



Neural Networks MLSS 2015 Summer School

Rob Fergus

Facebook AI Research



[On leave from Courant Institute, New York University]



Overview

- Look at some of the recent progress with Neural Net / Deep Learning models
 Assume familiarity with basic neural nets
- Non-exhaustive coverage
 Huge number of recent papers
- Draw on material from FAIR colleagues

 Experts in NLP, speech etc.

Thanks!

Facebook AI Research colleagues & NYU PhD students:







Yaniv Taigman







Emily Denton







Jason Weston Tomas Mikolov Ronan Collobert

Marc'Aurelio Sainbayar Ranzato Sukhbaatar

Schedule

- Overview
- Vision

[~15m] [~1h50m]

- Convnets for image recognition
- Other applications
- Speech
- NLP (RNNs, LSTMs)
- Unsupervised Learning
- Memory in neural nets

[~5m]

[~20m] [~1h] [~30m]

Motivation

- A lot of recent successes with large neural networks, trained with supervision.
- Feature learning crucial to performance in many tasks
- Still many open problems however!

Traditional ML Approach

For classification:



- Features are not learned
- Trainable classifier is often generic (e.g. SVM)

Case Study: Object Recognition ~2010

- Multitude of hand-designed features currently in use – SIFT, HOG, LBP, MSER, Color-SIFT.....
- Where next? Better classifiers? Or keep building more features?



Felzenszwalb, Girshick, McAllester and Ramanan, PAMI 2007



Yan & Huang (Winner of PASCAL 2010 classification competition)

Hand-Crafted Features + SVM

- LP-β Multiple Kernel Learning
 - Gehler and Nowozin, On Feature Combination for Multiclass Object Classification, ICCV'09
- 39 different kernels
 PHOG, SIFT, V1S+, Region Cov. Etc.
- MKL only gets few % gain over averaging features
- → Features are doing the work



What Limits Performance?

- Ablation studies on Deformable Parts Model
 - Felzenszwalb, Girshick, McAllester, Ramanan, PAMI'10
- Replace each part with humans (Amazon Turk):



What about Learning the Features?

- Learn hierarchy
- All the way from pixels \rightarrow classifier
- One layer extracts features from output of previous layer



• Train all layers jointly

Deep Learning



Deep Learning is BIG

Main types of deep architectures



Deep Learning is **BIG**

Main types of learning protocols

- Purely supervised
 - Backprop + SGD
 - Good when there is lots of labeled data.
- Layer-wise unsupervised + superv. linear classifier
 - Train each layer in sequence using regularized auto-encoders or RBMs
 - Hold fix the feature extractor, train linear classifier on features
 - Good when labeled data is scarce but there is lots of unlabeled data.
- Layer-wise unsupervised + supervised backprop
 - Train each layer in sequence
 - Backprop through the whole system
 - Good when learning problem is very difficult.



Deep Learning is BIG

Main types of learning protocols

- Purely supervised
 - Backprop + SGD
 - Good when there is lots of labeled data.
- Layer-wise unsupervised + superv. linear classifier
 - Train each layer in sequence using regularized auto-encoders or RBMs
 - Hold fix the feature extractor, train linear classifier on features
 - Good when labeled data is scarce but there is lots of unlabeled data.
- Layer-wise unsupervised + supervised backprop
 - Train each layer in sequence
 - Backprop through the whole system
 - Good when learning problem is very difficult.



Deep Learning for Computer Vision



Convolutional Neural Networks

- LeCun et al. 1989
- Neural network with specialized connectivity structure
 - [Everyone OK with basic NN?]





Multistage Hubel-Wiesel Architecture

- Stack multiple stages of simple cells / complex cells layers
- Higher stages compute more global, more invariant features
- Classification layer on top

History:

- Neocognitron [Fukushima 1971-1982]
- Convolutional Nets [LeCun 1988-2007]
- HMAX [Poggio 2002-2006]
- Many others....



Overview of Convnets

Feature maps

Pooling

Non-linearity

Convolution (Learned)

- Feed-forward:
 - Convolve input
 - Non-linearity (rectified linear)
 - Pooling (local max)
- Supervised
- Train convolutional filters by back-propagating classification error



Convnet Successes

- Handwritten text/digits
 - MNIST (0.17% error [Ciresan et al. 2011])
 - Arabic & Chinese [Ciresan et al. 2012]
- Simpler recognition benchmarks
 - CIFAR-10 (9.3% error [Wan et al. 2013])
 - Traffic sign recognition
 - 0.56% error vs 1.16% for humans [Ciresan et al. 2011]
- But less good at more complex datasets
 - E.g. Caltech-101/256 (few training examples)





Application to ImageNet



- ~14 million labeled images, 20k classes
- Images gathered from Internet
- Human labels via Amazon Turk

ImageNet Classification with Deep Convolutional Neural Networks [NIPS 2012]

Alex Krizhevsky University of Toronto kriz@cs.utoronto.ca Ilya Sutskever University of Toronto ilya@cs.utoronto.ca Geoffrey E. Hinton University of Toronto hinton@cs.utoronto.ca

Goal

Image Recognition
 – Pixels → Class Label

 $\bullet \bullet \bullet \bullet \bullet \bullet \bullet \bullet$

lens cap	abacus	slug	hen
reflex camera	abacus	slug	hen
Polaroid camera	typewriter keyboard	zucchini	cock
pencil sharpener	space bar	ground beetle	cocker spaniel
switch	computer keyboard	common newt	partridge
combination lock	accordion	water snake	English setter

[Krizhevsky et al. NIPS 2012]

• • •

Krizhevsky et al. [NIPS2012]

- Same model as LeCun'98 but:
 - Bigger model (8 layers)
 - More data $(10^6 \text{ vs } 10^3 \text{ images})$
 - GPU implementation (50x speedup over CPU)
 - Better regularization (DropOut)



- 7 hidden layers, 650,000 neurons, 60,000,000 parameters
- Trained on 2 GPUs for a week

ImageNet Classification (2010 – 2015)



Examples

• From Clarifai.com



Predicted Tags:

. . .

food	(16.00%)
dinner	(3.10%)
bbq	(2.90%)
market	(2.50%)
meal	(1.40%)
turkey	(1.40%)
grill	(1.30%)
pizza	(1.30%)
eat	(1.10%)
holiday	(1.00%)

Stats:

Size: 247.24 KB Time: 110 ms

Examples

• From Clarifai.com



Predicted Tags:

ship	(2.30%)
helsinki	(1.80%)
fish	(1.40%)
port	(1.10%)
istanbul	(1.10%)
beach	(1.00%)
denmark	(1.00%)
copenhagen	(0.90%)
sea	(0.80%)
boat	(0.80%)

Examples

• From Clarifai.com



Predicted Tags:

 \bullet \bullet \bullet

barcelona	(6.50%)
street	(3.00%)
cave	(2.20%)
sagrada	(1.90%)
old	(1.80%)
night	(1.40%)
familia	(1.40%)
jerusalem	(1.40%)
guanajuato	(1.10%)
alley	(1.00%)

Stats:

Size: 278.96 KB Time: 113 ms

Using Features on Other Datasets

- Train model on ImageNet 2012 training set
- Re-train classifier on new dataset
 Just the top layer (softmax)

• Classify test set of new dataset

Caltech 256

Zeiler & Fergus, Visualizing and Understanding Convolutional Networks, arXiv 1311.2901, 2013



Caltech 256

Zeiler & Fergus, Visualizing and Understanding Convolutional Networks, arXiv 1311.2901, 2013



The Details

- Operations in each layer
- Architecture
- Training
- Results



Filtering

• Convolution

- Filter is learned during training
- Same filter at each location







Filtering

- Local
 - Each unit layer above
 look at local window
 - But no weight tying





• E.g. face recognition



Filtering

- Tiled
 - Filters repeat every n
 - More filters than convolution for given # features





Filters





Feature maps
Non-Linearity

- Rectified linear function
 - Applied per-pixel
 - output = max(0,input)

Input feature map





Output feature map



Non-Linearity

- Other choices:
 - Tanh
 - Sigmoid: $1/(1+\exp(-x))$

– PReLU

[Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification, Kaiming He et al. arXiv:1502.01852v1.pdf, Feb 2015]







Pooling

- Spatial Pooling
 - Non-overlapping / overlapping regions
 - Sum or max
 - Boureau et al. ICML'10 for theoretical analysis







Pooling

- Pooling across feature groups
 - Additional form of inter-feature competition
 - MaxOut Networks [Goodfellow et al. ICML 2013]



Role of Pooling

- Spatial pooling
 - Invariance to small transformations
 - Larger receptive fields (see more of input)

Visualization technique from [Le et al. NIPS'10]:





Zeiler, Fergus [arXiv 2013]

Videos from: http://ai.stanford.edu/~quocle/TCNNweb



Normalization

- Contrast normalization across features
 - See Divisive Normalization in Neuroscience





Filters

Normalization

Contrast normalization (across feature maps)
 Local mean = 0, local std. = 1, "Local" → 7x7 Gaussian
 Equalizes the features maps



Feature Maps

Feature Maps After Contrast Normalization

Role of Feature Normalization

- Introduces local competition between features
 - "Explaining away" in graphical models
 - Just like top-down models
 - But more local mechanism
- Also helps to scale activations at each layer better for learning
 - Makes energy surface more isotropic
 - So each gradient step makes more progress

- Empirically, seems to help a bit (1-2%) on ImageNet
- Most recent models don't seem to have use though

Normalization across Data

• Batch Normalization

[Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift, Sergey Ioffe, Christian Szegedy, arXiv:1502.03167]

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_{1...m}\}$; Parameters to be learned: γ, β Output: $\{y_i = BN_{\gamma,\beta}(x_i)\}$ $\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i$ // mini-batch mean $\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2$ // mini-batch variance $\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$ // normalize $y_i \leftarrow \gamma \hat{x}_i + \beta \equiv BN_{\gamma,\beta}(x_i)$ // scale and shift

Algorithm 1: Batch Normalizing Transform, applied to activation *x* over a mini-batch.



Figure 2: Single crop validation accuracy of Inception and its batch-normalized variants, vs. the number of training steps.

Overview of Convnets

Feature maps

Pooling

Non-linearity

Convolution (Learned)

- Feed-forward:
 - Convolve input
 - Non-linearity (rectified linear)
 - Pooling (local max)
- Supervised
- Train convolutional filters by back-propagating classification error



Architecture

- Big issue: how to select
 - Manual tuning of features → manual tuning of architechtures
- Depth
- Width
- Parameter count

How to Choose Architecture

- Many hyper-parameters: - # layers, # feature maps
- Cross-validation
- Grid search (need lots of GPUs)
- Smarter strategies:
 - Random [Bergstra & Bengio JMLR 2012]
 - Gaussian processes [Hinton??]

How important is Depth

- "Deep" in Deep Learning
- Ablation study
- Tap off features

- 8 layers total
- Trained on Imagenet dataset [Deng et al. CVPR'09]
- 18.2% top-5 error
- Our reimplementation: 18.1% top-5 error



- Remove top fully connected layer
 Layer 7
- Drop 16 million parameters
- Only 1.1% drop in performance!



- Remove both fully connected layers
 Layer 6 & 7
- Drop ~50 million parameters
- 5.7% drop in performance



- Now try removing upper feature extractor layers: – Layers 3 & 4
- Drop ~1 million parameters
- 3.0% drop in performance



- Now try removing upper feature extractor layers & fully connected: – Layers 3, 4, 6,7
- Now only 4 layers
- 33.5% drop in performance

 \rightarrow Depth of network is key



Tapping off Features at each Layer

Plug features from each layer into linear SVM or soft-max

	Cal-101	Cal-256		
	(30/class)	(60/class)		
SVM (1)	44.8 ± 0.7	24.6 ± 0.4		
SVM (2)	66.2 ± 0.5	39.6 ± 0.3		
SVM (3)	72.3 ± 0.4	46.0 ± 0.3		
SVM (4)	76.6 ± 0.4	51.3 ± 0.1		
SVM (5)	$\bf 86.2 \pm 0.8$	65.6 ± 0.3		
SVM (7)	85.5 ± 0.4	71.7 ± 0.2		
Softmax (5)	82.9 ± 0.4	65.7 ± 0.5		
Softmax (7)	85.4 ± 0.4	72.6 ± 0.1		

Translation (Vertical)



Scale Invariance



Rotation Invariance



Very Deep Models (1)

[Very Deep Convolutional Networks for Large-Scale Image Recognition, Karen Simonyan & Andrew Zisserman, arXiv:1409.1556, 2014]

ConvNet Configuration						
А	A-LRN	В	C	D	Е	
11 weight	11 weight	13 weight	16 weight	16 weight	19 weight	
layers	layers	layers	layers	layers	layers	
	i	nput (224×2	24 RGB image	e)		
conv3-64	conv3-64	conv3-64	conv3-64	conv3-64	conv3-64	
	LRN	conv3-64	conv3-64	conv3-64	conv3-64	
		max	pool			
conv3-128	conv3-128	conv3-128	conv3-128	conv3-128	conv3-128	
		conv3-128	conv3-128	conv3-128	conv3-128	
		max	pool			
conv3-256	conv3-256	conv3-256	conv3-256	conv3-256	conv3-256	
conv3-256	conv3-256	conv3-256	conv3-256	conv3-256	conv3-256	
			conv1-256	conv3-256	conv3-256	
					conv3-256	
		max	pool			
conv3-512	conv3-512	conv3-512	conv3-512	conv3-512	conv3-512	
conv3-512	conv3-512	conv3-512	conv3-512	conv3-512	conv3-512	
			conv1-512	conv3-512	conv3-512	
					conv3-512	
		max	pool			
conv3-512	conv3-512	conv3-512	conv3-512	conv3-512	conv3-512	
conv3-512	conv3-512	conv3-512	conv3-512	conv3-512	conv3-512	
			conv1-512	conv3-512	conv3-512	
					conv3-512	
	maxpool					
FC-4096						
FC-4096						
FC-1000						
soft-max						

Table 2: Number of parameters (in milli	ons)
---	------

Network	A,A-LRN	В	С	D	E
Number of parameters	133	133	134	138	144

- Lots of 3x3 conv layers: more non-linearity than single 7x7 layer
- Close to SOA results on Imagenet: 6.8% top-5 val
- Can be hard to train

Table 3: ConvNet performance at a single test scale.				
ConvNet config. (Table 1)	smallest image side		top-1 val. error (%)	top-5 val. error (%)
	train (S)	test (Q)		
А	256	256	29.6	10.4
A-LRN	256	256	29.7	10.5
В	256	256	28.7	9.9
	256	256	28.1	9.4
С	384	384	28.1	9.3
	[256;512]	384	27.3	8.8
	256	256	27.0	8.8
D	384	384	26.8	8.7
	[256;512]	384	25.6	8.1
	256	256	27.3	9.0
E	384	384	26.9	8.7
	[256;512]	384	25.5	8.0

Very Deep Models (2)

[Going Deep with Convolutions, Szegedy et al., arXiv:1409.4842, 2014]

GoogLeNet inception module:

- 1. Multiple filter scales at each layer
- 2. Dimensionality reduction to keep computational requirements down



GoogLeNet vs Previous Models

[Going Deep with Convolutions, Szegedy et al., arXiv:1409.4842, 2014]



Zeiler-Fergus Architecture (1 tower)

[From http://image-net.org/challenges/ LSVRC/2014/slides/GoogLeNet.pptx]



Width of inception modules ranges from 256 filters (in early modules) to 1024 in top inception modules.

Can remove fully connected layers on top completely

Number of parameters is reduced to 5 million

6.7% top-5 validation error on Imagnet

[From http://image-net.org/challenges/ LSVRC/2014/slides/GoogLeNet.pptx] Computional cost is increased by less than 2X compared to Krizhevsky's network. (<1.5Bn operations/ evaluation)

Visualizing Convnets

- Want to know what they are learning
- Raw coefficients of learned filters in higher layers difficult to interpret
- Two classes of method:
 - 1. Project activations back to pixel space
 - 2. Optimize input image to maximize a particular feature map or class

Visualizing Convnets

- Projection from higher layers back to input
 - Several similar approaches:
 - Visualizing and Understanding Convolutional Networks, Matt Zeiler & Rob Fergus, ECCV 2014
 - Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps, Karen Simonyan, Andrea Vedaldi, Andrew Zisserman, arXiv 1312.6034, 2013
 - Object Detectors Emerge in Deep Scene CNNs, Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, Antonio Torralba, ICLR 2015



Details of Operation



Unpooling Operation



Layer 1 Filters



Visualizations of Higher Layers

- Use ImageNet 2012 validation set
- Push each image through network



- Take max activation from feature map associated with each filter
- Use Deconvnet to project back to pixel space
- Use pooling "switches" peculiar to that activation

Layer 1: Top-9 Patches



Layer 2: Top-1














		1000				4		No.	1	2	1 2	Sec.		10					-	10	1	A.K.	
TAN	MA	- All			1	X					2	· · ·							0		Y.C.		
TAN	1	The	100			26							X	Sr.			Č.	•	0	(6)	N.X.	Ch.	
2	- Ally		Con la	<u>(0)</u>	(A.	J.	9 0	3me	and the second s	-	194		æ	~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~		1	Deop St	D WINT VISTAN	O.S.	Ċ	1	000
-	装	1 de la				.M	*	-	Sec	N.		9	-10	4	-	1		CMANEL	- AT GANT FREE France	SAC		300	
	-	1			0	-		- 78	-	- A		- Ala		ý,	00) 00)		1		milite ecom	Wettook	-	1000	٩
P	E I		1,6	de		1	. Alex	N.		10	0		(A)			THE REAL					1	19	(A)
R	Ø	11	-	×.	2	and the second	inter .		10	(***) (**)	1		- Ale		(iii)		#		0		3	7	
েরা	(B)		al a	sit.	- And	1. 1 9 (_{1.1}		N.Nr.										10	(a)		S.	1	7
-	X					Ø	30	9			3		0		0	۲	0	(63)	2.	-	TTTT	Phys.	
						11 20									11510			10000					
٢	· · ·	0	-		145	0	No.	AN A	-	Ø.				Ø	Ø	0	(aligned)	68	90	-	WW		÷
		** ©			(4) (4) (4) (4) (4) (4) (4) (4) (4) (4)	•	ê. 9	100 N	A A	<u>بې</u>	ت چ		0			0	(E)	(* 		(C)	W. M	Q. 24	and an
•		* *			(A 12 (C))))	1. (h) 1. (h) 1. (h)	× 3: *	1 A 3	E s	*		000000000000000000000000000000000000000		() () () () ()	0		(e) (c) (c) (c) (c) (c) (c) (c) (c) (c) (c			E au m	11 M	
• • • • • • • • • • • • • • • • • • •		* * *					1 (k) (k)	je de ≫	18 (K) (K) (K)		(E) (E) 😽 💊		0 0 0				 3 3 4 4			100 (100 (100 (100 (100 (100 (100 (100	😵 😵 🕖	 2 3 3 4 4 4 	
		 ≫ ≫ ≫ >> <li< td=""><td></td><td></td><td> (2) 1/2 (3) 1/2 (4) (5) (5) </td><td></td><td></td><td>× 3:</td><td>* * 3 3 3</td><td></td><td>ke (fe (fe 🔖 👞</td><td></td><td></td><td> () <</td><td>00000</td><td></td><td>0 0 2 0 1</td><td></td><td></td><td> (a) (b) (b) (c) (c) (c) (c) (c) (c) (c) (c) (c) (c</td><td>10 A 20 20 20</td><td>11 11 11 11 11</td><td></td></li<>			 (2) 1/2 (3) 1/2 (4) (5) (5) 			× 3:	* * 3 3 3		ke (fe (fe 🔖 👞			 () <	00000		0 0 2 0 1			 (a) (b) (b) (c) (c) (c) (c) (c) (c) (c) (c) (c) (c	10 A 20 20 20	11 11 11 11 11	
					× & 1/2 1/2 1/2			* 3: * 6 *													10 A 20 20 A	11 11 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	
					14 × 15 1/2 1/2 1/2 1/2					 ▲ ▲												2/2 (2) (2) (2) (2) (2) (2) (2) (2) (2) (2	









Visualizing Convnets

- Optimize input to maximize particular ouput
 - Lots of approaches, e.g. Erhan et al. [Tech Report 2009], Le et al. [NIPS 2010].
 - Depend on initialization



- Google DeepDream [http://googleresearch.blogspot.ch/2015/06/inceptionism-going-deeperinto-neural.html]
 - Maximize "banana" output



Google DeepDream



https://photos.google.com/share/ F1QipPX0SCl7OzWilt9LnuQliattX4OUCj_8EP65_cTVnBmS1jnYgsGQAieQUc1VQWdgQ/photo/ AF1QipMYTXpt0TvZ0Q5kubkGw8VAq2isxBuL02wKZafB? key=aVBxWjhwSzg2RjJWLWRuVFBBZEN1d205bUdEMnhB

Training Big ConvNets

- Stochastic Gradient Descent
 - Compute (noisy estimate of) gradient on small batch of data & make step
 - Take as many steps as possible (even if they are noisy)
 - Large initial learning rate
 - Anneal learning rate

- Momentum
 - Variants [Sutskever ICML 2012]

Annealing of Learning Rate

- Start large, slowly reduce
- Explore different scales of energy surface



Evolution of Features During Training



Evolution of Features During Training

 $\bullet \bullet \bullet \bullet \bullet \bullet \bullet$

.

)	X	八	TA	À	*		-	0	*	*		ø	0	0
	No.	-	(internet	0		1	-			1	ý	\$	×.	ħ	
de.	Ţ	N i	3	¥.	¥	N.	¥		13	4	4	-	ø	Ø	¢
	and the	No.	and the	Q		(2)				1	*	*	a la	*	×.
	Ĭ	(internet)	(J)	Z		À	-			S.W.	The second		19	100	Ø
4	19	All and a second	Sec.	S.	R	R.	R.	. No	NO.	1	-	0.	(B)	- AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA	<i>sig</i>
d a) }	~	K	C	Ô	C	C			<i>I</i>	1	Ø	U	1	(13)
illy	- Mark	Jellin	01(10	SIRUSI	(OU) SUM (P	(000)	(000)			4	16	*	NC.	×	S)
1. Contraction of the second s	A NUMBER	1	X	N.	X	X.				1.110	***	1	22	123	122
r#	alle .	N.	\$	and an	Site	Sille.				(S)	k	8	(M	0	
	4	and a	10		ŧ		- af-				(22)			0	3
at.	1			-	1920	600	100			*	1	Č\$	<i>tif</i>	i	淌
1	>	-	*	La	ayer 4	×	1			-	N	۶) La	ayer 5	-9	-9

Normalization across Data

• Batch Normalization

[Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift, Sergey Ioffe, Christian Szegedy, arXiv:1502.03167]

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_{1...m}\}$; Parameters to be learned: γ, β Output: $\{y_i = BN_{\gamma,\beta}(x_i)\}$ $\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i$ // mini-batch mean $\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2$ // mini-batch variance $\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$ // normalize $y_i \leftarrow \gamma \hat{x}_i + \beta \equiv BN_{\gamma,\beta}(x_i)$ // scale and shift

Algorithm 1: Batch Normalizing Transform, applied to activation *x* over a mini-batch.



Figure 2: Single crop validation accuracy of Inception and its batch-normalized variants, vs. the number of training steps.

Automatic Tuning of Learning Rate?

• ADAGRAD

J. Duchi, E. Hazan, and Y. Singer, "Adaptive subgradient methods for online leaning and stochastic optimization," in COLT, 2010.

$$\Delta x_t = -\frac{\eta}{\sqrt{\sum_{\tau=1}^t g_\tau^2}} g_t$$

 $\Delta x_t = -\frac{\mathrm{RMS}[\Delta x]_{t-1}}{\mathrm{RMS}[a]_t} g_t$

• ADADELTA

ADADELTA: An Adaptive Learning Rate Method, Matthew D. Zeiler, arXiv 1212.5701, 2012.

• No more pesky learning rates

 $\Delta x_t = -\frac{1}{|\text{diag}(H_t)|} \frac{E[g_{t-w:t}]^2}{E[g_{t-w:t}^2]} g_t$

T. Schaul, S. Zhang, and Y. LeCun, "No more pesky learning rates," arXiv:1206.1106, 2012.

Local Minima?

[The Loss Surfaces of Multilayer Networks Choromanska et al. http://arxiv.org/pdf/1412.0233v3.pdf]



What about 2nd order methods?

- Newton's method: $\Delta x_t = H_t^{-1}g_t$
- Full Hessian impractical to compute
- Approximations:
 - Diagonal [Becker & Lecun '88]
 - Truncated CG [Martens, ICML'10]
 - Per-batch low-rank [Sohl-Dickstien et al., ICML'14]
 - Saddle free (|H|) [Dauphin et al. NIPS'14]
- Generally, extra computation needed seems not worth it: take more (dumb) steps instead!

$$\Delta x_t = -\frac{1}{|\mathrm{diag}(H_t)| + \mu} g_t$$

Saddle Point Perspective

[Identifying and attacking the saddle point problem in high-dimensional nonconvex optimization, Dauphin et al., NIPS 2014]

- During optimization Hessian has both +ve and -ve eigenvalues
 - and maybe some zeros too (flat directions)
 - At minimum, all are +ve
- Cause problems for SGD
- Saddle Free Newton (SFN)
 - Use |H| (matrix where take absolute value of each eigenvalue of H)



Improving Generalization

- Data Augmentation (jitter, peturb)
- Weight decay (L1/2 penalty on weights)
- Weight sharing (reduces # parameters)
- Multi-task learning
- Inject Noise into network
 - DropOut [Hinton et al. 2012]
 - DropConnect [Wan et al. ICML 2012]
 - Stochastic Pooling [Zeiler & Fergus ICLR'13]

Big Model + Regularize vs Small Model



Fooling Convnets

- Search for images that are misclassified by the network
- Intriguing properties of neural networks, Christian Szegedy et al. arXiv 1312.6199, 2013
- Deep Neural Networks are Easily Fooled: High Confidence Predictions for Unrecognizable Images, Anh Nguyen, Jason Yosinski, Jeff Clune, arXiv 1412.1897.
- Problem common to any discriminative method



Figure 1. Evolved images that are unrecognizable to humans, but that state-of-the-art DNNs trained on ImageNet believe with $\geq 99.6\%$ certainty to be a familiar object. This result highlights differences between how DNNs and humans recognize objects.

DropOut

- G. E. Hinton, N. Srivastava, A. Krizhevsky, I. Sutskever and R. R. Salakhutdinov, *Improving neural networks by preventing co-adaptation of feature detectors*, arXiv:1207.0580 2012
- Fully connected layers only
- Randomly set activations in layer to zero
- Gives ensemble of models
- Similar to bagging [Breiman'94], but differs in that parameters are shared.



DropConnect

- Wan et al. ICML 2013
- Fully-connected layers only
- Random binary mask on weights



Stochastic Pooling

• For conv layers [Zeiler and Fergus, ICLR 2013]

 d_{i}

- Compute activations a_i : (≥ 0)
- Normalize to sum to 1 $\rightarrow p_i = \frac{1}{2}$
- Sample location, *l*, from multinomial $\sum_{k \in R_j} a_k$
- Use activation from the location: $s = a_l$



Check gradients numerically by finite differences

Visualize features (feature maps need to be uncorrelated) and have high variance.



hidden unit

Good training: hidden units are sparse across samples and across features.



Check gradients numerically by finite differences

Visualize features (feature maps need to be uncorrelated) and have high variance.



hidden unit

Bad training: many hidden units ignore the input and/or exhibit strong correlations.



Check gradients numerically by finite differences

Visualize features (feature maps need to be uncorrelated) and have high variance.

Visualize parameters





- Check gradients numerically by finite differences
- Visualize features (feature maps need to be uncorrelated) and have high variance.
- Visualize parameters
- Measure error on both training and validation set.
- Test on a small subset of the data and check the error \rightarrow 0.



WHAT IF IT DOES NOT WORK?

Training diverges:

- Learning rate may be too large \rightarrow decrease learning rate
- \blacksquare BPROP is buggy \rightarrow numerical gradient checking
- Parameters collapse / loss is minimized but accuracy is low
 - Check loss function:
 - Is it appropriate for the task you want to solve?
 - Does it have degenerate solutions? Check "pull-up" term.
- Network is underperforming
 - Compute flops and nr. params. \rightarrow if too small, make net larger
 - $\hfill \label{eq:stable}$ $\hfill \label{eq:stable}$ Visualize hidden units/params \rightarrow fix optmization
- Network is too slow
 - Compute flops and nr. params. → GPU,distrib. framework, make net smaller



Industry Deployment

- Used in Facebook, Google, Microsoft
- Face recognition, image search, photo organization....
- Very fast at test time (~100 images/sec/GPU)



[Taigman et al. DeepFace: Closing the Gap to Human-Level Performance in Face Verification, CVPR'14]

Labeled Faces in Wild Dataset

• Task: given pair of images, same person or not?





Funneled













[Tagman et al. CVPR'14]

Detection with ConvNets

- So far, all about classification
- What about localizing objects within the scene?



Groundtruth: tv or monitor tv or monitor (2) tv or monitor (3) person remote control remote control (2)
Two General Approaches

1. Examine very position / scale

- E.g. Overfeat: Integrated recognition, localization and detection using convolutional networks, Sermanet et al., ICLR 2014
- 2. Use some kind of proposal mechanism to attend to a set of possible regions
 - E.g. Region-CNN [Rich feature hierarchies for accurate object detection and semantic segmentation, Girshick et al., CVPR 2014]

Sliding Window with ConvNet



Sliding Window with ConvNet



Input Window

Sliding Window with ConvNet



No need to compute two separate windows --- Just one big input window

Multi-Scale Sliding Window ConvNet











Multi-Scale Sliding Window ConvNet











OverFeat – Output before NMS



Overfeat Detection Results

[Sermanet et al. ICLR 2014]



Top predictions: trombone (confidence 26.8) oboe (confidence 17.5) oboe (confidence 11.5) ILSVRC2012_val_00000614.JPEG



Groundtruth:

person hat with a wide brim hat with a wide brim (2) hat with a wide brim (3) oboe oboe (2) saxophone trombone person (2) person (3) person (4)



Top predictions: watercraft (confidence 72.2) watercraft (confidence 2.1)



Top predictions: tennis ball (confidence 3.5) banana (confidence 2.4) banana (confidence 2.1) hotdog (confidence 2.0)

banana (confidence 1.9)

ILSVRC2012_val_00000320.JPEG



Top predictions: microwave (confidence 5.6) refrigerator (confidence 2.5) ILSVRC2012_val_00000519.JPEG



Groundtruth: bowl microwave

Groundtruth:

strawberrv strawberry (2) strawberry (3) strawberry (4) strawberry (5) strawberry (6) strawberry (7) strawberry (8) strawberry (9) strawberry (10) apple apple (2) apple (3)

Groundtruth: watercraft watercraft (2)



R-CNN Approach

[Rich feature hierarchies for accurate object detection and semantic segmentation, Girshick et al., CVPR 2014]

- Bottom-up proposa mechanism
- Scored by classifier
- Current best detection approach on PASCAL VOC



Figure 1: Object detection system overview. Our system (1) takes an input image, (2) extracts around 2000 bottom-up region proposals, (3) computes features for each proposal using a large convolutional neural network (CNN), and then (4) classifies each region using class-specific linear SVMs. R-CNN achieves a mean average precision (mAP) of **53.7% on PASCAL VOC 2010**. For comparison, [34] reports 35.1% mAP using the same region proposals, but with a spatial pyramid and bag-of-visual-words approach. The popular deformable part models perform at 33.4%.

- Further work combines proposal mechanism with classification network:
 - Fast R-CNN, Ross Girshick, arXiv 1504.08083, 2015.
 - Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks, Shaoqing Ren et al., arXiv 1506.01497, 2015

Video Classification

- Want to capture temporal structure
- 3D convolutions & 3D max-pooling
- E.g. C3D model



8 convolution, 5 pool, 2 fully-connected layers
3x3x3 convolution kernels
2x2x2 pooling kernels

[Learning Spatiotemporal Features with 3D Convolutional Networks, Tran et al., arXiv: 1412.0767, 2014] [Slide: Manohar Paluri]

Action Recognition – UCF101 dataset



[Slide: Manohar Paluri]

Action Recognition Results

	Method	Accuracy (%)
Baselines	Imagenet	68.8
	iDT	76.2
Use raw pixel inputs	Deep networks [19]	65.4
	Spatial stream network [36]	72.6
	LRCN [7]	71.1
	LSTM composite model [39]	75.8
	C3D (1 net)	82.3
	C3D (3 nets)	85.2
Use optical flows	iDT with Fisher vector [31]	87.9
	Temporal stream network [36]	83.7
	Two-stream networks [36]	88.0
	LRCN [7]	82.9
	LSTM composite model [39]	84.3
	Multi-skip feature stacking [26]	89.1
	C3D (3 nets) + iDT	90.4

[Slide: Manohar Paluri]

2D vs 3D Convnets

• UCF101 training





t-SNE visualization

[Slide: Manohar Paluri]

Sport Classification Results



.

1 ice_skating:0.98 2 speed_skating:0.01

Method	Number of Nets	Clip hit@1	Video hit@1	Video hit@5
Deep Video's Single-Frame + Multires [19]	3 nets	42.4	60.0	78.5
Deep Video's Slow Fusion [19]	1 net	41.9	60.9	80.2
C3D (trained from scratch)	1 net	44.9	60.0	84.4
C3D (fine-tuned from I380K pre-trained model)	1 net	46.1	61.1	85.2

Dense Scene Labeling

- Classification: pixels -> label
- Detection: pixels -> boxes
- Use Convnets to do pixels -> pixels
 - Segmentation of image
 - Image processing tasks (denoising etc.)
 - Don't want pooling

Dense Scene Labeling



• Convnet output is per-pixel label map



Convnet output is per-pixel depth map



• Convnet output is per-pixel normal map

Eigen et al. architecture

Input: 320x240



Architecture



Multi-Scale Convolutional Architecture, Eigen et al., arXiv 1411.4734, 2014]

Multi-Scale Convnets

Input: 320x240



Multi-Scale Convolutional Architecture, Eigen et al., arXiv 1411.4734, 2014]

Use Appropriate Loss Functions

Depth: $d = D - D^*$ D = log predicted depth, D* = log true depth $L_{depth}(D, D^*) = \frac{1}{n} \sum_{i} d_i^2 - \frac{1}{2n^2} \left(\sum_{i} d_i\right)^2 + \frac{1}{n} \sum_{i} [(\nabla_x d_i)^2 + (\nabla_y d_i)^2]$

[Predicting Depth, Surface Normals and Semantic Labels with a Common Multi-Scale Convolutional Architecture, Eigen et al., arXiv 1411.4734, 2014]

Depths Comparison



[Predicting Depth, Surface Normals and Semantic Labels with a Common Multi-Scale Convolutional Architecture, Eigen et al., arXiv 1411.4734, 2014]

Surface Normals



[Predicting Depth, Surface Normals and Semantic Labels with a Common Multi-Scale Convolutional Architecture, Eigen et al., arXiv 1411.4734, 2014]

Scene Parsing

• Farabet et al. "Learning hierarchical features for scene labeling" PAMI 2013



Segmentation

- Ciresan et al. "DNN segment neuronal membranes..." NIPS 2012
- Turaga et al. "Maximin learning of image segmentation" NIPS 2009



Denoising with ConvNets

• Burger et al. "Can plain NNs compete with BM3D?" CVPR 2012



Deblurring with Convnets

- Blind deconvolution
 - Learning to Deblur, Schuler et al., arXiv 1406.7444, 2014







Result of [Zho+13]

PSNR 23.17





Deblurring result w. noise *agnostic* training PSNR 23.29

Deblurring result w. noise *specific* training **PSNR 23.41**

Inpainting with Convnets

- Image Denoising and Inpainting with Deep Neural Networks, Xie et al. NIPS 2012.
- Mask-specific inpainting with deep neural networks, Köhler et al., Pattern Recognition 2014

nd Sirius form a nearly equilateral triangle. Thes Naos, in the Ship, and Phaet, in the Dove, form a hu known as the Eqyptian "X." From earliest times Siri been known as the Dog of Orion. It is 324 times brig the average sixth-magnitude star, and is the nearest earth of all the stars in this latitude, its distance be 8.7 light years. At this distance the Sun star a little brighter than the Pole Star. appe CANIS MAJOR] ARGO NAVIS (ŤrÅ 'go i ARGO. (Face South.) LOCATION.-Argo is Canis Major. If a line joining Betelgeuze a prolonged 18Ű southeast, if will point out the scouther is the state of the state the second magnitude in the ro in the southeast corner of i of a deep yellow o above it, two of w companion, which i a double for an opera M.). The star Markel stars near it. The Fr natians believed th that bore Osiris a contains two noted cts invisible in Canopus, the second est star. an e remari variable star Ε. [Illust ONOCER (mÅ□-nos´-e-ros)--T ith.) L(Monoceros is to be foun een Can Minor. Three of its naaniti t line northeast and 9A° ea Rete ze, and about the san th of A Gen ne reaion around the rich ewed with an opera-a he variable S, and a cluster about midw to stars about 7° apart in the tail of th field ϱ and nter stars to Procyon. These stars ar

)riginal



Removing Local Corruption

.

Restoring An Image Taken Through a Window Covered with Dirt or Rain

Rain Sequence

Each frame processed independently

David Eigen, Dilip Krishnan and Rob Fergus ICCV 2013

Removing Local Corruption

• Restoring An Image Taken Through a Window Covered with Dirt or Rain, Eigen et al., ICCV 2013.



Convnet + Structured Learning

 Gradient-based learning applied to document recognition, Yann LeCun, Leon Bottou, Yoshua Bengio and Patrick Haffner, Proc. IEEE, Nov 1998.





Convnet + Structured Learning

- Learning Deep Structured Models, Liang-Chieh Chen, Alexander G. Schwing, Alan L. Yuille, Raquel Urtasun, arXiv 1407.2538, 2014
- Joint Training of a Convolutional Network and a Graphical Model for Human Pose Estimation, J. Tompson, A. Jain, Y. LeCun, C. Bregler, NIPS 2014
- Lots more recently.....

BODY TRACKING

• Joint Training of a Convolutional Network and a Graphical Model for Human Pose Estimation

J. Tompson, A. Jain, Y. LeCun, C. Bregler, NIPS 2014





BODY TRACKING: PART DETECTOR

Simplified multi-resolution efficient model:





BODY TRACKING: SPATIAL MODEL

Start with MRF formulation

"Convolutional priors"

Sum-product belief propagation






BODY TRACKING: SPATIAL MODEL

Implement it as a network (no longer MRF)!

Speech

Brief Aside - Speech

- Also huge impact by neural nets
- Traditional approach (pre-2009):



• Very incremental gains in performance

Deep Learning Technical Revolution



[Slide: Tara Sainath, Google, Advancements in Deep Learning, SLT Keynote, Dec 2014.]

DNN Acoustic Modeling Results

- DNNs provide between a 8-25% relative improvement in word error rate over GMM/HMM systems across a variety of tasks and languages
- Results confirmed by many, many research labs

	300 hour SWB Conversational Telephony	400 hour Broadcast News	2000 hour Voice Search
GMM/HMM	14.3	16.5	16.0
DNN	12.2	15.2	12.2
% Relative Improvement	14.7	7.9	23.8

[Slide: Tara Sainath, Google, Advancements in Deep Learning, SLT Keynote, Dec 2014.]

CNN vs DNN Results

 CNNs trained with vtln-warped log-mel fb features offer between a 4-12% relative improvement over DNNs trained with speaker-adapted features (VTLN +fMLLR)

Model	BN-50	BN-400	SWB-300
Baseline GMM/HMM	18.1	13.8	14.5
DNN	15.8	13.3	12.2
CNN	15.0	12.0	11.5

[Sainath et al, ICASSP 2013]

[Slide: Tara Sainath, Google, Advancements in Deep Learning, SLT Keynote, Dec 2014.]

End-to-End Recognition

• Go directly from raw waveform



- Convolutional Neural Networks-based Continuous Speech Recognition using Raw Speech Signal, Palaz, Magimai-Doss, Collobert, ICASSP 2015.
- Superior results on TIMIT (phoneme recog), comparable results on WSJ.

Natural Language Processing

Language modeling

- Natural language is a sequence of sequences
- Some sentences are more likely than others:
 - "How are you ?" has a high probability
 - "How banana you ? " has a low probability

Neural Network Language Models



Bengio, Y., Schwenk, H., Sencal, J. S., Morin, F., & Gauvain, J. L. (2006). Neural probabilistic language models. In Innovations in Machine Learning (pp. 137-186). Springer Berlin Heidelberg.

[Slide: Antoine Border & Jason Weston, EMNLP Tutorial 2014]

Recurrent Neural Network Language Models

Key idea: *input to predict next word is current word plus context fed-back from previous word (i.e. remembers the past with recurrent connection).*



Figure: Recurrent neural network based LM

Recurrent neural network based language model. Mikolov et al., Interspeech, '10.

[Slide: Antoine Border & Jason Weston, EMNLP Tutorial 2014]

Recurrent neural networks - schema



Backpropagation through time

- The intuition is that we unfold the RNN in time
- We obtain deep neural network with shared weights U and W



Backpropagation through time

- We train the unfolded RNN using normal backpropagation + SGD
- In practice, we limit the number of unfolding steps to 5 – 10
- It is computationally more efficient to propagate gradients after few training examples (batch mode)

s(t) v U w(t-1)U w(t-2)s(t-1) y(t-1) s(t-2) W y(t-2) s(t-3) $\mathbf{y}(t-3)$ 100

w(t)

 $\mathbf{y}(t)$

Tomas Mikolov, COLING 2014

NNLMS vs. RNNS: Penn Treebank Results (Mikolov)

Model	Weight	PPL
3-gram with Good-Turing smoothing (GT3)	0	165.2
5-gram with Kneser-Ney smoothing (KN5)	0	141.2
5-gram with Kneser-Ney smoothing + cache	0.0792	125.7
Maximum entropy model	0	142.1
Random clusterings LM	0	170.1
Random forest LM	0.1057	131.9
Structured LM	0.0196	146.1
Within and across sentence boundary LM	0.0838	116.6
Log-bilinear LM	0	144.5
Feedforward NNLM	0	140.2
Syntactical NNLM	0.0828	131.3
Combination of static RNNLMs	0.3231	102.1
Combination of adaptive RNNLMs	0.3058	101.0
ALL	1	83.5

Recent uses of NNLMs and RNNs to improve machine translation: Fast and Robust NN Joint Models for Machine Translation, Devlin et al, ACL '14. Also Kalchbrenner '13, Sutskever et al., '14., Cho et al., '14.

[Slide: Antoine Border & Jason Weston, EMNLP Tutorial 2014]

Language modelling – RNN samples

the meaning of life is that only if an end would be of the whole supplier. widespread rules are regarded as the companies of refuses to deliver. in balance of the nation's information and loan growth associated with the carrier thrifts are in the process of slowing the seed and commercial paper.

More depth gives more power



LSTM - Long Short Term Memory

[Hochreiter and Schmidhuber, Neural Computation 1997]

- Ad-hoc way of modelling long dependencies
- Many alternative ways of modelling it
- Next hidden state is modification of previous hidden state (so information doesn't decay too fast).



For simple explanation, see [Recurrent Neural Network Regularization, Wojciech Zaremba, Ilya Sutskever, Oriol Vinyals, arXiv 1409.2329, 2014]

RNN-LSTMs for Machine Translation



Sequence to Sequence Learning with Neural Networks, Ilya Sutskever, Oriol Vinyals, Quoc Le, NIPS 2014

Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation, Kyunghyun Cho, Bart van Merrienboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, Yoshua Bengio, EMNLP 2014

Visualizing Internal Representation

t-SNE projection of network state at end of input sentence



Ilya Sutskever, Oriol Vinyals, Quoc Le, NIPS 2014

Translation - examples

- FR: Les avionneurs se querellent au sujet de la largeur des sièges alors que de grosses commandes sont en jeu
- Google Translate: Aircraft manufacturers are quarreling about the seat width as large orders are at stake
- LSTM: Aircraft manufacturers are concerned about the width of seats while large orders are at stake
- Ground Truth: Jet makers feud over seat width with big orders at stake

[Sequence to Sequence Learning with Neural Networks, Ilya Sutskever, Oriol Vinyals, Quoc Le, NIPS 2014]

Image Captioning: Vision + NLP

- Generate short text descriptions of image, given just picture.
- Use Convnet to extract image features
- RNN or LSTM model takes image features as input, generates text



Many recent works on this:

- Baidu/UCLA: Explain Images with Multimodal Recurrent Neural Networks
- Toronto: Unifying Visual-Semantic Embeddings with Multimodal Neural Language Models
- Berkeley: Long-term Recurrent Convolutional Networks for Visual Recognition and Description
- Google: Show and Tell: A Neural Image Caption Generator
- Stanford: Deep Visual-Semantic Alignments for Generating Image Description
- UML/UT: Translating Videos to Natural Language Using Deep Recurrent Neural Networks
- Microsoft/CMU: Learning a Recurrent Visual Representation for Image Caption Generation
- Microsoft: From Captions to Visual Concepts and Back

Image Captioning Examples



[men (0.59)] [group (0.66)] [woman (0.64)] [people (0.89)] [holding (0.60)] [playing (0.61)] [tennis (0.69)] [court (0.51)] [standing (0.59)] [skis (0.58)] [street (0.52)] [man (0.77)] [skateboard (0.67)]

a group of people standing next to each other people stand outside a large ad for gap featuring a young boy



[person (0.55)] [street (0.53)] [holding (0.55)] [group (0.63)] [slope (0.51)] [standing (0.62)] [snow (0.91)] [skis (0.74)] [player (0.54)] [people (0.85)] [men (0.57)] [skiing (0.51)] [skateboard (0.89)] [riding (0.75)] [tennis (0.74)] [trick (0.53)] [skate (0.52)] [woman (0.52)] [man (0.86)] [down (0.61)]

a group of people riding skis down a snow covered slope a guy on a skate board on the side of a ramp



a courtyard full of poles pigeons and garbage cans also has benches on either side of it one of which shows the back of a large person facin g in the direction of the pigeons



[brown (0.68)] [baby (0.62)] [walking (0.57)] [laying (0.61)] [man (0.57)] [standing (0.79)] [field (0.65)] [water (0.83)] [large (0.71)] [dirt (0.65)] [river (0.58)] a baby elephant standing next to each other on a field elephants are playing together in a shallow watering hole

From Captions to Visual Concepts and Back, Hao Fang* Saurabh Gupta* Forrest Iandola* Rupesh K. Srivastava*, Li Deng Piotr Dollar, Jianfeng Gao Xiaodong He, Margaret Mitchell John C. Platt, C. Lawrence Zitnick, Geoffrey Zweig, CVPR 2015.

Facebook AI Research

- ~50 people working in ML/vision/NLP/speech/AI
 -1/3 are research engineers (some of FB's best coders)
 Yann LeCun is lab director
- Freedom to publish & open-source code
- Easy to productize (1.1B users)
- Labs in:
 - Menlo Park, California (Facebook HQ)
 - New York City
 - Paris
- We are hiring!

- [Slide 5]
- P. Felzenszwalb, R. Girshick, D. McAllester, D. Ramanan, Object Detection with Discriminatively Trained Part Based Models, IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol. 32, No. 9, September 2010
- Zheng Song^{*}, Qiang Chen^{*}, Zhongyang Huang, Yang Hua, and Shuicheng Yan. Con-tex-tual-iz-ing Ob-ject De-tec-tion and Clas-si-fi-ca-tion. In CVPR'11. (* in-di-cates equal contri-bu-tion) [No. 1 per-for-mance in VOC'10 clas-si-fi-ca-tion task]
- [Slide 6]
- Finding the Weakest Link in Person Detectors, D. Parikh, and C. L. Zitnick, CVPR, 2011.
- [Slide 7]
- Gehler and Nowozin, On Feature Combination for Multiclass Object Classification, ICCV'09
- [Slide 8]
- <u>http://www.amazon.com/Vision-David-Marr/dp/0716712849</u>
- [Slide 10]
- Yoshua Bengio and Yann LeCun: Scaling learning algorithms towards AI, in Bottou, L. and Chapelle, O. and DeCoste, D. and Weston, J. (Eds), Large-Scale Kernel Machines, MIT Press, 2007

- [Slide 11]
- S. Lazebnik, C. Schmid, and J. Ponce, Beyond Bags of Features: Spatial Pyramid Matching for Recognizing Natural Scene Categories, CVPR 2006
- [Slide 12]
- Christoph H. Lampert, Hannes Nickisch, Stefan Harmeling: "Learning To Detect Unseen Object Classes by Between-Class Attribute Transfer", IEEE Computer Vision and Pattern Recognition (CVPR), Miami, FL, 2009
- [Slide 14] Riesenhuber, M. & Poggio, T. (1999). Hierarchical Models of Object Recognition in Cortex. Nature Neuroscience 2: 1019-1025.
- <u>http://www.scholarpedia.org/article/Neocognitron</u>
- K. Fukushima: "Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position", Biological Cybernetics, 36[4], pp. 193-202 (April 1980).
- Y. LeCun, L. Bottou, Y. Bengio and P. Haffner: Gradient-Based Learning Applied to Document Recognition, Proceedings of the IEEE, 86(11):2278-2324, November 1998

- [Slide 30]
- Y-Lan Boureau, Jean Ponce, and Yann LeCun, A theoretical analysis of feature pooling in vision algorithms, Proc. International Conference on Machine learning (ICML'10), 2010
- [Slide 31]
- Q.V. Le, J. Ngiam, Z. Chen, D. Chia, P. Koh, A.Y. Ng, Tiled Convolutional Neural Networks. NIPS, 2010
- <u>http://ai.stanford.edu/~quocle/TCNNweb</u>
- Matthew D. Zeiler, Graham W. Taylor, and Rob Fergus, Adaptive Deconvolutional Networks for Mid and High Level Feature Learning, International Conference on Computer Vision(November 6-13, 2011)
- [Slide 32]
- Yuanhao Chen, Long Zhu, Chenxi Lin, Alan Yuille, Hongjiang Zhang. Rapid Inference on a Novel AND/OR graph for Object Detection, Segmentation and Parsing. NIPS 2007.

- [Slide 35]
- P. Smolensky, Parallel Distributed Processing: Volume 1: Foundations, D. E. Rumelhart, J. L. McClelland, Eds. (MIT Press, Cambridge, 1986), pp. 194–281.
- G. E. Hinton, Neural Comput. 14, 1711 (2002).
- [Slide 36]
- M. Ranzato, Y. Boureau, Y. LeCun. "Sparse Feature Learning for Deep Belief Networks". Advances in Neural Information Processing Systems 20 (NIPS 2007).
- [Slide 39]
- Hinton, G. E. and Salakhutdinov, R. R., Reducing the dimensionality of data with neural networks. Science, Vol. 313. no. 5786, pp. 504 507, 28 July 2006.
- [Slide 41]
- A. Torralba, K. P. Murphy and W. T. Freeman, Contextual Models for Object Detection using Boosted Random Fields, Adv. in Neural Information Processing Systems 17 (NIPS), pp. 1401-1408, 2005.

- [Slide 42]
- Ruslan Salakhutdinov and Geoffrey Hinton, Deep Boltzmann Machines, 12th International Conference on Artificial Intelligence and Statistics (2009).
- [Slide 44]
- P. Felzenszwalb, R. Girshick, D. McAllester, D. Ramanan, Object Detection with Discriminatively Trained Part Based Models, IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol. 32, No. 9, September 2010
- Long Zhu, Yuanhao Chen, Alan Yuille, William Freeman. Latent Hierarchical Structural Learning for Object Detection. CVPR 2010.
- [Slide 45]
- Matthew D. Zeiler, Graham W. Taylor, and Rob Fergus, Adaptive Deconvolutional Networks for Mid and High Level Feature Learning, International Conference on Computer Vision(November 6-13, 2011)

- [Slide 48]
- S.C. Zhu and D. Mumford, A Stochastic Grammar of Images, Foundations and Trends in Computer Graphics and Vision, Vol.2, No.4, pp 259-362, 2006.
- [Slide 49]
- R. Girshick, P. Felzenszwalb, D. McAllester, Object Detection with Grammar Models, NIPS 2011
- [Slide 50]
- P. Felzenszwalb, D. Huttenlocher, Pictorial Structures for Object Recognition, International Journal of Computer Vision, Vol. 61, No. 1, January 2005
- M. Fischler and R. Elschlager. The Representation and Matching of Pictoral Structures. (1973)
- [Slide 51]
- S. Fidler, M. Boben, A. Leonardis. A coarse-to-fine Taxonomy of Constellations for Fast Multi-class Object Detection. ECCV 2010.
- S. Fidler and A. Leonardis. Towards Scalable Representations of Object Categories: Learning a Hierarchy of Parts. CVPR 2007.

- [Slide 52]
- Long Zhu, Chenxi Lin, Haoda Huang, Yuanhao Chen, Alan Yuille. Unsupervised Structure Learning: Hierarchical Recursive Composition, Suspicious Coincidence and Competitive Exclusion. ECCV 2008.
- [Slide 53]
- Hinton, G. E., Krizhevsky, A. and Wang, S, Transforming Auto-encoders. ICANN-11: International Conference on Artificial Neural Networks, 2011
- Matthew D. Zeiler, Graham W. Taylor, and Rob Fergus, Adaptive Deconvolutional Networks for Mid and High Level Feature Learning, International Conference on Computer Vision(November 6-13, 2011)
- [Slide 54]
- Q.V. Le, M.A. Ranzato, R. Monga, M. Devin, K. Chen, G.S. Corrado, J. Dean, A.Y. Ng., Building high-level features using large scale unsupervised learning. ICML, 2012.
- [Slide 55]
- Ruslan Salakhutdinov and Geoffrey Hinton, Deep Boltzmann Machines, 12th International Conference on Artificial Intelligence and Statistics (2009).

- [Slide 56]
- http://www.image-net.org/challenges/LSVRC/2010/
- Q.V. Le, M.A. Ranzato, R. Monga, M. Devin, K. Chen, G.S. Corrado, J. Dean, A.Y. Ng., Building high-level features using large scale unsupervised learning. ICML, 2012.
- Q.V. Le, W.Y. Zou, S.Y. Yeung, A.Y. Ng., Learning hierarchical spatio-temporal features for action recognition with independent subspace analysis, CVPR 2011