

2019 Summer School on Human Language Technology

-- Computer Vision --





Jose Dolz ETS Montreal

June, 12th, 2019



Visual recognition until 2012

Hand-crafting features



Images from: Fei-Fei Li, Andrej Karpathy and Justin Johnson 2016, cs231n.



Biological inspiration



What's new?

Multi-layer perceptron (MLP)



Many parameters Limited depth CNN



- Shared parameters
- Local connectivity

Less parameters Deep architectures possible

Fully connected layer

Image: 28x28 pixels



Fully connected layer

Image: 28x28 pixels



Convolutional layer

Image: 28x28 pixels



Convolutional filter

30	3,	22	1	0	
02	02	1,0	3	1	
3.	1,	22	2	3	
2	0	0	2	2	
2	0	0	0	1	



Convolutional layer

Image: 28x28 pixels



It can be formulated as:

• The convolution of an image x with a kernel k is computed as follows:

$$(x * k)_{ij} = \sum\limits_{pq} x_{i+p,j+q} \; k_{r\text{-}p,r\text{-}q}$$

• Example:



.28

28

Convolutional layer

Image: 28x28 pixels

Image: 28x28x1 Filter: 5x5

- **Result:** Scalar product between the filter and the image region.

Convolutional layer

Image: 28x28 pixels



Feature map

Convolutional layer (We have several filters!!!)



Feature map

Convolutional layer (We have several filters!!!)



Convolutional neural network (simplified)









No padding (unit stride)

Reduced dimension along border

Padding (unit stride)

Preserved spatial dimension



Basics: Pooling



Goals

- Reduce the spatial resolution of feature maps
- Lower memory and computation requirements
- Provide partial invariance to position, scale and rotation

Pooling is typically done separately for each feature map

Convolutional neural network (simplified)



Convolutional neural network (simplified)



Basics: Activation functions



Ensures that neurons are always activated



 $\begin{array}{l} \textbf{Maxout} \\ \max(w_1^T x + b_1, w_2^T x + b_2) \end{array}$



Convolutional neural network (simplified)



Basics: Receptive field



- Region in the input space which can the neuron's output
- A large receptive is necessary to capture spatial context
- The receptive field of neurons increases when going deeper in the network.

Basics: Receptive field



Question: How to get a large receptive field without having too many network parameters ?

Basics: Dilated convolutions (à trous)



- Inflates the kernel by inserting spaces between the kernel elements
- Controlled by dilation parameter *d* (*d* = 1 gives a regular convolution)

More in optimization

'Simply' an optimization problem



Gradient descent

 $\theta = \theta - \eta \cdot \nabla_{\theta} J(\theta)$



Optimization

for i in range(nb_epochs):
params_grad = evaluate_gradient(loss_function, data, params)
params = params - learning_rate * params_grad

□ It is very slow

□ Intractable for data that do not fit in memory

Stochastic Gradient descent (SGD)

 $\theta = \theta - \eta \cdot \nabla_{\theta} J(\theta; x^{(n)}; y^{(n)})$

Optimization



for i in range(nb_epochs):
np.random.shuffle(data)
for batch in get_batches(data, batch_size=50):
 params_grad = evaluate_gradient(loss_function, batch, params)
 params = params - learning_rate * params_grad



Stochastic Gradient descent (SGD) Problems

'Condition number'



 $f(x,z) = x^2 + az^2$ $H_{xx} = \frac{\partial^2 f}{\partial x^2} = 2$ $H_{xz} = \frac{\partial^2 f}{\partial x \partial z} = 0$ $H_{zx} = \frac{\partial^2 f}{\partial z \partial x} = 0 \quad H_{zz} = \frac{\partial^2 f}{\partial z^2} = 2a$ $\frac{H_{zz}}{H_{xx}} = a$

Optimization



Stochastic Gradient descent (SGD) Problems

Optimization



'Condition number'





Stochastic Gradient descent (SGD) Problems Optimization



'Condition number'



Evolution towards the minimum becomes slow



Stochastic Gradient descent (SGD) Problems







Batch gradient descent
Mini-batch gradient Descent
Stochastic gradient descent

SGD + Momentum

Standard gradient

 $\theta^{t+1} = \theta^t - \eta \cdot \nabla J(\theta^t)$

 $\theta^{t+1} = \theta^t - \eta v^{t+1}$

With momentum

Optimization

 $v^{t+1} = \rho v^t + \nabla J(\theta^t)$

Benefits:

- Accelerates learning when gradient direction is stable
- Reduces oscillations during training
- Faster convergence, possibly better solution
SGD + Momentum

Optimization

Standard gradient

for i in range (nb_epochs):
 dx = compute_grad(x)
 x = x - learning_rate * dx

With momentum

```
vel = 0
for i in range (nb_epochs):
    dx = compute_grad(x)
    vel = vel * rho + dx
    x = x - learning_rate * vel
```



Nesterov Momentum

Optimization





Source: http://cs231n.github.io/assets/nn3/nesterov.jpeg

Adaptative moment estimation (Adam)

Optimization

$$m_t = \gamma_1 m_{t-1} + (1 - \gamma_1) \frac{\partial \mathcal{L}}{\partial w}$$

Decaying average of gradient <u>first</u> moment (mean)

$$v_t = \gamma_2 v_{t-1} + (1 - \gamma_2) \frac{\partial \mathcal{L}^2}{\partial w}$$

Decaying average of gradient <u>second</u> moment (~variance)

$$\hat{v}_t = \frac{v_t}{1 - \gamma_2^t} \qquad \hat{m}_t = \frac{m_t}{1 - \gamma_1^t}$$

Correct values to avoid bias to zero in first steps

$$heta^{t+1} = heta^t - rac{\eta}{\sqrt{v^t + \epsilon}} m^t$$



Second moment

$$v_t = \gamma_2 v_{t-1} + (1 - \gamma_2) \frac{\partial \mathcal{L}^2}{\partial w}$$



Comparison



Optimization

Batch normalization

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_{1...m}\}$; Parameters to be learned: γ, β **Output:** $\{y_i = BN_{\gamma,\beta}(x_i)\}$

 $\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^{m} x_i \qquad // \text{ mini-batch mean}$ $\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^{m} (x_i - \mu_{\mathcal{B}})^2 \qquad // \text{ mini-batch variance}$ $\widehat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \qquad // \text{ normalize}$ $y_i \leftarrow \gamma \widehat{x}_i + \beta \equiv \mathbf{BN}_{\gamma,\beta}(x_i) \qquad // \text{ scale and shift}$

Key ideas:

- Normalize the output of a layer using batch mean and standard deviation (mean of 0, standard deviation of 1)
- Optimal scale and shift is learned by the network (parameters gamma and beta)
- Stabilizes the learning process, faster convergence

Dropout



Standard network

With dropout

 ∞

 \otimes

 \propto

 \otimes

 \otimes

 \otimes

 ∞

Key ideas:

- Regularization technique to avoid overfitting
- Ignore some neurons at each iteration (forward and backward passes)
- Dropped neurons chosen randomly with probability p
- Avoid co-adaptation of neurons: each neuron should work independently of others

How prediction is done?

Example: classification



CNN output (Logits)

Ciel 2.45

Chien **4.29**

Velo -1.24

How prediction is done?

Example: classification



How the objective is evaluated?
Example: classification
Cost function: cross-entropy
$$\checkmark$$
 Widely used
 $KL(Q||P) = \sum_{t} Q(t|\mathbf{x}) \log \frac{Q(t|\mathbf{x})}{P(t|\mathbf{x})} = -\sum_{k=1}^{K} I(k, t_{class}) \log P(t = k|\mathbf{x})$
Kullback-Leibler
Divergence \checkmark
 $P(t = k|\mathbf{x}) = \frac{e^{zk}}{\sum_{j} e^{z_j}}$

Back-propagating the error

The actual core of learning

- f(x, y, z) = (x + y)z
 - x=-2 y=5 z=-4 z^{-4}















COMPUTER VISION TASKS

- Object classification
- Semantic segmentation
- Object detection

OBJECT CLASSIFICATION



Goal: predict a single label (or a probability distribution over labels) for a given image.

Challenges

	37	49	43	43	63	45	51	56	65	59
	47	64	68	37	48	56	37	47	61	47
	56	67	64	39	80	66	31	48	49	33
	38	49	32	75	48	49	71	35	47	27
	61	62	33	64	60	10	35	10	70	10
	52	22	21	56	24	22	24	27	12	26
	52	52	20	20	34	22	54	27	45	50
	34	//	26	30	46	27	62	/6	70	65
	69	36	49	31	34	41	75	61	69	46
	31	40	62	30	67	43	54	77	72	72
	42	69	65	76	73	61	64	34	53	66
	39	52	55	64	45	78	34	76	60	57
	55	68	31	56	63	77	78	69	78	35
	31	47	60	43	42	51	55	57	50	78
	48	41	32	35	55	39	63	50	29	40
	31	28	33	60	71	68	28	40	73	76
and a second	32	63	31	80	58	67	70	67	60	38
	60	40	22	25	11	66	67	20	45	46
	42	49	35	35	44	00	22	20	45	40
	42	50	35	40	42	00	32	29	80	30
	59	59	36	47	50	31	54	68	38	61
	79	29	43	30	49	63	43	62	61	35

What humans see

What the computer sees

Other challenges



Easy for humans (and lesser animals). Hard for computers?

Image: http://cs231n.github.io

IM GENET Large Scale Visual Recognition Challenge The Image Classification Challenge: 1,000 object classes 1,431,167 images



Image: http://cs231n.stanford.edu/slides/2018

AlexNet [Krizhevsky et al, 2012]



AlexNet [Krizhevsky et al, 2012]



- First model to perform well on the challenging ImageNet dataset.
- Combined techniques used in today's architectures, like ReLU, data augmentation and dropout
- Largely responsible for the rise of deep learning in computer vision

Image: Krizhevsky et al. "Imagenet classification with deep convolutional neural networks," NIPS 2012.

ZF Net [Zeiler and Fergus, 2014]



Image: Zeiler and Fergus. "Visualizing and understanding convolutional networks." ECCV, 2014.

VGG [Simonyan and Zisserman, 2014]



VGG [Simonyan and Zisserman, 2014]



- **Simpler structure**: only 3x3 convolutions, ReLU and 2x2 max pooling
- **Deeper network**: 16 and 19 layers (compared to 8 for AlexNet)
- **Key idea**: cascading two 3x3 convolutions gives the same receptive field as a 5x5 convolution, with much less parameters

Image: http://cs231n.stanford.edu/slides/2018

GoogLeNet (inception v1) [Szegedy et al, 2014]





GoogLeNet (inception v1) [Szegedy et al, 2014]

Let's take a closer look at Inception modules...



Choice for each layer:

- Convolution or pooling ?
- If convolution, what kernel size ?

GoogLeNet (inception v1) [Szegedy et al, 2014]

Let's take a closer look at Inception modules...



Key idea:

- Compute all in parallel
- Concatenate results
- Let the learning decide

Problem: this gives too many outputs and parameters

GoogLeNet (inception v1) [Szegedy et al, 2014]

Let's take a closer look at Inception modules...



Solution: Reduce dimensionality using *bottleneck layers* composed of 1x1 convolutions

Efficiency





~5M parameters

GoogLeNet has 12x fewer parameters than AlexNet !

GoogLeNet (inception v1) [Szegedy et al, 2014]



Image: http://cs231n.stanford.edu/slides/2018






• Enables very deep networks (over 100 layers)



SEMANTIC SEGMENTATION





Image

Segmentation

Goal: Assign the correct class label to each pixel of a given image

Can be seen as a <u>dense</u> and <u>structured</u> classification problem

Predictions are not independent

Possibly millions of pixels at the same time



Key ideas:

- Replace FC layers by 1x1 convolutions
- Treat network as a non-linear filter that can be applied to any image size

Problem: Produces very coarse segmentation maps

Solution: Add upsampling operations at the end of the network

Upsampling

Convolution as a matrix operation:

$$\begin{pmatrix} f_1 & f_2 & f_3 \end{pmatrix} * \begin{pmatrix} g_1 \\ g_2 \\ g_3 \\ g_4 \\ g_5 \end{pmatrix} = \begin{pmatrix} f_1 & f_2 & f_3 & 0 & 0 \\ 0 & f_1 & f_2 & f_3 & 0 \\ 0 & 0 & f_1 & f_2 & f_3 \end{pmatrix} \begin{pmatrix} g_1 \\ g_2 \\ g_3 \\ g_4 \\ g_5 \end{pmatrix}$$

Transposed convolution:



Image: http://deeplearning.net/software/theano/tutorial



Problem: Spatial resolution is lost while downsampling

Solution: Add skip connections to copy high-resolution feature maps of the encoder to the decoder

Image adapted from: Badrinarayanan et al. "Segnet: A deep convolutional encoder-decoder architecture for image segmentation." PAMI, 2017.



One of the most popular networks for medical image segmentation



Multi-modal segmentation

[Dolz et al, 2018]



CNN

Problem: Ignores low or high level dependencies between information in different image modalities

Image: Dolz et al. "IVD-Net: Intervertebral disc localization and segmentation in MRI with a multi-modal UNet." MICCAI Workshops, 2018.

Multi-modal segmentation

HyperDense-Net [Dolz et al, 2018]

T1 from T1





T1 from T2

Key ideas:

- Dense connections across layers of different modalities
- Let training decide how to best combine information across modalities (see image on the right) ٠

OBJECT DETECTION



Goal: Find <u>objects</u> (or people, animals) in a given image, and their <u>location</u>.

Object detection

R-CNN [Girshick et al, 2014]



- **Problems**:
- Multiple training stages: CNN, SVM, bounding box regressors
- Very slow (~53s per image)

- 5. Non-maxima suppression (Removes redundant overlapping regions)
- 4. Classification & localization
- 3. Compute CNN features (AlexNet)
- 2. Warp proposals to fixed size
- 1. Extract region proposals (~2K)
- Uses the Selective Search method
- Finds regions likely to contain objects

88

Object detection

Fast R-CNN [Girshick et al, 2015]



Key idea:

- Feed the input image to the CNN (instead of region proposals)
- Enables sharing CNN computations between different proposals

Problem: Still requires extracting proposals in a separate step

How to avoid this ?

Image: Girshick, Ross. "Fast R-CNN." ICCV. 2015.

Object detection

Faster R-CNN [Ren et al, 2015]



Images: Ren et al. "Faster R-CNN: Towards real-time object detection with region proposal networks." NIPS. 2015.