n=1}^N \alpha_n \phi(x_n)^T \right) \phi(x) = \sum_{n=1}^N \alpha_n (\phi(x_n)^T \phi(x))$$

Therefore we do not really need to know a nonlinear mapping  $\phi$ , only inner products in the new feature space matter!

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

## 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:

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

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

July 13, 2020 35 / 56

July 13, 2020 36 / 56

## 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

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

July 13, 2020 37 / 56

July 13, 2020 38 / 56

## Example of the kernel trick

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  $\phi(\mathbf{x}')$ ?

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

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

Count the number of multiplications.

July 13, 2020 39 / 56

## More complicated example

Based on the previous example mapping  $\phi_\theta$ , we define a new one  $\phi_L : \mathbb{R}^D \rightarrow \mathbb{R}^{2D(L+1)}$  as follows:

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

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

$$\begin{aligned}\phi_L(\mathbf{x})^T \phi_L(\mathbf{x}') &= \sum_{\ell=0}^L \phi_{\frac{2\pi\ell}{L}}(\mathbf{x})^T \phi_{\frac{2\pi\ell}{L}}(\mathbf{x}') \\ &= \sum_{\ell=0}^L \sum_{d=1}^D \cos\left(\frac{2\pi\ell}{L}(x_d - x'_d)\right)\end{aligned}$$

## 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  $\phi_\theta(\mathbf{x}')$ ?

$$\begin{aligned}\phi_\theta(\mathbf{x})^T \phi_\theta(\mathbf{x}') &= \sum_{d=1}^D \cos(\theta x_d) \cos(\theta x'_d) + \sin(\theta x_d) \sin(\theta x'_d) \\ &= \sum_{d=1}^D \cos(\theta(x_d - x'_d)) = k(\mathbf{x}, \mathbf{x}')\end{aligned}$$

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

July 13, 2020 40 / 56

## Infinite dimensional mapping

Let us set  $L \rightarrow \infty$ . This means that  $\phi_L(\mathbf{x})$  vector has infinite dimension. Clearly we cannot compute  $\phi_L(\mathbf{x})$ , but we can still compute the inner

product:

$$\begin{aligned}\phi_\infty(\mathbf{x})^T \phi_\infty(\mathbf{x}') &= \int_0^{2\pi} \sum_{d=1}^D \cos(\theta(x_d - x'_d)) d\theta \\ &= \sum_{d=1}^D \frac{\sin(2\pi(x_d - x'_d))}{x_d - x'_d}\end{aligned}$$

Again, a simple function of the original features.

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

July 13, 2020 41 / 56

July 13, 2020 42 / 56

## 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  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^D$ ,

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$$

Kernel functions are used to quantify similarity between a pair of points  $\mathbf{x}$  and  $\mathbf{x}'$ .

Examples we have seen

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

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}')^2$$

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

July 13, 2020 43 / 56

## Using kernel functions

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

$$\begin{aligned} \mathbf{w}^{*T} \phi(\mathbf{x}) &= \left( \sum_{n=1}^N \alpha_n \phi(\mathbf{x}_n)^T \right) \phi(\mathbf{x}) = \sum_{n=1}^N \alpha_n (\phi(\mathbf{x}_n)^T \phi(\mathbf{x})) \\ &= \sum_{n=1}^N \alpha_n k(\mathbf{x}_n, \mathbf{x}) \end{aligned}$$

July 13, 2020 45 / 56

## 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:

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

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

July 13, 2020 44 / 56

## 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.

July 13, 2020 46 / 56

## 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!

## Kernelizing NNC

Regular KNN algorithm has two shortcomings.

- all neighbors receive equal weight
- the number of neighbors must be chosen globally.

Kernel addresses these issues.

Instead of selected nearest neighbors, all neighbors are used, but with different weights.

Closer neighbors receive higher weight.

The weighting function is a kernel.

One of the most common way is the Gaussian kernel

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

## Kernelizing ML algorithms

There are two **main aspects** of kernelized algorithms:

- the solution is expressed as a linear combination of training examples
- algorithm relies only on inner products between data points

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

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

## Kernelizing NNC

### Kernel Binary Classification Algorithm.

Given

- training data  $\mathcal{D} = \{(\mathbf{x}_n, y_n), n = 1, 2, \dots, N\}$ , where  $y_n \in \{-1, 1\}$
- kernel function  $k(\mathbf{x}_i, \mathbf{x}_j)$
- input  $\mathbf{x}$  to classify

Return the class given by

$$\text{sign} \left( \sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n) y_n \right)$$

## Kernelizing NNC

For NNC with **L2 distance**,  $\|\mathbf{x} - \mathbf{x}'\|_2^2$  is not a valid kernel.

But we can convert the norm

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

into the following kernel function

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

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 Perceptron

### Kernelized Perceptron Algorithm:

- Pick  $(\mathbf{x}_n, y_n)$  randomly
  - Compute  $y^* = \text{sign}(\mathbf{w}^\top \phi(\mathbf{x}_n))$
  - If  $y^* \neq y_n$  then
    - ▶  $\mathbf{w} = \mathbf{w} + y_n \phi(\mathbf{x}_n)$
    - ▶  $\alpha_n = \alpha_n + y_n$
  - Solution  $\mathbf{w} = \sum_n \alpha_n \phi(\mathbf{x}_n)$  is a linear combination of features.
- 
- Pick  $(\mathbf{x}_n, y_n)$  randomly
  - Compute  $y^* = \text{sign}(\sum_i \alpha_i k(\mathbf{x}_i, \mathbf{x}_n))$
  - If  $y^* \neq y_n$  then
    - ▶  $\alpha_n = \alpha_n + y_n$
  - Solution  $\mathbf{w} = \sum_n \alpha_n k(\mathbf{x}_n, \mathbf{x})$ .

## Kernelizing Perceptron

### Perceptron Algorithm:

- Pick  $(\mathbf{x}_n, y_n)$  randomly
- Compute  $y^* = \text{sign}(\mathbf{w}^\top \phi(\mathbf{x}_n))$
- If  $y^* \neq y_n$  then
  - ▶  $\mathbf{w} = \mathbf{w} + y_n \phi(\mathbf{x}_n)$
  - ▶  $\alpha_n = \alpha_n + y_n$
- Solution  $\mathbf{w}$  is a linear combination of features.

## Kernelizing Perceptron

### XOR example

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

Four 2-dimensional training points:

Gram matrix  $\mathbf{K}$ :

$$\begin{pmatrix} x_1 & x_2 & y \\ 1 & 1 & 1 \\ -1 & 1 & -1 \\ -1 & -1 & 1 \\ 1 & -1 & -1 \end{pmatrix}$$

$$\begin{pmatrix} 4 & 0 & 4 & 0 \\ 0 & 4 & 0 & 4 \\ 4 & 0 & 4 & 0 \\ 0 & 4 & 0 & 4 \end{pmatrix}$$

$$k(\mathbf{x}_1, \mathbf{x}_1) = ((1, 1)^\top (1, 1))^2 = 4 \quad k(\mathbf{x}_1, \mathbf{x}_2) = ((1, 1)^\top (-1, 1))^2 = 0$$

$$k(\mathbf{x}_1, \mathbf{x}_3) = ((1, 1)^\top (-1, -1))^2 = 4 \quad k(\mathbf{x}_1, \mathbf{x}_4) = ((1, 1)^\top (1, -1))^2 = 0$$

**XOR example**

Initialization:  $\alpha = (0, 0, 0, 0)$ .

First round

$$\mathbf{x}_1: \text{compute } y^* = \text{sign}\left(\sum_{i=1}^4 \alpha_i k(\mathbf{x}_i, \mathbf{x}_1)\right) = \text{sign}(0) = -1 \neq y_1$$

Thus,  $\alpha_1 = y_1 = 1$ .

$$\mathbf{x}_2: \text{compute } y^* = \text{sign}((1, 0, 0, 0)^T(0, 4, 0, 4)) = \text{sign}(0) = -1 = y_2$$

$$\mathbf{x}_3: \text{compute } y^* = \text{sign}((1, 0, 0, 0)^T(4, 0, 4, 0)) = \text{sign}(4) = 1 = y_3$$

$$\mathbf{x}_4: \text{compute } y^* = \text{sign}((1, 0, 0, 0)^T(0, 4, 0, 4)) = \text{sign}(0) = -1 = y_4$$

**XOR example**

$$\alpha = (1, 0, 0, 0).$$

Second round.

$$\mathbf{x}_1: \text{compute } y^* = \text{sign}((1, 0, 0, 0)^T(4, 0, 4, 0)) = \text{sign}(4) = 1 = y_1$$

$$\mathbf{x}_2: \text{compute } y^* = \text{sign}((1, 0, 0, 0)^T(0, 4, 0, 4)) = \text{sign}(0) = -1 = y_2$$

$$\mathbf{x}_3: \text{compute } y^* = \text{sign}((1, 0, 0, 0)^T(4, 0, 4, 0)) = \text{sign}(4) = 1 = y_3$$

$$\mathbf{x}_4: \text{compute } y^* = \text{sign}((1, 0, 0, 0)^T(0, 4, 0, 4)) = \text{sign}(0) = -1 = y_4$$

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

$$\sum_{i=1}^4 \alpha_i k(\mathbf{x}_i, \mathbf{x}) = k(\mathbf{x}_1, \mathbf{x}) = ((1, 1)^T \mathbf{x})^2 = (x_1 + x_2)^2$$