

CSCI567 Machine Learning (Spring 2025)

Haipeng Luo

University of Southern California

Feb 14, 2025

1 / 42

Outline

- 1 Convolutional neural networks (ConvNets/CNNs)
- 2 Kernel methods

3 / 42

Administration

HW2 will be released today. Due on Feb 27th.

2 / 42

Convolutional neural networks (ConvNets/CNNs)

Outline

- 1 Convolutional neural networks (ConvNets/CNNs)
 - Motivation
 - Architecture
- 2 Kernel methods

4 / 42

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.

Also, demo from <https://poloclub.github.io/cnn-explainer/>

Image Classification: A core task in Computer Vision

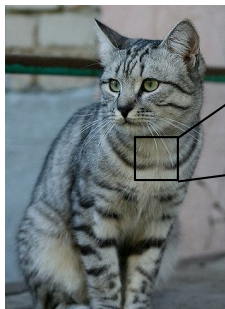


This image by Nikita is licensed under CC-BY 2.0

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

→ cat

The Problem: Semantic Gap



This image by Nikita is licensed under CC-BY 2.0

```

[[185 112 200 111 184 99 106 99 96 183 112 110 184 97 93 87]
 [ 91 98 182 186 184 79 98 183 99 185 123 116 118 185 94 85]
 [ 78 85 98 185 128 180 87 98 95 99 113 112 180 183 99 83]
 [ 99 81 81 93 128 131 127 188 95 98 182 99 96 93 181 94]
 [ 188 91 81 84 89 81 88 85 181 187 189 98 75 84 96 93]
 [114 188 85 55 55 89 64 54 64 87 112 129 98 74 84 91]
 [123 127 147 183 85 81 88 62 54 74 84 182 93 85 82]
 [128 127 144 180 89 85 88 62 54 74 84 182 93 85 82]
 [125 123 148 137 110 111 117 78 65 79 88 65 54 64 72 88]
 [127 125 131 147 113 117 111 111 111 88 75 61 64 84]
 [115 114 188 123 158 148 131 118 113 189 188 92 74 65 72 88]
 [ 89 93 98 97 188 187 111 118 113 114 113 189 186 95 77 88]
 [ 83 77 86 81 77 79 182 123 117 115 117 125 125 118 115 87]
 [ 82 66 82 89 78 71 88 181 124 126 118 181 187 114 113 113]
 [ 83 65 75 88 88 91 62 81 128 138 135 185 81 88 118 118]
 [ 87 65 71 87 186 95 69 65 78 138 138 187 92 84 185 122]
 [118 97 82 86 117 123 116 86 41 51 95 93 89 95 182 187]
 [184 146 132 88 82 128 124 184 78 68 65 88 181 182 189]
 [157 178 157 128 93 86 114 112 112 97 69 55 78 82 99 94]
 [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]
 [122 187 96 86 83 123 149 122 189 184 75 88 187 112 95]
 [122 121 182 88 82 86 94 117 145 148 153 182 58 78 92 187]
 [122 144 148 143 71 56 78 83 93 183 118 128 182 61 69 84]

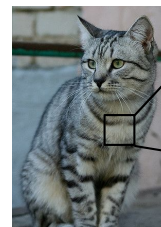
```

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



```

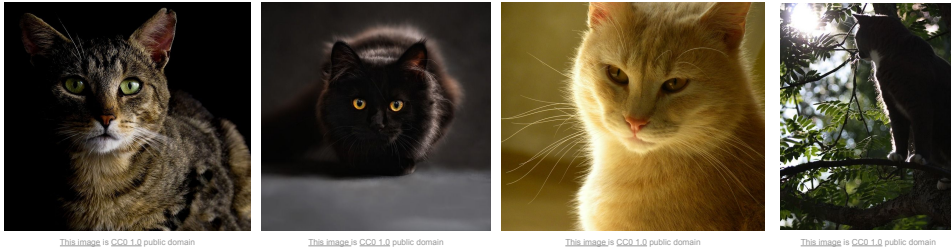
[[185 112 200 111 184 99 106 99 96 183 112 110 184 97 93 87]
 [ 91 98 182 186 184 79 98 183 99 185 123 116 118 185 94 85]
 [ 78 85 98 185 128 180 87 98 95 99 113 112 180 183 99 83]
 [ 99 81 81 93 128 131 127 188 95 98 182 99 96 93 181 94]
 [ 188 91 81 84 89 81 88 85 181 187 189 98 75 84 96 93]
 [114 188 85 55 55 89 64 54 64 87 112 129 98 74 84 91]
 [123 127 147 183 85 81 88 62 54 74 84 182 93 85 82]
 [128 127 144 180 89 85 88 62 54 74 84 182 93 85 82]
 [125 123 148 137 110 111 117 78 65 79 88 65 54 64 72 88]
 [127 125 131 147 113 117 111 111 111 88 75 61 64 84]
 [115 114 188 123 158 148 131 118 113 189 188 92 74 65 72 88]
 [ 89 93 98 97 188 187 111 118 113 114 113 189 186 95 77 88]
 [ 83 77 86 81 77 79 182 123 117 115 117 125 125 118 115 87]
 [ 82 66 82 89 78 71 88 181 124 126 118 181 187 114 113 113]
 [ 83 65 75 88 88 91 62 81 128 138 135 185 81 88 118 118]
 [ 87 65 71 87 186 95 69 65 78 138 138 187 92 84 185 122]
 [118 97 82 86 117 123 116 86 41 51 95 93 89 95 182 187]
 [184 146 132 88 82 128 124 184 78 68 65 88 181 182 189]
 [157 178 157 128 93 86 114 112 112 97 69 55 78 82 99 94]
 [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]
 [122 187 96 86 83 123 149 122 189 184 75 88 187 112 95]
 [122 121 182 88 82 86 94 117 145 148 153 182 58 78 92 187]
 [122 144 148 143 71 56 78 83 93 183 118 128 182 61 69 84]

```

All pixels change when the camera moves!

This image by Nikita is licensed under CC-BY 2.0

Challenges: Illumination



Challenges: Deformation



Challenges: Occlusion



Challenges: Background Clutter



Fundamental problems in vision

Challenges: Intraclass variation



This image is CC0.1.0 public domain

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 2 - 13

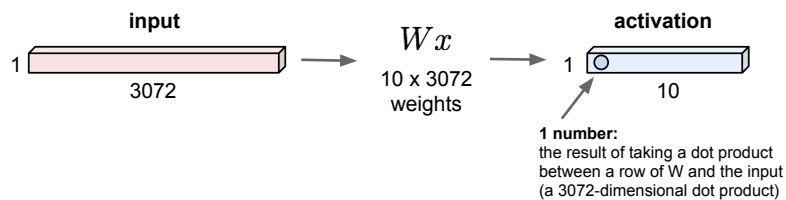
April 6, 2017

6 / 42

Issues of standard NN for image inputs

Fully Connected Layer

32x32x3 image \rightarrow stretch to 3072 x 1



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 27

April 18, 2017

Spatial structure is lost!

7 / 42

Fundamental problems in vision

The key challenge

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

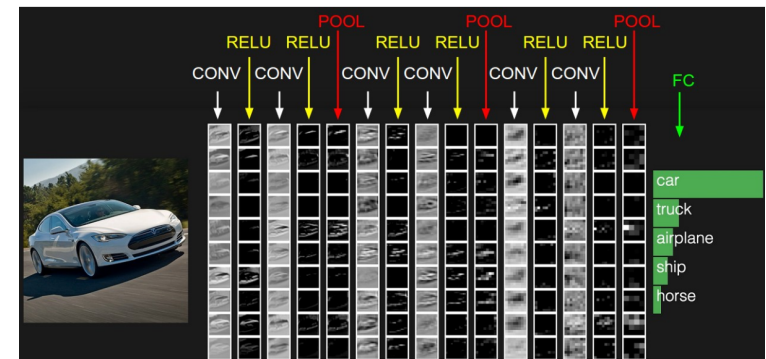
Main ideas

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

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*



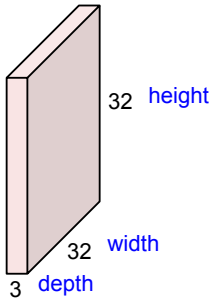
8 / 42

Convolution layer

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

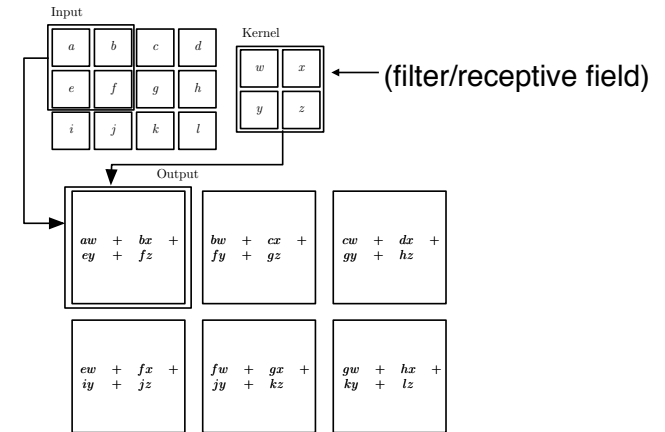
Convolution Layer

32x32x3 image -> preserve spatial structure



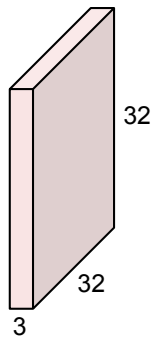
Convolution

2D Convolution



Convolution Layer

32x32x3 image



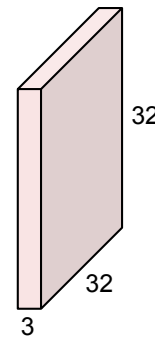
5x5x3 filter



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

Convolution Layer

32x32x3 image



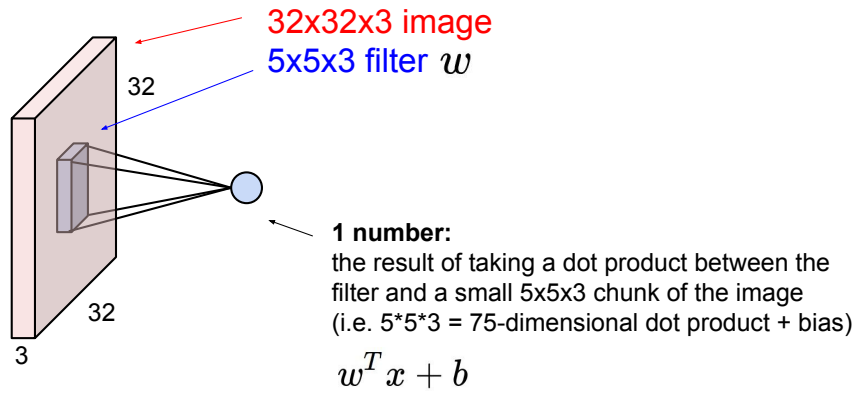
Filters always extend the full
depth of the input volume

5x5x3 filter



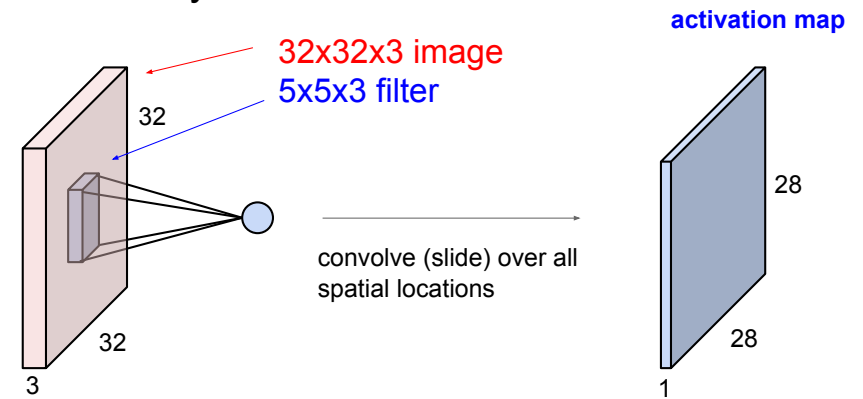
Convolve the filter with the image
i.e. "slide over the image spatially,
computing dot products"

Convolution Layer



Fei-Fei Li & Justin Johnson & Serena Yeung Lecture 5 - 31 April 18, 2017

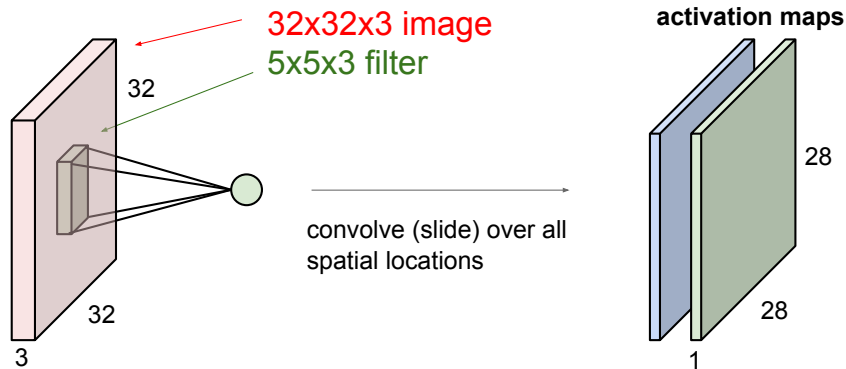
Convolution Layer



Fei-Fei Li & Justin Johnson & Serena Yeung Lecture 5 - 32 April 18, 2017

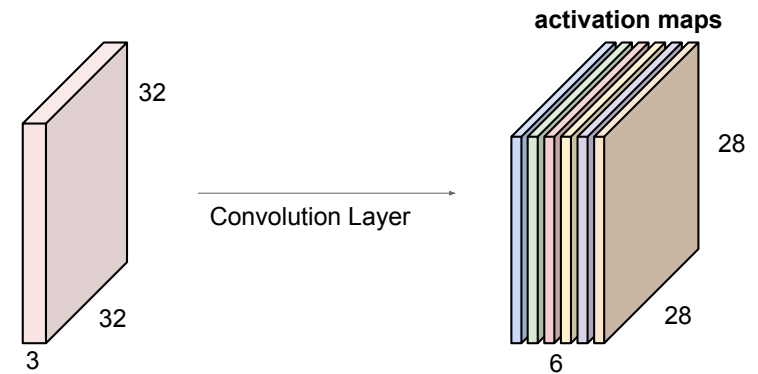
Convolution Layer

consider a second, **green** filter



Fei-Fei Li & Justin Johnson & Serena Yeung Lecture 5 - 33 April 18, 2017

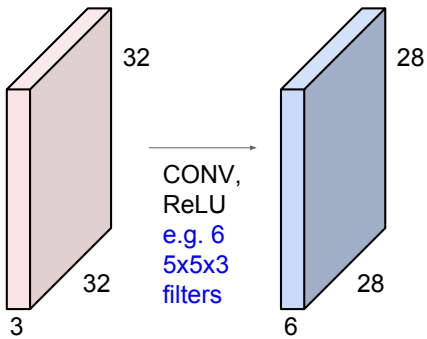
For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:



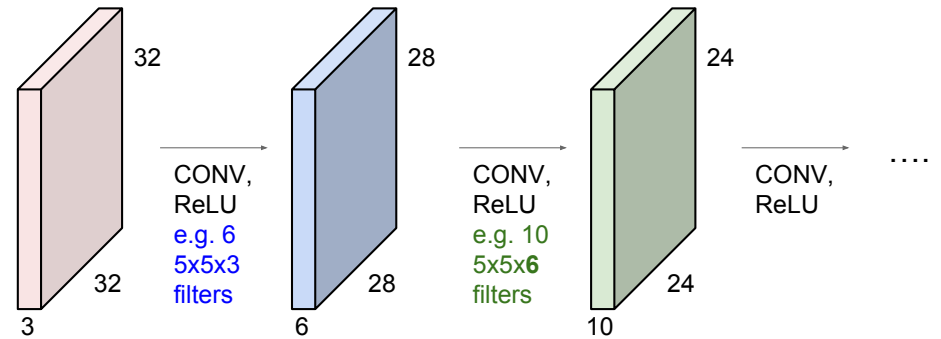
We stack these up to get a "new image" of size 28x28x6!

Fei-Fei Li & Justin Johnson & Serena Yeung Lecture 5 - 34 April 18, 2017

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



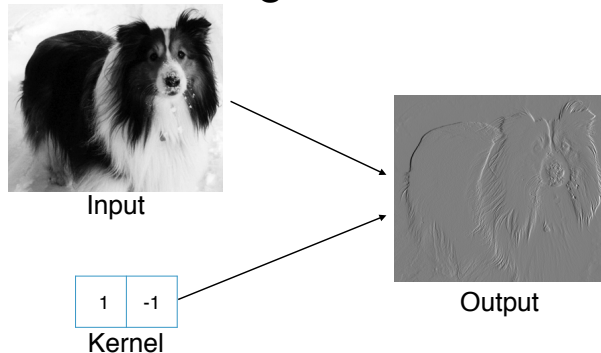
Preview: ConvNet is a sequence of Convolutional Layers, interspersed with activation functions



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



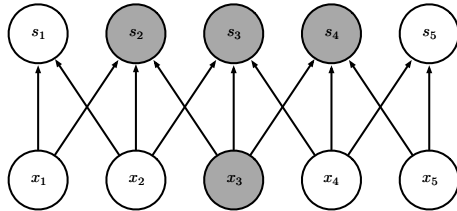
Connection to fully connected NNs

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

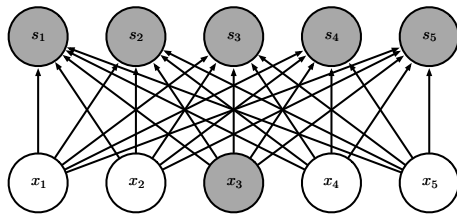
- filter = weights with **sparse connection**

Local Receptive Field Leads to Sparse Connectivity (affects less)

Sparse connections due to small convolution kernel



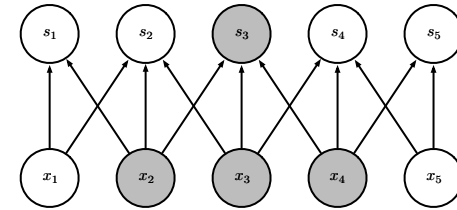
Dense connections



(Goodfellow 2016)

Sparse connectivity: being affected by less

Sparse connections due to small convolution kernel



Dense connections

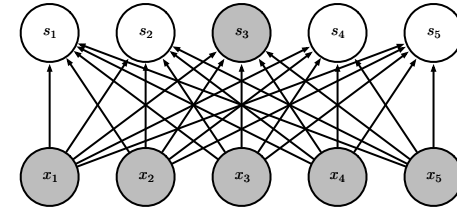


Figure 9.3

(Goodfellow 2016)

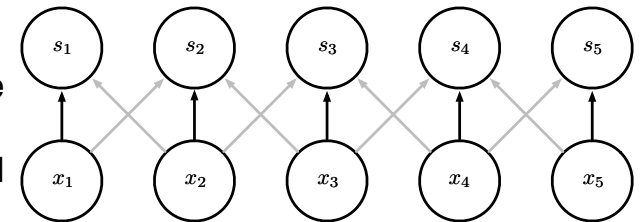
Connection to fully connected NNs

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

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

Parameter Sharing

Convolution shares the same parameters across all spatial locations



Traditional matrix multiplication does not share any parameters

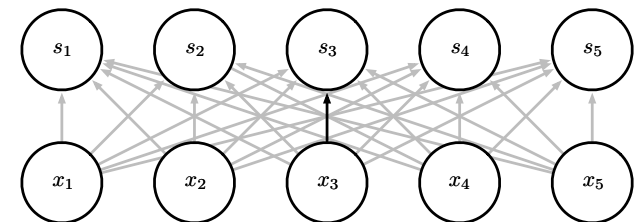


Figure 9.5

(Goodfellow 2016)

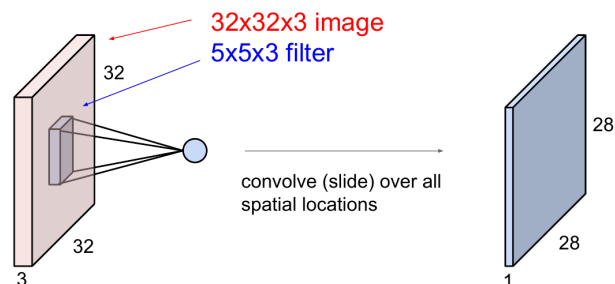
Connection to fully connected NNs

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

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

Much fewer parameters! Example (ignore bias terms):

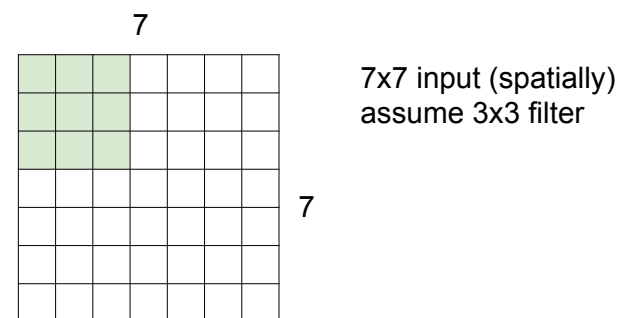
- FC: $(32 \times 32 \times 3) \times (28 \times 28) \approx 2.4M$
- CNN: $5 \times 5 \times 3 = 75$



14 / 42

Spatial arrangement: stride and padding

A closer look at spatial dimensions:



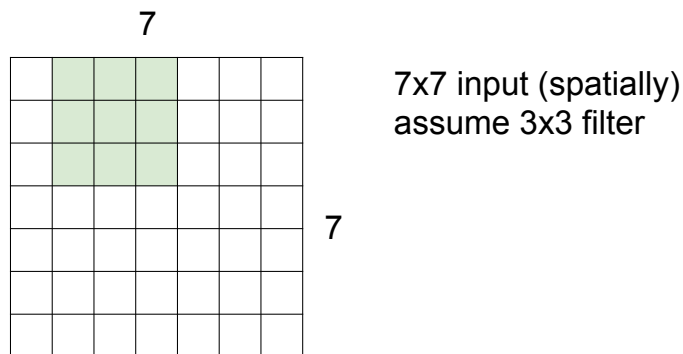
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 42

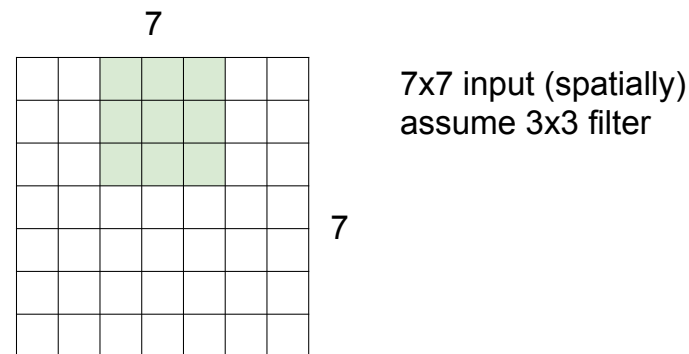
April 18, 2017

15 / 42

A closer look at spatial dimensions:



A closer look at spatial dimensions:



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 43

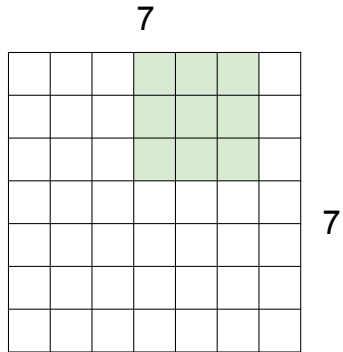
April 18, 2017

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 44

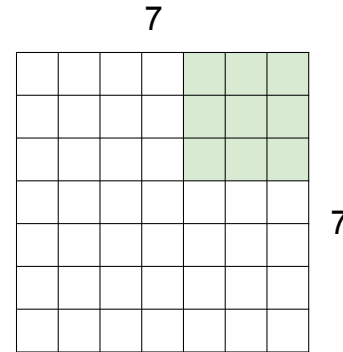
April 18, 2017

A closer look at spatial dimensions:



7x7 input (spatially)
assume 3x3 filter

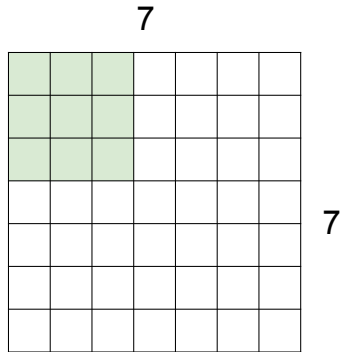
A closer look at spatial dimensions:



7x7 input (spatially)
assume 3x3 filter

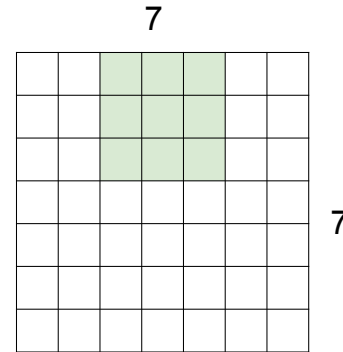
=> 5x5 output

A closer look at spatial dimensions:



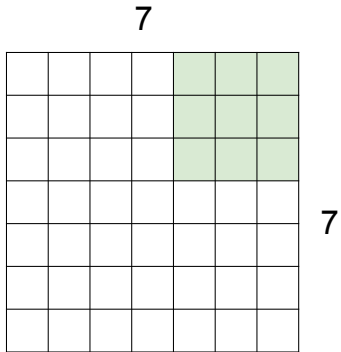
7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**

A closer look at spatial dimensions:



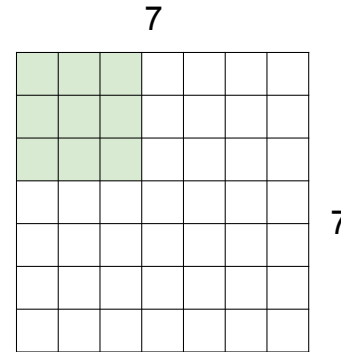
7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**

A closer look at spatial dimensions:



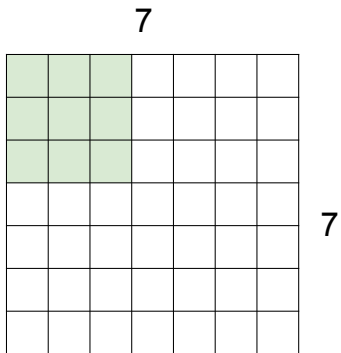
7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**
=> **3x3 output!**

A closer look at spatial dimensions:



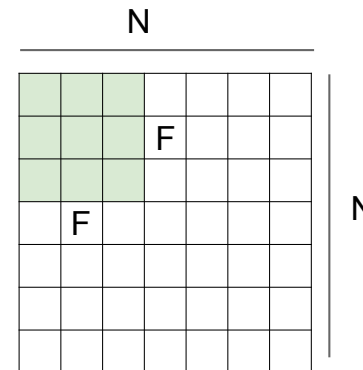
7x7 input (spatially)
assume 3x3 filter
applied **with stride 3?**

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

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

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							
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

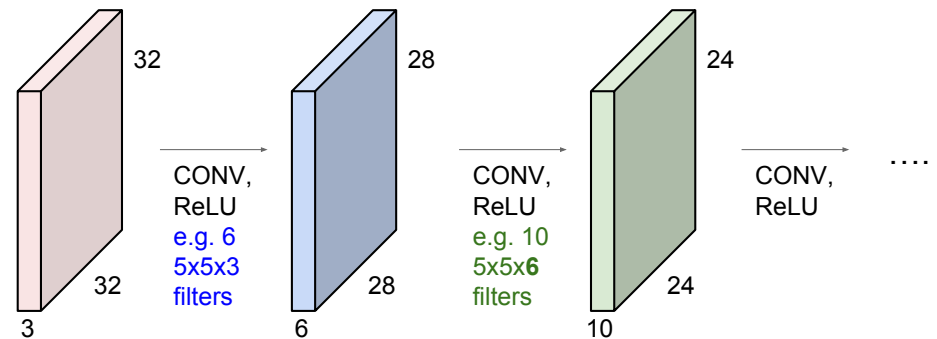
0	0	0	0	0	0		
0							
0							
0							
0							

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 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.



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

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

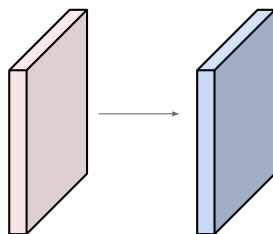
#parameters: $(F \times F \times D_1 + 1) \times K$ weights

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

16 / 42

Examples time:

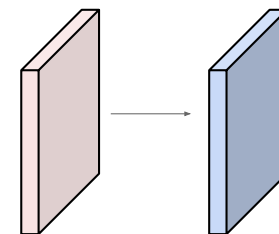
Input volume: **32x32x3**
10 5x5 filters with stride 1, pad 2



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

Examples time:

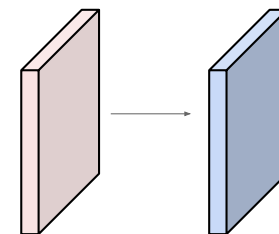
Input volume: **32x32x3**
10 5x5 filters with stride 1, pad 2



Output volume size: ?

Examples time:

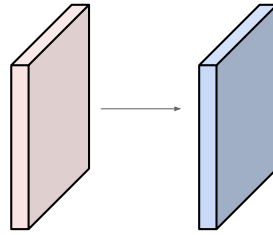
Input volume: **32x32x3**
10 5x5 filters with stride 1, pad 2



Number of parameters in this layer?

Examples time:

Input volume: $32 \times 32 \times 3$
 10 5×5 filters with stride 1, pad 2

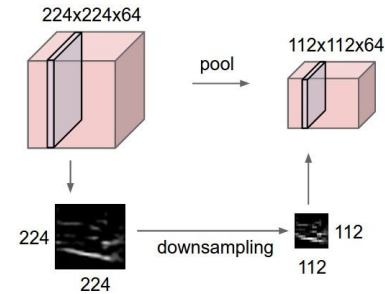


Number of parameters in this layer?
 each filter has $5 \times 5 \times 3 + 1 = 76$ params (+1 for bias)
 $\Rightarrow 76 \times 10 = 760$

Another element: pooling

Pooling layer

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



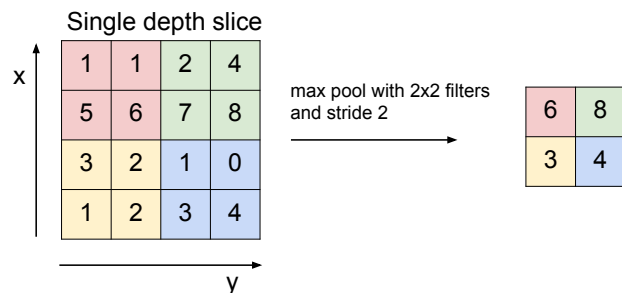
Pooling

Similar to a filter, except

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

Max pooling with 2×2 filter and stride 2 is very common

MAX POOLING



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?

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

20 / 42

Outline

- ① Convolutional neural networks (ConvNets/CNNs)
- ② Kernel methods
 - Motivation
 - Dual formulation of linear regression
 - Kernel Trick

21 / 42

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 ϕ from data
- **kernel method** is another one: sidestep the issue of choosing ϕ by using *kernel functions*

22 / 42

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!*

23 / 42

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.

Why is this helpful?

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

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

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 ϕ* .

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

How to find $\boldsymbol{\alpha}$?

Plugging in $\mathbf{w} = \Phi^T\boldsymbol{\alpha}$ into $F(\mathbf{w})$ gives

$$\begin{aligned} G(\boldsymbol{\alpha}) &= F(\Phi^T\boldsymbol{\alpha}) \\ &= \|\Phi\Phi^T\boldsymbol{\alpha} - \mathbf{y}\|_2^2 + \lambda\|\Phi^T\boldsymbol{\alpha}\|_2^2 \\ &= \|\mathbf{K}\boldsymbol{\alpha} - \mathbf{y}\|_2^2 + \lambda\boldsymbol{\alpha}^\top \mathbf{K}\boldsymbol{\alpha} \quad (\mathbf{K} = \Phi\Phi^T \in \mathbb{R}^{N \times N}) \end{aligned}$$

\mathbf{K} is called **Gram matrix** or **kernel matrix** where the (i, j) entry is

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

Example 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

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

Gram matrix vs covariance matrix

	dimensions	entry (i, j)	property
$\Phi\Phi^T$	$N \times N$	$\phi(\mathbf{x}_i)^T\phi(\mathbf{x}_j)$	both are symmetric and positive semidefinite
$\Phi^T\Phi$	$M \times M$	$\sum_{n=1}^N \phi(\mathbf{x}_n)_i\phi(\mathbf{x}_n)_j$	

28 / 42

How to find α ?

Minimize (the so-called *dual formulation*)

$$G(\alpha) = \|\mathbf{K}\alpha - \mathbf{y}\|_2^2 + \lambda\alpha^T\mathbf{K}\alpha$$

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

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

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

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

Exercise: *are there other minimizers? and are there other \mathbf{w}^* 's?*

29 / 42

Comparing two solutions

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

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

Note \mathbf{I} has different dimensions in these two formulas.

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

They have to be the same because $F(\mathbf{w})$ has a unique minimizer!

And they are:

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

30 / 42

Then what is the difference?

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

More importantly, computing $\alpha = (\mathbf{K} + \lambda\mathbf{I})^{-1}\mathbf{y}$ 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 ϕ it is indeed possible to compute $\phi(\mathbf{x})^T\phi(\mathbf{x}')$ without computing/knowing ϕ . This is the *kernel trick*.

31 / 42

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.*

32 / 42

Another example

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

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

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

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

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

33 / 42

More complicated example

Based on ϕ_θ , 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 $\phi_L(\mathbf{x}')$?

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

34 / 42

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:

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

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!*

35 / 42

Kernel functions

Definition: a function $k : \mathbb{R}^D \times \mathbb{R}^D \rightarrow \mathbb{R}$ is called a *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

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

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

Examples that are not kernels

Function

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

is *not a kernel*, why?

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?*

Using kernel functions

Choosing a nonlinear basis ϕ 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:

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

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**).

- useful for proving that a function is not a kernel

More examples of kernel functions

Two most commonly used kernel functions in practice:

Polynomial kernel

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}' + c)^d$$

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 ϕ is* for each kernel.

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

- **conical 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!

Predicting with a kernel function

As long as $\mathbf{w}^* = \sum_{n=1}^N \alpha_n \phi(\mathbf{x}_n)$, prediction on a new example \mathbf{x} becomes

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

This is a non-parametric method!

Kernelizing other ML algorithms

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

- nearest neighbor classifier
- Perceptron (HW2)
- logistic regression
- SVM (next week)
- ...