



10-301/10-601 Introduction to Machine Learning

Machine Learning Department School of Computer Science Carnegie Mellon University

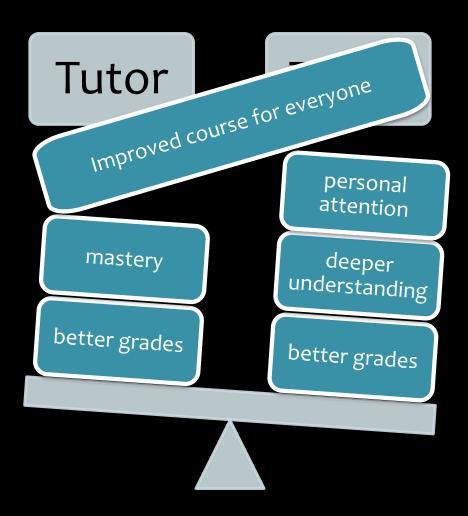
Deep Learning: RNNs & CNNs

Matt Gormley Lecture 14 Mar. 2, 2023

Reminders

- Exit Poll: Exam 1
- Homework 5: Neural Networks
 - Out: Sun, Feb 26
 - Due: Fri, Mar 17 at 11:59pm

Peer Tutoring



Backpropagation and Deep Learning

Convolutional neural networks (CNNs) and recurrent neural networks (RNNs) are simply fancy computation graphs (aka. hypotheses or decision functions).

Our recipe also applies to these models and (again) relies on the **backpropagation algorithm** to compute the necessary gradients.

BACKGROUND: HUMAN LANGUAGE TECHNOLOGIES

Human Language Technologies



Machine Translation

기계 번역은 특히 영어와 한국어와 같은 언어 쌍의 경우 매우 어렵습니다.

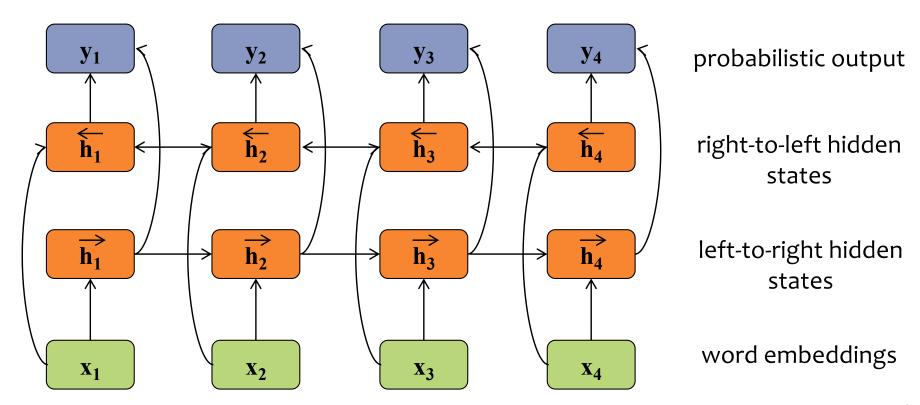
Summarization

```
Lorem ipsum dolor sit amet,
cor

lab Lorem ipsum dolor sit amet,
nith eiu Lorem ipsum dolor sit amet,
viol lab Lorem ipsum dolor sit amet,
lor
```

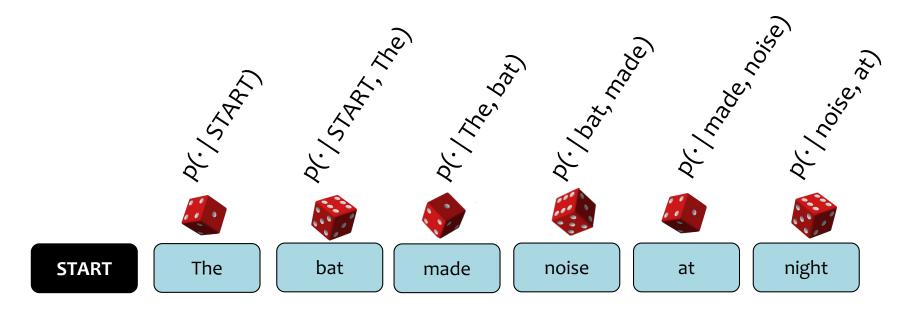
Bidirectional RNN

RNNs are a now commonplace backbone in deep learning approaches to natural language processing

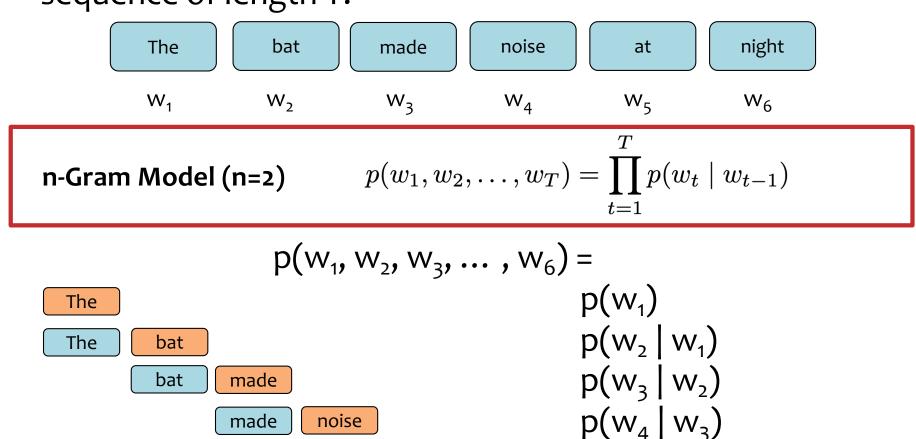


BACKGROUND: N-GRAM LANGUAGE MODELS

- Goal: Generate realistic looking sentences in a human language
- Key Idea: condition on the last n-1 words to sample the nth word



<u>Question</u>: How can we **define** a probability distribution over a sequence of length T?



at

at

night

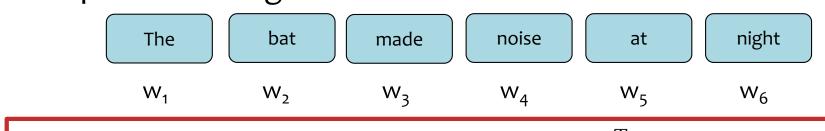
noise

 $p(w_5 | w_4)$

 $p(w_6 | w_5)$

<u>Question</u>: How can we **define** a probability distribution over a sequence of length T?

night



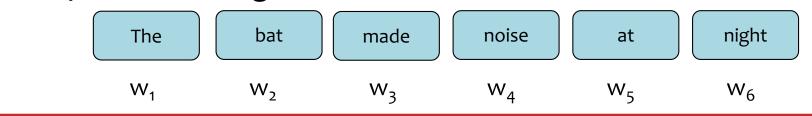
$$p(w_1, w_2, \dots, w_T) = \prod_{t=1}^{T} p(w_t \mid w_{t-1}, w_{t-2})$$

$$p(w_{1}, w_{2}, w_{3}, ..., w_{6}) = \\ p(w_{1}) \\ p(w_{1}) \\ p(w_{2} \mid w_{1}) \\ p(w_{3} \mid w_{2}, w_{2}, w_{1}) \\ p(w_{3} \mid w_{2}, w_{2}, w_{1}) \\ p(w_{3} \mid w_{2}, w_{2}, w_{2}, w_{2}, w_{2}) \\ p(w_{3} \mid w_{2}, w_{2}, w_{2}, w_{2}, w_{2}, w_{2}, w_{2}) \\ p(w_{3} \mid w_{2}, w_$$

$$p(w_4 | w_3, w_2)$$

 $p(w_5 | w_4, w_3)$
 $p(w_6 | w_5, w_4)$

Question: How can we define a probability distribution over a sequence of length T?



n-Gram Model (n=3)
$$p(w_1, w_2, \dots, w_T) = \prod_{t=1}^{T} p(w_t \mid w_{t-1}, w_{t-2})$$

$$p(w_1, w_3, ..., w_6) = p(w_1)$$

The

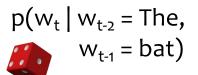
The

The

Note: This is called a **model** because we made some **assumptions** about how many previous words to condition on (i.e. only n-1 words)

Learning an n-Gram Model

Question: How do we learn the probabilities for the n-Gram Model?



vv t	P(' ',')
ate	0.015
•••	
flies	0.046
•••	
zebra	0.000

$$p(w_t | w_{t-2} = made, w_{t-1} = noise)$$

W _t	p(· ·,·)
at	0.020
•••	
pollution	0.030
•••	
zebra	0.000

$$p(w_t | w_{t-2} = cows, w_{t-1} = eat)$$

w _t	p(· ·,·)
corn	0.420

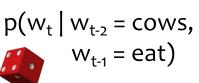
•••		
grass	0.510	
•••		
zebra	0.000	

Learning an n-Gram Model

<u>Question</u>: How do we **learn** the probabilities for the n-Gram Model?

Answer: From data! Just count n-gram frequencies

```
... the cows eat grass...
... our cows eat hay daily...
... factory-farm cows eat corn...
... on an organic farm, cows eat hay and...
... do your cows eat grass or corn?...
... what do cows eat if they have...
... cows eat corn when there is no...
... which cows eat which foods depends...
... if cows eat grass...
... when cows eat corn their stomachs...
... should we let cows eat corn?...
```

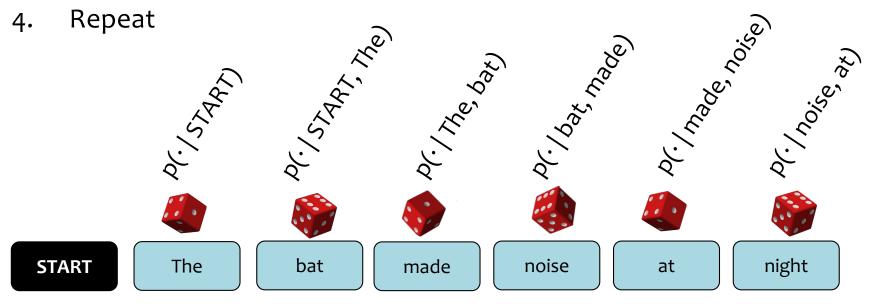


W _t	p(· ·,·)
corn	4/11
grass	3/11
hay	2/11
if	1/11
which	1/11

Sampling from a Language Model

<u>Question</u>: How do we sample from a Language Model? <u>Answer</u>:

- 1. Treat each probability distribution like a (50k-sided) weighted die
- 2. Pick the die corresponding to $p(w_t | w_{t-2}, w_{t-1})$
- 3. Roll that die and generate whichever word w_t lands face up



Sampling from a Language Model

Question: How do we sample from a Language Model?

Answer:

- 1. Treat each probability distribution like a (50k-sided) weighted die
- 2. Pick the die corresponding to $p(w_t | w_{t-2}, w_{t-1})$
- 3. Roll that die and generate whichever word w_t lands face up
- 4. Repeat

Training Data (Shakespeaere)

I tell you, friends, most charitable care ave the patricians of you. For your wants, Your suffering in this dearth, you may as well Strike at the heaven with your staves as lift them Against the Roman state, whose course will on The way it takes, cracking ten thousand curbs Of more strong link asunder than can ever Appear in your impediment. For the dearth, The gods, not the patricians, make it, and Your knees to them, not arms, must help.

5-Gram Model

Approacheth, denay. dungy
Thither! Julius think: grant,—0
Yead linens, sheep's Ancient,
Agreed: Petrarch plaguy Resolved
pear! observingly honourest
adulteries wherever scabbard
guess; affirmation—his monsieur;
died. jealousy, chequins me.
Daphne building. weakness: sun—
rise, cannot stays carry't,
unpurposed. prophet—like drink;
back—return 'gainst surmise
Bridget ships? wane; interim?
She's striving wet;

RECURRENT NEURAL NETWORK (RNN) LANGUAGE MODELS

Recurrent Neural Networks (RNNs)

inputs: $\mathbf{x} = (x_1, x_2, \dots, x_T), x_i \in \mathcal{R}^I$

hidden units: $\mathbf{h} = (h_1, h_2, \dots, h_T), h_i \in \mathcal{R}^J$

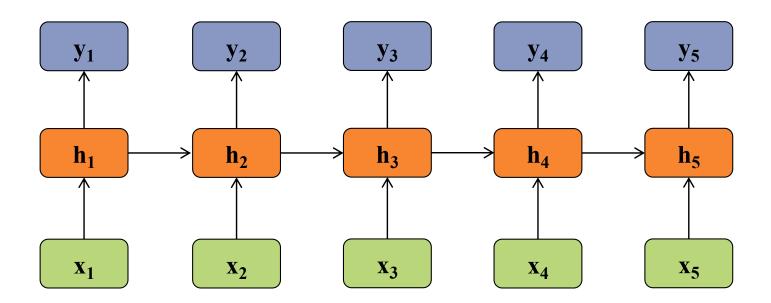
outputs: $\mathbf{y} = (y_1, y_2, \dots, y_T), y_i \in \mathcal{R}^K$

nonlinearity: \mathcal{H}

Definition of the RNN:

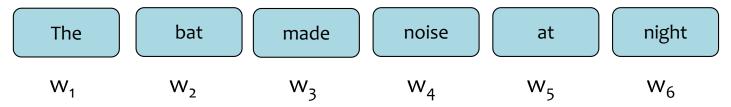
$$h_t = \mathcal{H}\left(W_{xh}x_t + W_{hh}h_{t-1} + b_h\right)$$

$$y_t = W_{hy}h_t + b_y$$



The Chain Rule of Probability

<u>Question</u>: How can we **define** a probability distribution over a sequence of length T?



Chain rule of probability: $p(w_1, w_2, \ldots, w_T) = \prod_{t=1}^{T} p(w_t \mid w_{t-1}, \ldots, w_1)$

$$p(w_1, w_3, \dots, w_6) = p(w_1)$$
The Note: This is called the chain **rule** because it is **always** true for every probability distribution

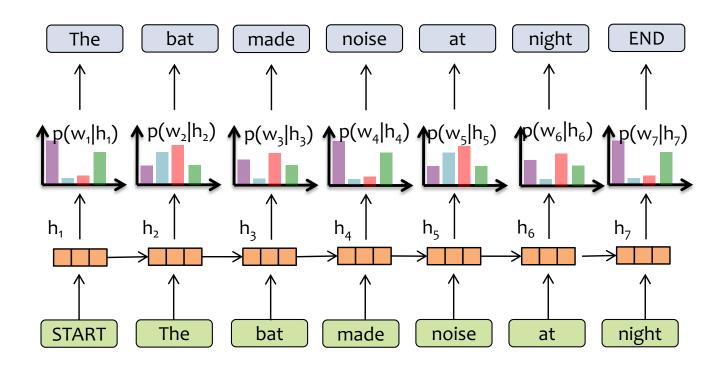
The Other of the Note: The Note: This is called the chain v_1 and v_2 are the Note: The Note:

Recall...

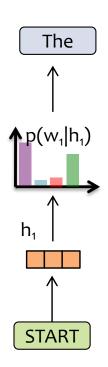
RNN Language Model:
$$p(w_1, w_2, \ldots, w_T) = \prod_{t=1}^T p(w_t \mid f_{\boldsymbol{\theta}}(w_{t-1}, \ldots, w_1))$$

$$p(w_{1},w_{2},w_{3},\ldots,w_{6}) = \\ p(w_{1}) \\ The & bat \\ made \\ p(w_{2} \mid f_{\theta}(w_{1})) \\ p(w_{3} \mid f_{\theta}(w_{2},w_{1})) \\ The & bat \\ made & noise \\ at \\ p(w_{4} \mid f_{\theta}(w_{3},w_{2},w_{1})) \\ The & bat \\ made & noise \\ at \\ night \\ p(w_{6} \mid f_{\theta}(w_{5},w_{4},w_{3},w_{2},w_{1})) \\ The & bat \\ made & noise \\ at \\ night \\ p(w_{6} \mid f_{\theta}(w_{5},w_{4},w_{3},w_{2},w_{1})) \\ The & bat \\ made & noise \\ at \\ night \\ p(w_{6} \mid f_{\theta}(w_{5},w_{4},w_{3},w_{2},w_{1})) \\ The & bat \\ made & noise \\ at \\ night \\ p(w_{6} \mid f_{\theta}(w_{5},w_{4},w_{3},w_{2},w_{1})) \\ The & bat \\ made & noise \\ at \\ night \\ p(w_{6} \mid f_{\theta}(w_{5},w_{4},w_{3},w_{2},w_{1})) \\ The & bat \\ made & noise \\ at \\ night \\ p(w_{6} \mid f_{\theta}(w_{5},w_{4},w_{3},w_{2},w_{1})) \\ The & bat \\ made & noise \\ at \\ night \\ p(w_{6} \mid f_{\theta}(w_{5},w_{4},w_{3},w_{2},w_{1})) \\ The & bat \\ made & noise \\ at \\ night \\ p(w_{6} \mid f_{\theta}(w_{5},w_{4},w_{3},w_{2},w_{1})) \\ The & bat \\ made & noise \\ at \\ night \\ p(w_{6} \mid f_{\theta}(w_{5},w_{4},w_{3},w_{2},w_{1})) \\ The & bat \\ made & noise \\ at \\ night \\ p(w_{6} \mid f_{\theta}(w_{5},w_{4},w_{3},w_{2},w_{1})) \\ The & bat \\ made & noise \\ at \\ night \\ p(w_{6} \mid f_{\theta}(w_{5},w_{4},w_{3},w_{2},w_{1})) \\ The & bat \\ made & noise \\ at \\ night \\ p(w_{6} \mid f_{\theta}(w_{5},w_{4},w_{3},w_{2},w_{1})) \\ The & bat \\ made & noise \\ at \\ night \\ p(w_{6} \mid f_{\theta}(w_{5},w_{4},w_{3},w_{2},w_{1})) \\ The & bat \\ made & noise \\ at \\ night \\ p(w_{6} \mid f_{\theta}(w_{5},w_{4},w_{3},w_{2},w_{1})) \\ The & bat \\ night \\ p(w_{6} \mid f_{\theta}(w_{5},w_{4},w_{3},w_{2},w_{1})) \\ The & bat \\ night \\ p(w_{6} \mid f_{\theta}(w_{5},w_{4},w_{3},w_{2},w_{1})) \\ The & bat \\ night \\ p(w_{6} \mid f_{\theta}(w_{5},w_{4},w_{3},w_{2},w_{1})) \\ The & bat \\ night \\ p(w_{6} \mid f_{\theta}(w_{5},w_{4},w_{3},w_{2},w_{1})) \\ The & bat \\ night \\ p(w_{6} \mid f_{\theta}(w_{5},w_{4},w_{3},w_{2},w_{1})) \\ The & bat \\ night \\ p(w_{6} \mid f_{\theta}(w_{5},w_{4},w_{3},w_{2},w_{1})) \\ The & bat \\ night \\ p(w_{6} \mid f_{\theta}(w_{5},w_{4},w_{3},w_{2},w_{1})) \\ The & bat \\ night \\ p(w_{6} \mid f_{\theta}(w_{5},w_{4},w_{3},w_{2},w_{1})) \\ The & bat \\ night \\ p(w_{6} \mid f_{$$

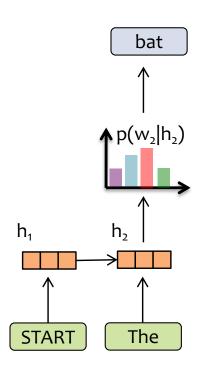
- (1) convert all previous words to a fixed length vector
- (2) define distribution $p(w_t | f_{\theta}(w_{t-1}, ..., w_1))$ that conditions on the vector



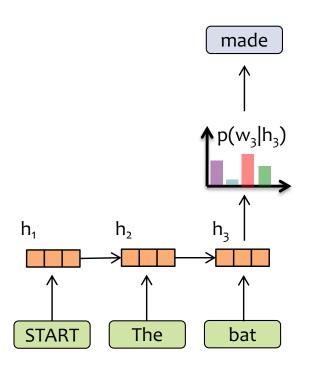
- (1) convert all previous words to a fixed length vector
- (2) define distribution $p(w_t | f_{\theta}(w_{t-1}, ..., w_1))$ that conditions on the vector $\mathbf{h}_t = f_{\theta}(w_{t-1}, ..., w_1)$



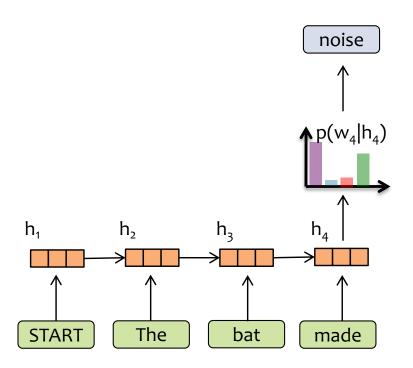
- (1) convert all previous words to a fixed length vector
- (2) define distribution $p(w_t | f_{\theta}(w_{t-1}, ..., w_1))$ that conditions on the vector $\mathbf{h}_t = f_{\theta}(w_{t-1}, ..., w_1)$



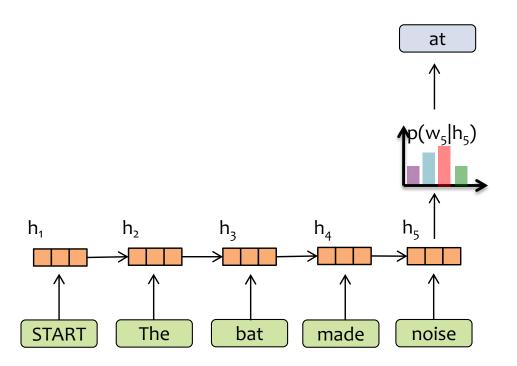
- (1) convert all previous words to a fixed length vector
- (2) define distribution $p(w_t | f_{\theta}(w_{t-1}, ..., w_1))$ that conditions on the vector $\mathbf{h}_t = f_{\theta}(w_{t-1}, ..., w_1)$



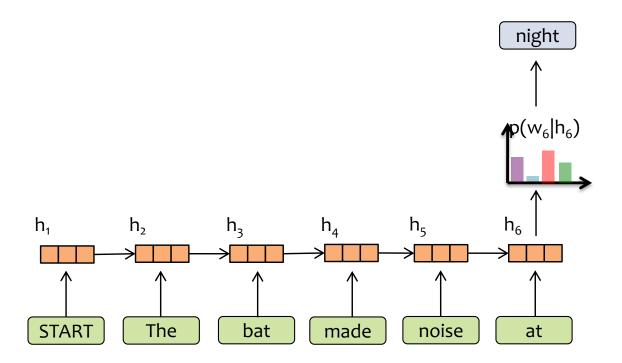
- (1) convert all previous words to a fixed length vector
- (2) define distribution $p(w_t | f_{\theta}(w_{t-1}, ..., w_1))$ that conditions on the vector $\mathbf{h}_t = f_{\theta}(w_{t-1}, ..., w_1)$



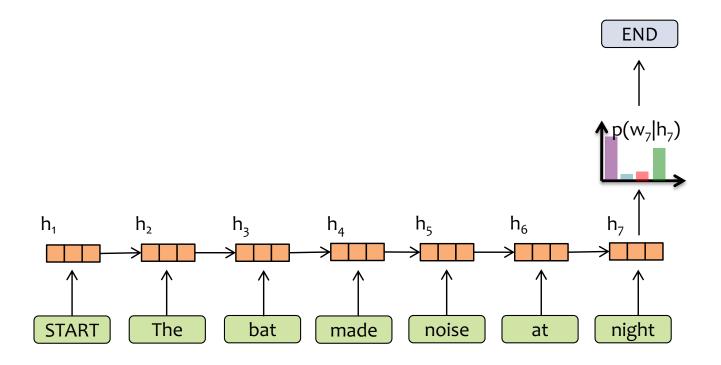
- (1) convert all previous words to a fixed length vector
- (2) define distribution $p(w_t | f_{\theta}(w_{t-1}, ..., w_1))$ that conditions on the vector $\mathbf{h}_t = f_{\theta}(w_{t-1}, ..., w_1)$



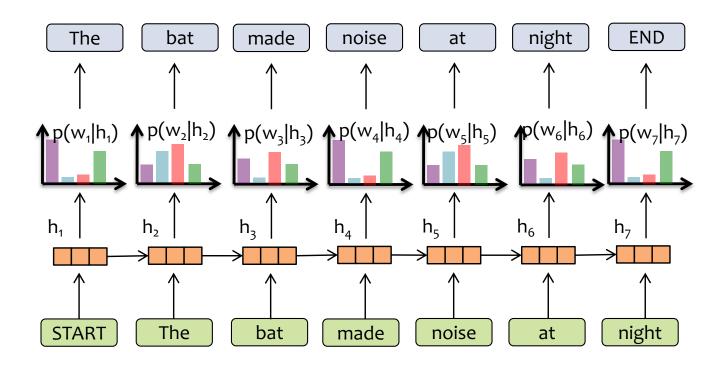
- (1) convert all previous words to a fixed length vector
- (2) define distribution $p(w_t | f_{\theta}(w_{t-1}, ..., w_1))$ that conditions on the vector $\mathbf{h}_t = f_{\theta}(w_{t-1}, ..., w_1)$



- (1) convert all previous words to a fixed length vector
- (2) define distribution $p(w_t | f_{\theta}(w_{t-1}, ..., w_1))$ that conditions on the vector $\mathbf{h}_t = f_{\theta}(w_{t-1}, ..., w_1)$



- (1) convert all previous words to a fixed length vector
- (2) define distribution $p(w_t | f_{\theta}(w_{t-1}, ..., w_1))$ that conditions on the vector $\mathbf{h}_t = f_{\theta}(w_{t-1}, ..., w_1)$



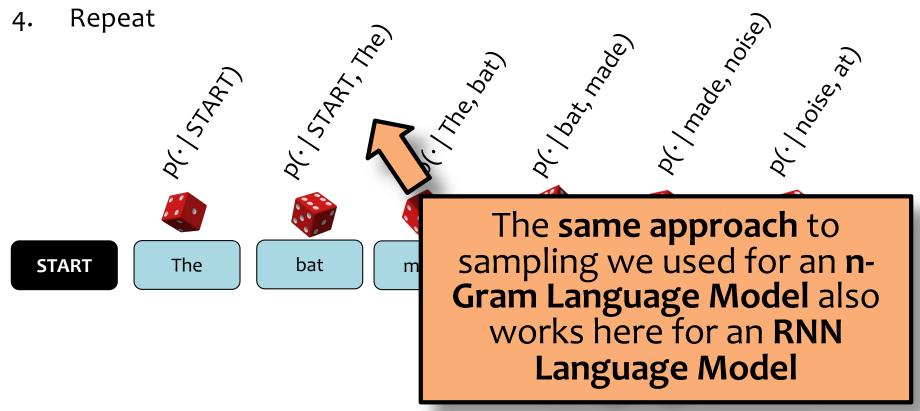
$$p(w_1, w_2, w_3, ..., w_T) = p(w_1 | h_1) p(w_2 | h_2) ... p(w_2 | h_T)$$

Sampling from a Language Model

Question: How do we sample from a Language Model?

Answer:

- 1. Treat each probability distribution like a (50k-sided) weighted die
- 2. Pick the die corresponding to $p(w_t | w_{t-2}, w_{t-1})$
- 3. Roll that die and generate whichever word w_t lands face up



??

VIOLA: Why, Salisbury must find his flesh and thought That which I am not aps, not a man and in fire, To show the reining of the raven and the wars To grace my hand reproach within, and not a fair are hand, That Caesar and my goodly father's world; When I was heaven of presence and our fleets, We spare with hours, but cut thy council I am great, Murdered a master's ready there My powe so much as hell: Some service i bondman here, Would show hi

KING LEAR: O, if you we feeble sight, the courtesy of your law, Your sight and several breath, will wear the gods With his heads, and my hands are wonder'd at the deeds, So drop upon your lordship's head, and your opinion Shall be against your honour.

??

CHARLES: Marry, do I, sir; and I came to acquaint you with a matter. I am given, sir, secretly to understand that your younger brother Orlando hath a disposition to come in disguised against me to try a fall. To-morrow, sir, I wrestle for my credit; and he that escapes proken limb shall acquit him is but young and tender; and, uld be loath to foil him, as I honour, if he come in: by love to you, I came hither to acquaint you with that either you might

to acquaint you with that either you might stay him from his intended or brook such disgrace well as he shared into, in that it is a thing of his own search and altogether against my will.

Shakespeare's As You Like It

VIOLA: Why, Salisbury must find his flesh and thought That which I am not aps, not a man and in fire, To show the reining of the raven and the wars To grace my hand reproach within, and not a fair are hand, That Caesar and my goodly father's world; When I was heaven of presence and our fleets, We spare with hours, but cut thy council I am great, Murdered and by thy master's ready there My power to give thee but so much as hell: Some service in the noble bondman here, Would show him to her wine.

KING LEAR: O, if you were a feeble sight, the courtesy of your law, Your sight and several breath, will wear the gods With his heads, and my hands are wonder'd at the deeds, So drop upon your lordship's head, and your opinion Shall be against your honour.

RNN-LM Sample

CHARLES: Marry, do I, sir; and I came to acquaint you with a matter. I am given, sir, secretly to understand that your younger brother Orlando hath a disposition to come in disguised against me to try a fall. To-morrow, sir, I wrestle for my credit; and he that escapes me without some broken limb shall acquit him well. Your brother is but young and tender; and, for your love, I would be loath to foil him, as I must, for my own honour, if he come in: therefore, out of my love to you, I came hither to acquaint you withal, that either you might stay him from his intendment or brook such disgrace well as he shall run into, in that it is a thing of his own search and altogether against my will.

RNN-LM Sample

VIOLA: Why, Salisbury must find his flesh and thought That which I am not aps, not a man and in fire, To show the reining of the raven and the wars To grace my hand reproach within, and not a fair are hand, That Caesar and my goodly father's world; When I was heaven of presence and our fleets, We spare with hours, but cut thy council I am great, Murdered and by thy master's ready there My power to give thee but so much as hell: Some service in the noble bondman here, Would show him to her wine.

KING LEAR: O, if you were a feeble sight, the courtesy of your law, Your sight and several breath, will wear the gods With his heads, and my hands are wonder'd at the deeds, So drop upon your lordship's head, and your opinion Shall be against your honour.

Shakespeare's As You Like It

CHARLES: Marry, do I, sir; and I came to acquaint you with a matter. I am given, sir, secretly to understand that your younger brother Orlando hath a disposition to come in disguised against me to try a fall. To-morrow, sir, I wrestle for my credit; and he that escapes me without some broken limb shall acquit him well. Your brother is but young and tender; and, for your love, I would be loath to foil him, as I must, for my own honour, if he come in: therefore, out of my love to you, I came hither to acquaint you withal, that either you might stay him from his intendment or brook such disgrace well as he shall run into, in that it is a thing of his own search and altogether against my will.

??

VIOLA: Why, Salisbury must find his flesh and thought That which I am not aps, not a man and in fire, To show the reining of the raven and the wars To grace my hand reproach within, and not a fair are hand, That Caesar and my goodly father's world; When I was heaven of presence and our fleets, We spare with hours, but cut thy council I am great, Murdered a master's ready there My powe so much as hell: Some service i

bondman here, Would show hi

KING LEAR: O, if you we feeble sight, the courtesy of your law, Your sight and several breath, will wear the gods With his heads, and my hands are wonder'd at the deeds, So drop upon your lordship's head, and your opinion Shall be against your honour.

??

CHARLES: Marry, do I, sir; and I came to acquaint you with a matter. I am given, sir, secretly to understand that your younger brother Orlando hath a disposition to come in disguised against me to try a fall. To-morrow, sir, I wrestle for my credit; and he that escapes but cut thy me without some broken limb shall acquit him is but young and tender; and, uld be loath to foil him, as I honour, if he come in: by love to you, I came hither to acquaint you with a matter. I am given, sir, secretly to understand that your younger brother Orlando hath a disposition to come in disguised against me to try a fall. To-morrow, sir, I wrestle for my credit; and he that escapes broken limb shall acquit him is but young and tender; and, uld be loath to foil him, as I honour, if he come in: by love to you, I came hither to acquaint you with a matter. I am given, sir, secretly to understand that your younger brother Orlando hath a disposition to come in disguised against me to try a fall. To-morrow, sir, I wrestle for my credit; and he that escapes broken limb shall acquit him is but young and tender; and, uld be loath to foil him, as I honour, if he come in: by love to you, I came hither to acquaint you with a matter. I am given, sir, secretly to understand that your younger broken limb shall acquit him is but young and tender; and, uld be loath to foil him, as I honour, if he come in:

to acquaint you will that either you might stay him from his intent or brook such disgrace well as he shared into, in that it is a thing of his own search and altogether against my will.

SEQUENCE TO SEQUENCE MODELS

Sequence to Sequence Model



Machine Translation

기계 번역은 특히 영어와 한국어와 같은 언어 쌍의 경우 매우 어렵습니다.

Summarization

```
Lorem ipsum dolor sit amet,
cor

lab Lorem ipsum dolor sit amet,
nith eiu Lorem ipsum dolor sit amet,
viol lab Lorem ipsum dolor sit amet,
lor
```

Sequence to Sequence Model

Now suppose you want generate a sequence conditioned on another input

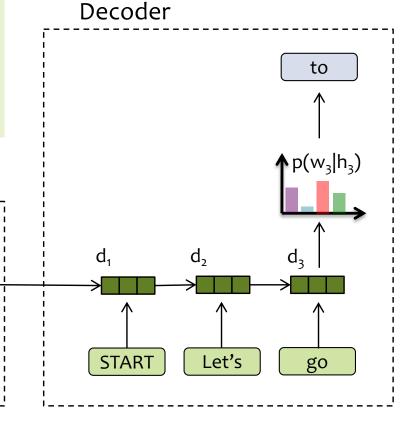
Key Idea:

Encoder

- Use an encoder model to generate a vector representation of the input
- 2. Feed the output of the encoder to a **decoder** which will generate the **output**

Applications:

- translation:
 Spanish → English
- summarization: article → summary
- speech recognition: speech signal → transcription



BACKGROUND: COMPUTER VISION

Example: Image Classification

- ImageNet LSVRC-2011 contest:
 - Dataset: 1.2 million labeled images, 1000 classes
 - Task: Given a new image, label it with the correct class
 - Multiclass classification problem
- Examples from http://image-net.org/

Not logged in. Login I Signup

Bird

IM. GENET

Warm-blooded egg-laying vertebrates characterized by feathers and forelimbs modified as wings

2126 pictures 92.85% Popularity Percentile



marine animal, marine creature, sea animal, sea creature (1)		1.1.1	120
scavenger (1)	Treemap Visualization	Images of the Synset	Downloads
- biped (0)			The second second
predator, predatory animal (1)		Maria A	F
- larva (49)			
- acrodont (0)			
- feeder (0)	No.		3
- stunt (0)			
r- chordate (3087)			
tunicate, urochordate, urochord (6)			
rephalochordate (1)			
vertebrate, craniate (3077)	725,704	X	
mammal, mammalian (1169)			
bird (871)	A STATE OF THE STA		一个
- dickeybird, dickey-bird, dickybird, dicky-bird (0)			
r cock (1)			
- hen (0)			
- nester (0)			
i- night bird (1)		355	
- bird of passage (0)	To State of the second	453	C (4) 3
- protoavis (0)			
- archaeopteryx, archeopteryx, Archaeopteryx lithographi			
- Sinornis (0)			
- Ibero-mesornis (0)	Sele Holding		Mark Report and the Control of the C
- archaeornis (0)	the state of the s	W/	1. 3.5
ratite, ratite bird, flightless bird (10)		W	-//
- carinate, carinate bird, flying bird (0)			
passerine, passeriform bird (279)			200
nonpasserine bird (0)	The state of the s	. 1	
i⊸ bird of prey, raptor, raptorial bird (80)			
gallinaceous bird, gallinacean (114)	The same of the sa		

Not logged in. Login I Signup

German iris, Iris kochii

Iris of northern Italy having deep blue-purple flowers; similar to but smaller than Iris germanica

469 pictures 49.6% Popularity Percentile



halophyte (0)	
succulent (39)	
cultivar (0)	
 cultivated plan 	t (0)
weed (54)	
- evergreen, eve	rgreen plant (0)
- deciduous plan	t (0)
vine (272)	
creeper (0)	
woody plant, li	gneous plant (1868)
geophyte (0)	
	erophyte, xerophytic plant, xerophile, xerophile esophytic plant (0)
	water plant, hydrophyte, hydrophytic plant (11
tuberous plant	
bulbous plant (
iris, flag,	fleur-de-lis, sword lily (19)
. beard	led iris (4)
Flo	orentine iris, orris, Iris germanica florentina, Iris
- Ge	erman iris, Iris germanica (0)
Ge	erman iris, Iris kochii (0)
- Da	ılmatian iris, Iris pallida (0)
⊩ beard	lless iris (4)
bulbo	us iris (0)
dwarf	iris, Iris cristata (0)
stinki	ng iris, gladdon, gladdon iris, stinking gladwyn,
Persia	an iris, Iris persica (0)
- yellov	v iris, yellow flag, yellow water flag, Iris pseuda
- dwarf	iris, vernal iris, Iris verna (0)
- blue f	lag, Iris versicolor (0)



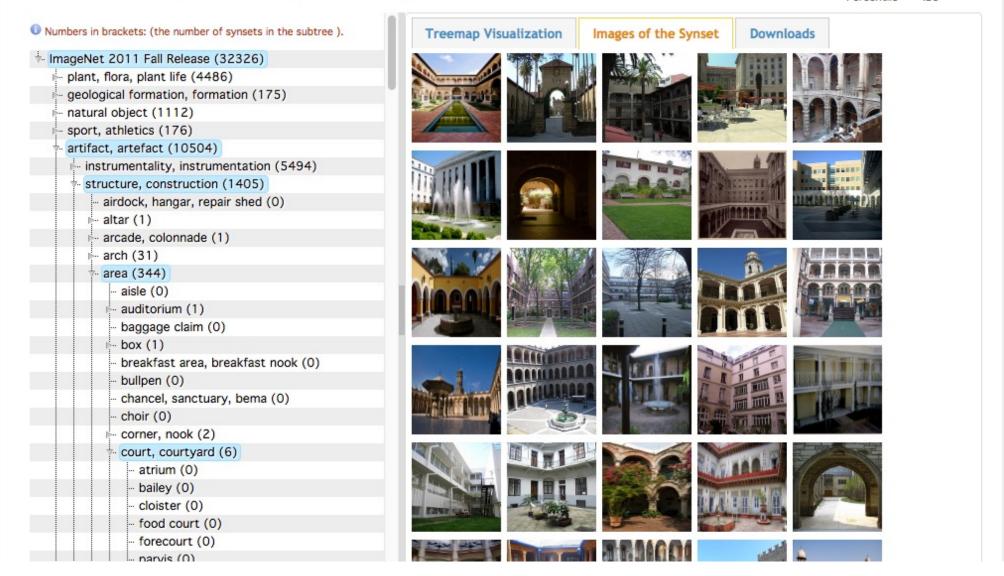
Not logged in. Login I Signup

Court, courtyard

An area wholly or partly surrounded by walls or buildings; "the house was built around an inner court"

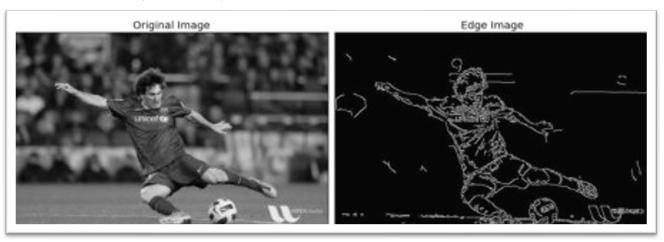
165 pictures 92.61% Popularity Percentile



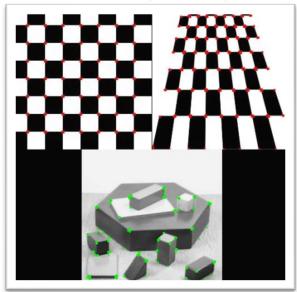


Feature Engineering for CV

Edge detection (Canny)

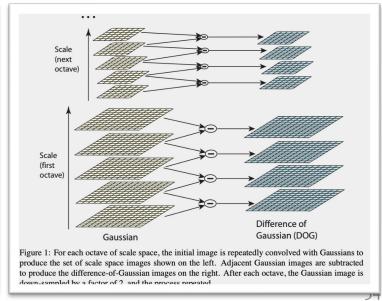


Corner Detection (Harris)



Scale Invariant Feature Transform (SIFT)





Figures from http://opencv.org

Figure from Lowe (1999) and Lowe (2004)

Example: Image Classification

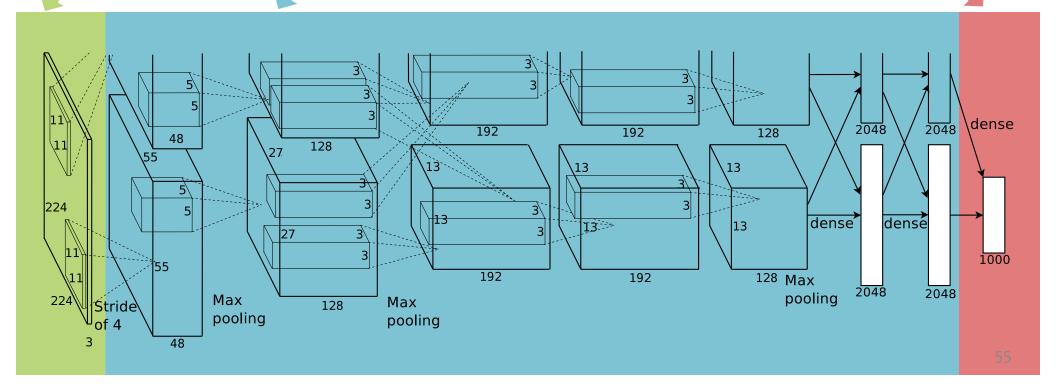
CNN for Image Classification

(Krizhevsky, Sutskever & Hinton, 2012) 15.3% error on ImageNet LSVRC-2012 contest

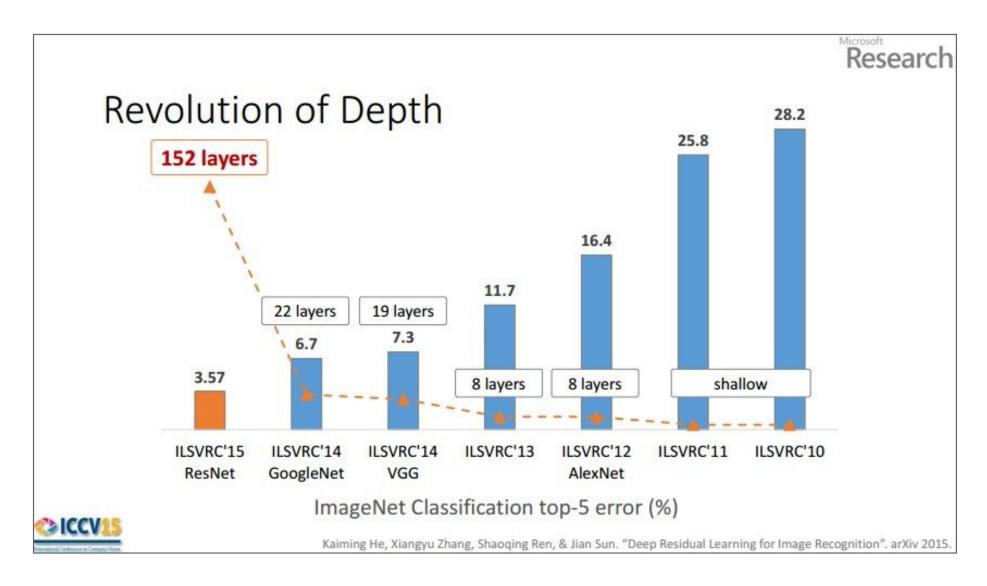
Input image (pixels)

- Five convolutional layers (w/max-pooling)
- Three fully connected layers

1000-way softmax



CNNs for Image Recognition



Backpropagation and Deep Learning

Convolutional neural networks (CNNs) and recurrent neural networks (RNNs) are simply fancy computation graphs (aka. hypotheses or decision functions).

Our recipe also applies to these models and (again) relies on the **backpropagation algorithm** to compute the necessary gradients.

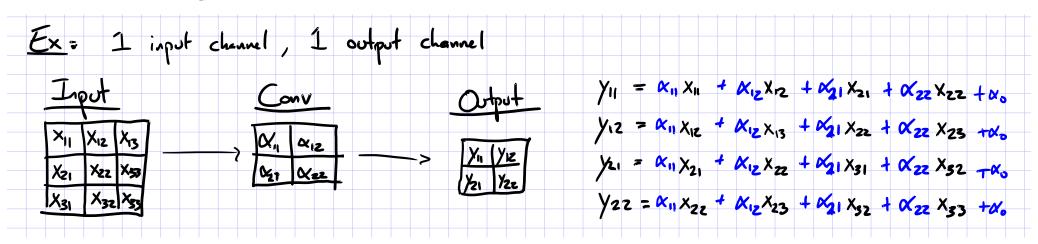
CONVOLUTION

Basic idea:

- Pick a 3x3 matrix F of weights
- Slide this over an image and compute the "inner product" (similarity) of F and the corresponding field of the image, and replace the pixel in the center of the field with the output of the inner product operation

Key point:

- Different convolutions extract different types of low-level "features" from an image
- All that we need to vary to generate these different features is the weights of F



A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image

О	0	0	0	0	0	О
0	1	1	1	1	1	0
О	1	0	0	1	0	О
0	1	0	1	0	0	0
0	1	1	0	0	0	О
О	1	0	0	0	0	О
0	0	0	0	0	0	0

Convolution

О	0	0
0	1	1
О	1	0

Convolved Image

1	1	1	1	1
1	0	0	1	0
1	0	1	0	0
1	1	0	0	0
1	0	0	0	0

Input Image

0	0	0	0	0	0	0
0	1	1	1	1	1	0
0	1	0	0	1	0	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	О
0	0	0	0	0	0	0

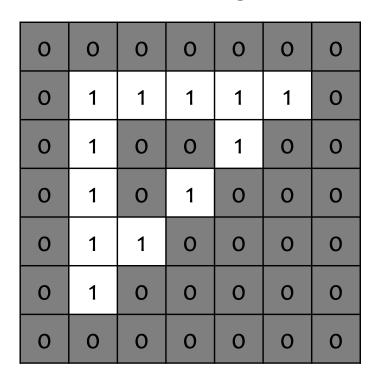


0	0	0
0	1	1
О	1	0

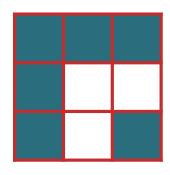
Convolved Image

3	2	2	3	1
2	0	2	1	0
2	2	1	0	0
3	1	0	0	0
1	0	0	0	0

Input Image





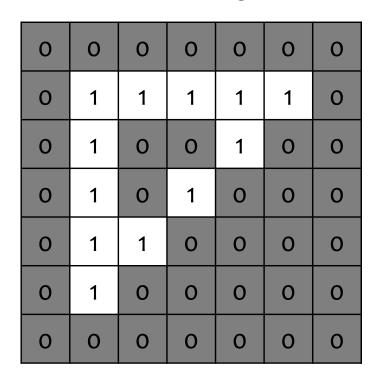


Convolved Image

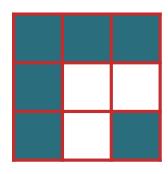
3	2	2	3	1
2	0	2	1	0
2	2	1	0	0
3	1	0	0	0
1	0	0	0	0

A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image



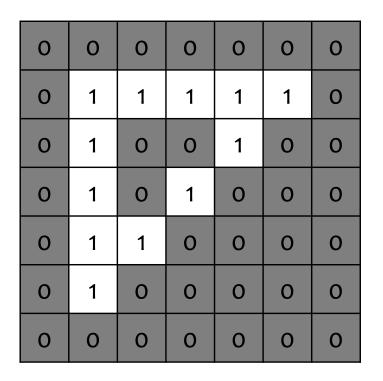




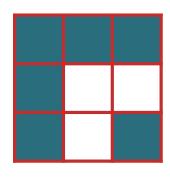
Convolved Image

3	2	2	3	1
2	0	2	1	0
2	2	1	0	0
3	1	0	0	0
1	0	0	0	0

Input Image





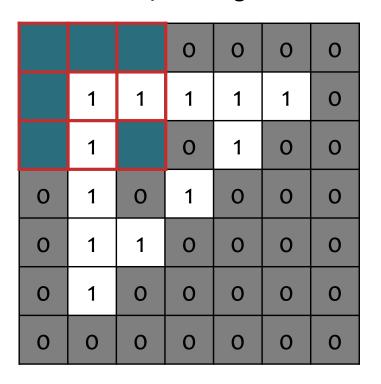


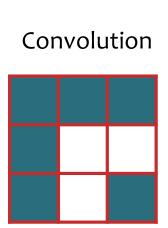
Convolved Image

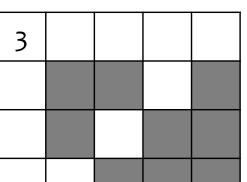
3	2	2	3	1
2	0	2	1	0
2	2	1	0	0
3	1	0	0	0
1	0	0	0	0

A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image

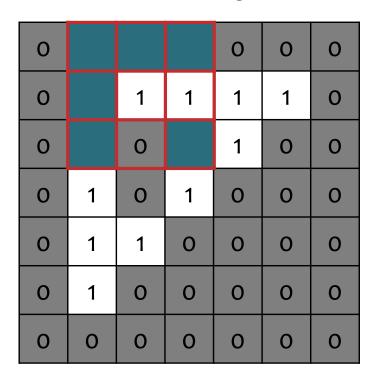


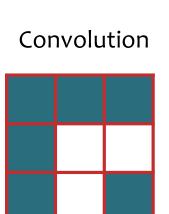




Convolved Image

Input Image

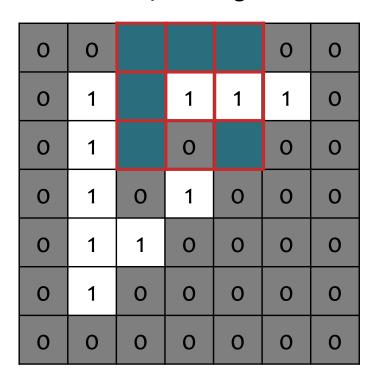


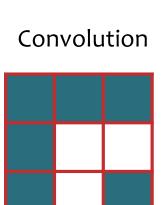




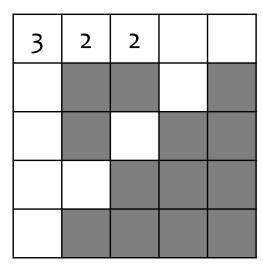
3	2		

Input Image

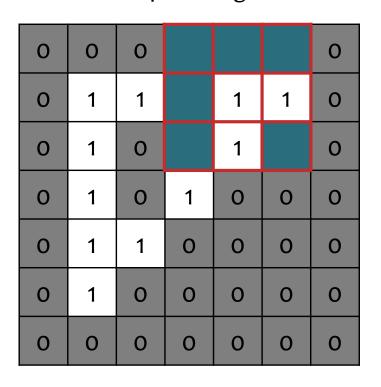




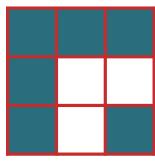




Input Image



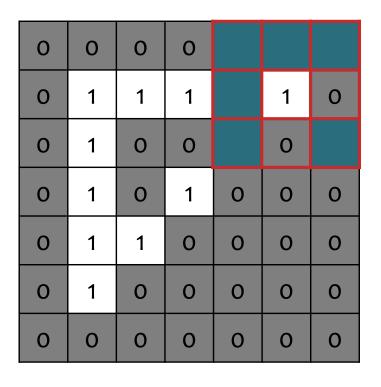




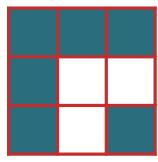
Convolved Image

3	2	2	3	

Input Image



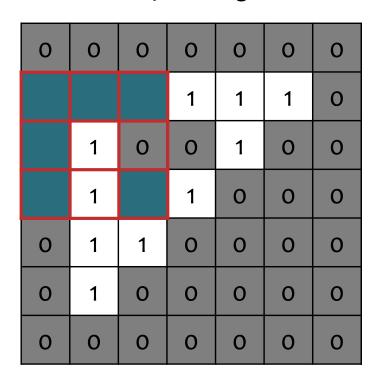




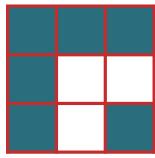
Convolved Image

3	2	2	3	1

Input Image



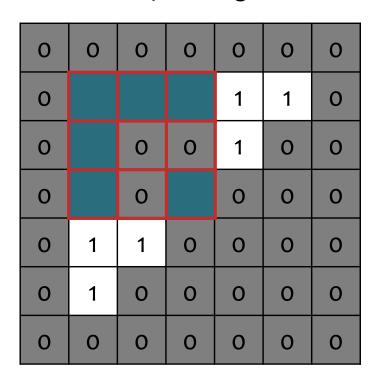




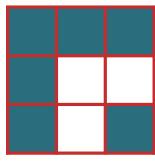
Convolved Image

3	2	2	3	1
2				

Input Image





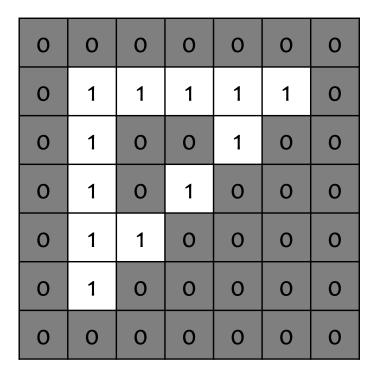


Convolved Image

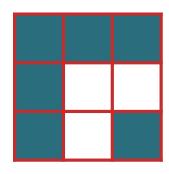
3	2	2	3	1
2	0			

A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image







Convolved Image

3	2	2	3	1
2	0	2	1	0
2	2	1	0	0
3	1	0	0	0
1	0	0	0	0

Input Image

0	0	0	0	0	0	0
0	1	1	1	1	1	0
0	1	0	0	1	0	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0

Identity Convolution

0	0	0
О	1	0
0	0	0

Convolved Image

1	1	1	1	1
1	0	0	1	0
1	0	1	0	0
1	1	0	0	0
1	0	0	0	0

Input Image

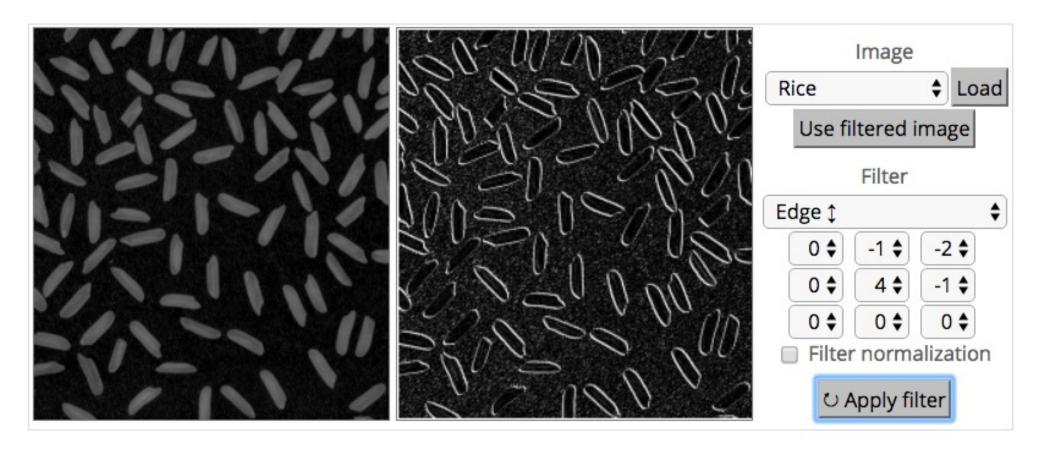
0	0	0	0	0	0	0
0	1	1	1	1	1	0
0	1	0	0	1	0	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	О
0	0	0	0	0	0	0

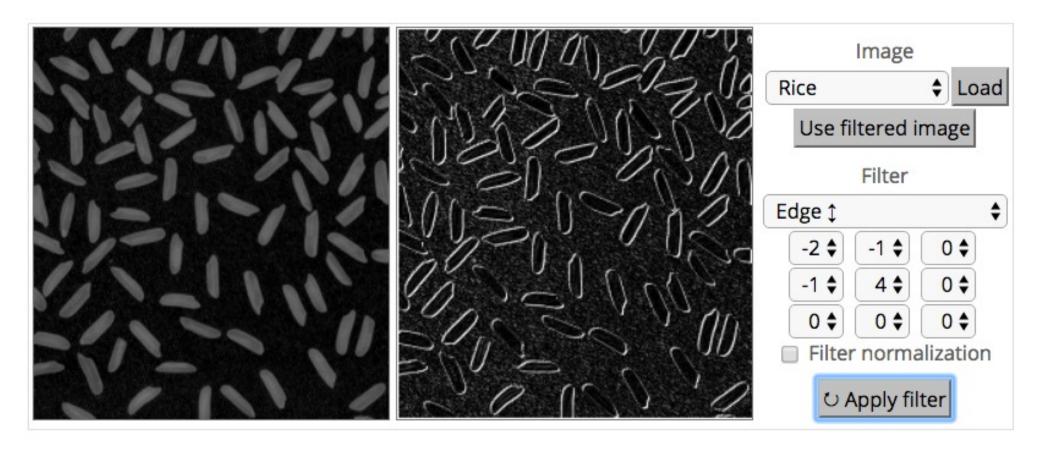
Blurring Convolution

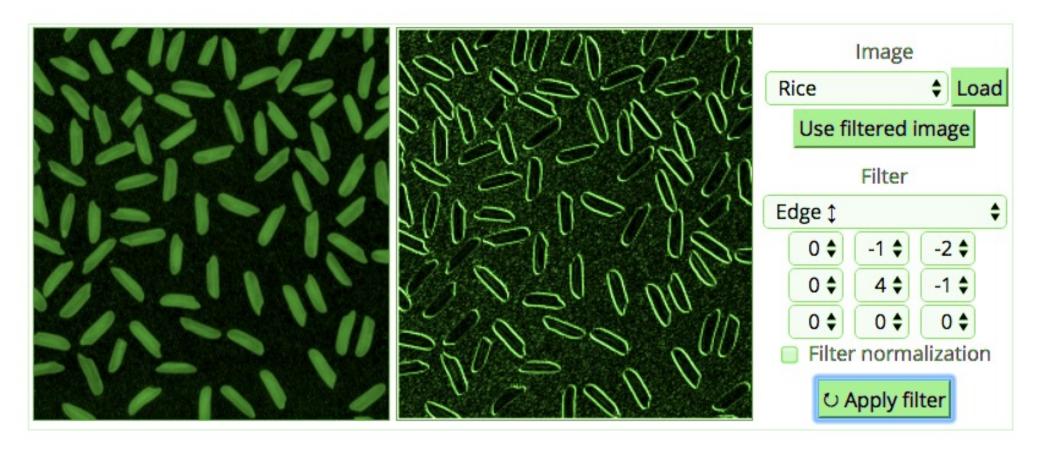
.1	.1	.1
.1	.2	.1
.1	.1	.1

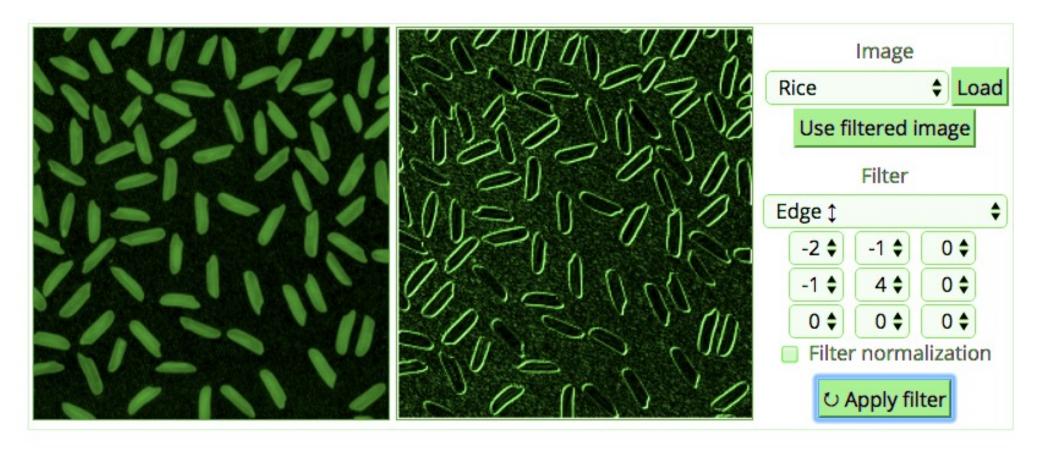
Convolved Image

.4	.5	.5	.5	.4
•4	.2	•3	.6	.3
.5	•4	.4	.2	.1
.5	.6	.2	.1	0
.4	.3	.1	0	0







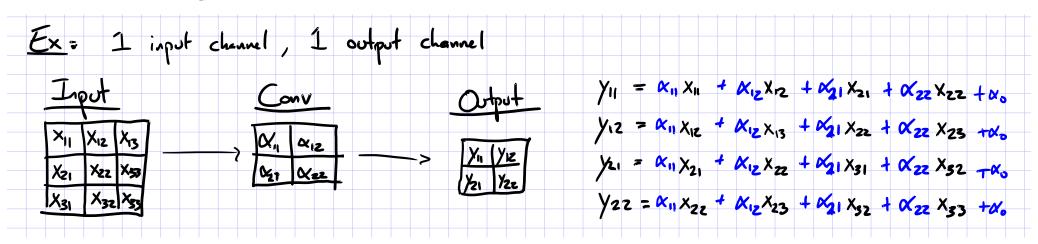


Basic idea:

- Pick a 3x3 matrix F of weights
- Slide this over an image and compute the "inner product" (similarity) of F and the corresponding field of the image, and replace the pixel in the center of the field with the output of the inner product operation

Key point:

- Different convolutions extract different types of low-level "features" from an image
- All that we need to vary to generate these different features is the weights of F



DOWNSAMPLING

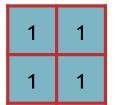
Downsampling

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

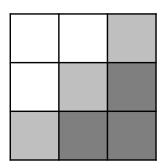
Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution



Convolved Image



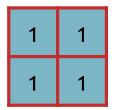
Downsampling

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

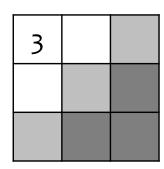
Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution



Convolved Image



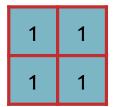
Downsampling

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

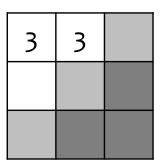
Input Image

1	1	1	1	1	0
1	0	О	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution



Convolved Image

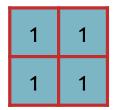


- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

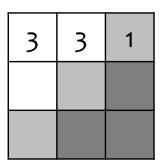
Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution



Convolved Image

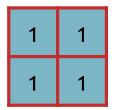


- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

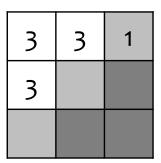
Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution



Convolved Image

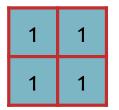


- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

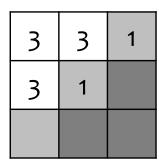
Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution



Convolved Image

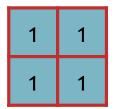


- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

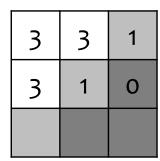
Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution



Convolved Image

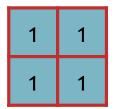


- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution



Convolved Image

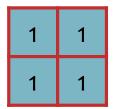
3	3	1
3	1	0
1		

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution



Convolved Image

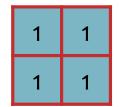
3	3	1
3	1	0
1	0	

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution



Convolved Image

3	3	1
3	1	0
1	0	0

Downsampling by Averaging

- Downsampling by averaging is a special case of convolution where the weights are fixed to a uniform distribution
- The example below uses a stride of 2

Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution

1/4	1/4
1/4	1/4

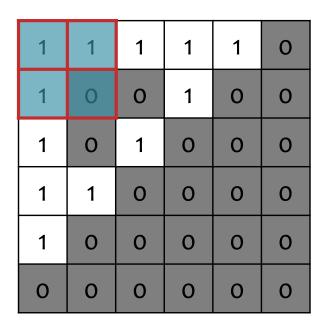
Convolved Image

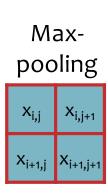
3/4	3/4	1/4
3/4	1/4	0
1/4	0	0

Max-Pooling

- Max-pooling is another form of downsampling
- Instead of averaging, we take the max value within the same range as the equivalently-sized convolution
- The example below uses a stride of 2

Input Image







1	1	1
1	1	0
1	0	0

$$y_{ij} = \max(x_{ij}, x_{i,j+1}, x_{i+1,j}, x_{i+1,j+1})$$

CONVOLUTIONAL NEURAL NETS

Background

A Recipe for Machine Learning

1. Given training data:

$$\{oldsymbol{x}_i,oldsymbol{y}_i\}_{i=1}^N$$

- 2. Choose each of these:
 - Decision function

$$\hat{\boldsymbol{y}} = f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)$$

Loss function

$$\ell(\hat{m{y}},m{y}_i)\in\mathbb{R}$$

3. Define goal:

$$oldsymbol{ heta}^* = rg\min_{oldsymbol{ heta}} \sum_{i=1}^N \ell(f_{oldsymbol{ heta}}(oldsymbol{x}_i), oldsymbol{y}_i)$$

4. Train with SGD:

(take small steps opposite the gradient)

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta_t \nabla \ell(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i), \boldsymbol{y}_i)$$

Background

A Recipe for Machine Learning

- Convolutional Neural Networks (CNNs) provide another form of decision function
 - Let's see what they look like...

2. Choose each of these:

Decision function

$$\hat{\boldsymbol{y}} = f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)$$

Loss function

$$\ell(\hat{m{y}}, m{y}_i) \in \mathbb{R}$$

Train with SGD:

ke small steps
opposite the gradient)

$$oldsymbol{ heta}^{(t+1)} = oldsymbol{ heta}^{(t)} - \eta_t
abla \ell(f_{oldsymbol{ heta}}(oldsymbol{x}_i), oldsymbol{y}_i)$$

Convolutional Layer

CNN key idea:

Treat convolution matrix as parameters and learn them!

Input Image

0	0	0	0	0	0	0
0	1	1	1	1	1	0
0	1	0	0	1	0	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	О
0	0	0	0	0	0	О



Learned Convolution

θ_{11}	θ_{12}	θ_{13}
θ_{21}	θ_{22}	θ_{23}
θ_{31}	θ_{32}	θ_{33}

Convolved Image

.4	.5	.5	.5	.4
•4	.2	•3	.6	.3
•5	.4	•4	.2	.1
.5	.6	.2	.1	0
.4	.3	.1	0	0

Convolutional Neural Network (CNN)

- Typical layers include:
 - Convolutional layer
 - Max-pooling layer
 - Fully-connected (Linear) layer
 - ReLU layer (or some other nonlinear activation function)
 - Softmax
- These can be arranged into arbitrarily deep topologies

Architecture #1: LeNet-5

PROC. OF THE IEEE, NOVEMBER 1998

7

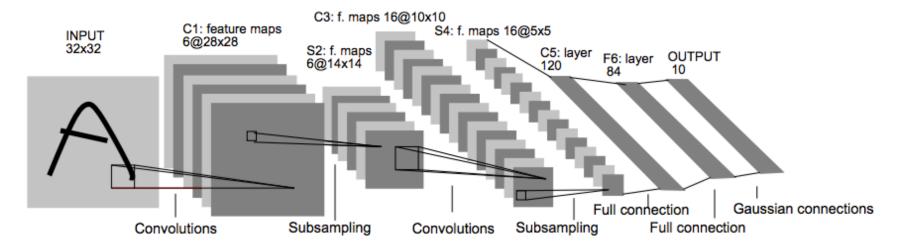


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

TRAINING CNNS

Background

A Recipe for Machine Learning

3. Define goal:

1. Given training data:

$$\{oldsymbol{x}_i, oldsymbol{y}_i\}_{i=1}^N$$

- $oldsymbol{ heta}^* = rg\min_{oldsymbol{ heta}} \sum_i \ell(f_{oldsymbol{ heta}}(oldsymbol{x}_i), oldsymbol{y}_i)$
- 2. Choose each of these:
 - Decision function

$$\hat{\boldsymbol{y}} = f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)$$

Loss function

$$\ell(\hat{m{y}},m{y}_i)\in\mathbb{R}$$

4. Train with SGD:

(take small steps opposite the gradient)

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta_t \nabla \ell(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i), \boldsymbol{y}_i)$$

Background

A Recipe for Machine Learning

1. Given training data:

$$\{oldsymbol{x}_i,oldsymbol{y}_i\}_{i=1}^N$$

- 2. Choose each of the
 - Decision function

$$\hat{\boldsymbol{y}} = f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)$$

Loss function

$$\ell(\hat{m{y}},m{y}_i)\in\mathbb{R}$$

3. Define goal:

- $\{\boldsymbol{x}_i,\boldsymbol{y}_i\}_{i=1}^N$ Q: Now that we have the CNN as a decision function, how do we compute the gradient?
 - A: Backpropagation of course!

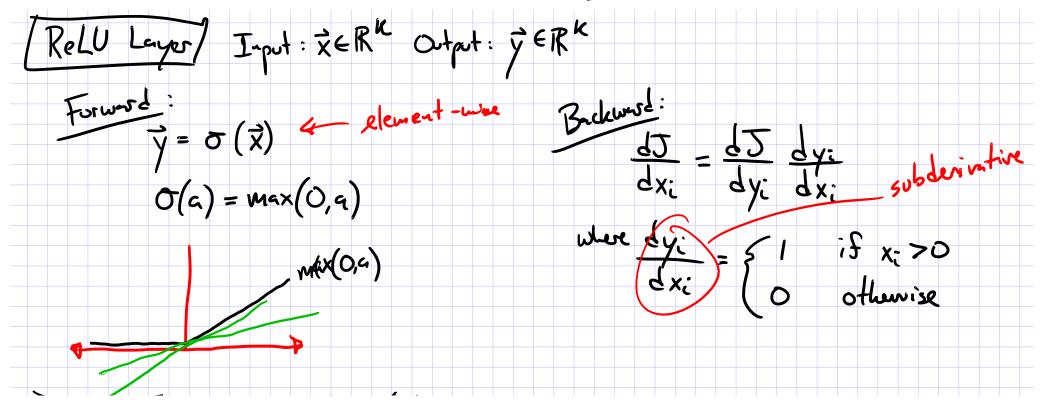
opposite the gradient)
$$\boldsymbol{\theta}^{(t)} = \boldsymbol{\eta}_t \nabla \ell(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i), \boldsymbol{y}_i)$$

SGD for CNNs

$$\begin{array}{lll}
\hline SGD & for CNN_{5} \\
\hline Ex: Architecture: & Given \vec{x}, \vec{y}^{*} \\
\hline J = l(y, y^{*}) \\
y = softmx(z^{(5)}) & Parameters $\vec{\Theta} = [\times, \beta, W] \\
z^{(5)} = linear(z^{(4)}, W) \\
z^{(4)} = relu(z^{(3)}) & SGD: \\
z^{(3)} = (conv(z^{(2)}, \beta) & DI_{4}; \vec{\Theta} \\
z^{(2)} = (conv(z^{(2)}, \beta) & DI_{4}; \vec{\Theta} \\
z^{(2)} = mx-pool(z^{(1)}) & Sample i \in \{1, ..., W\} \\
z^{(1)} = conv(\vec{x}, \infty) & Forward: y = h_{\Theta}(\vec{x}^{(1)}), J_{1}(\vec{\Theta}) = l(y, y^{*}) \\
Backward: V_{\vec{\Theta}}J_{1}(\vec{\Theta}) = ... \\
\hline
Vertice: $\vec{\Theta} \leftarrow \vec{\Theta} - NV_{\vec{\Theta}}J_{1}(\vec{\Theta})
\end{array}$$$$

LAYERS OF A CNN

ReLU Layer



Softmax Layer

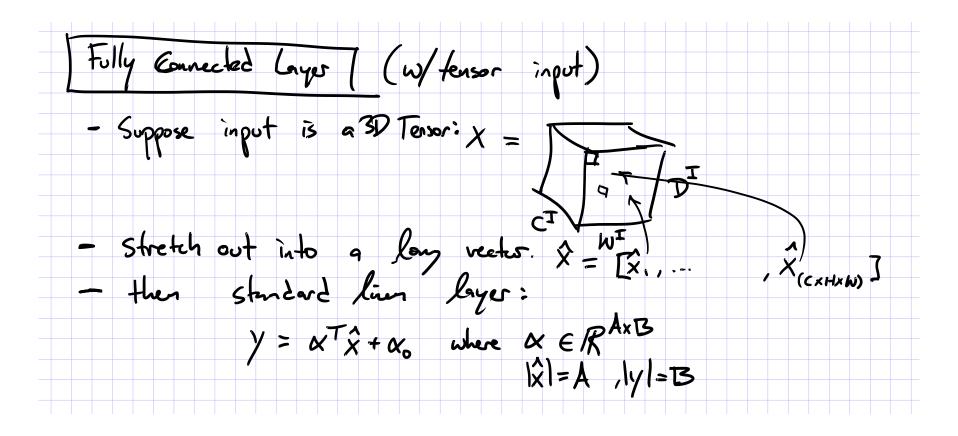
Softmax Layer

Input:
$$\vec{x} \in \mathbb{R}^{K}$$
 Dutput: $\vec{y} \in \mathbb{R}^{K}$

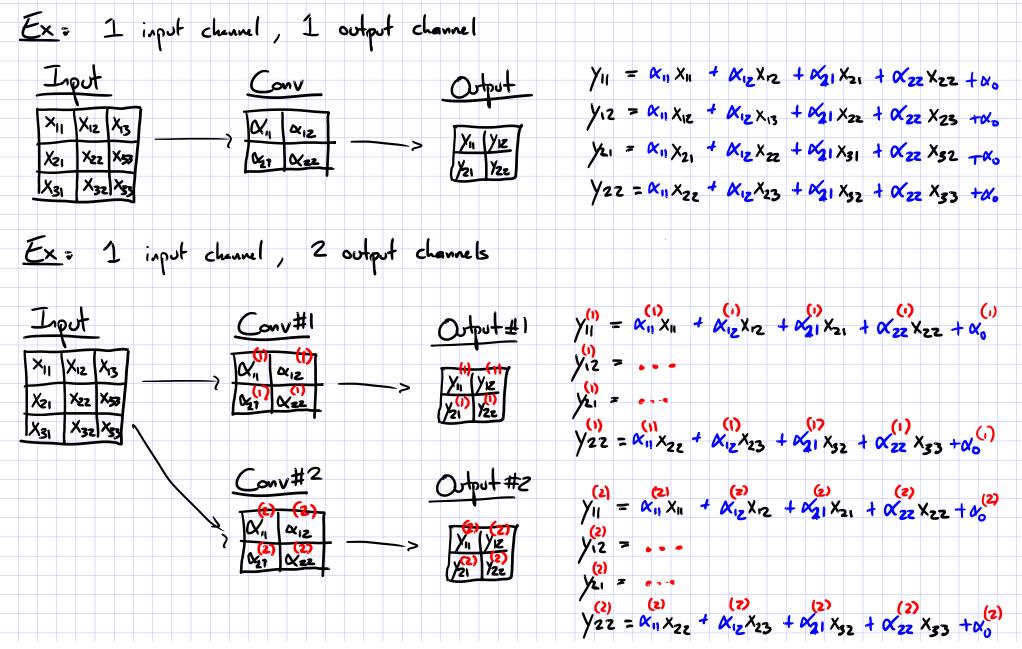
Forward:

 $y_i = \exp(x_i)$
 $\exists x_i \in \mathbb{R}^{K}$
 $\exists x_i \in \mathbb{R$

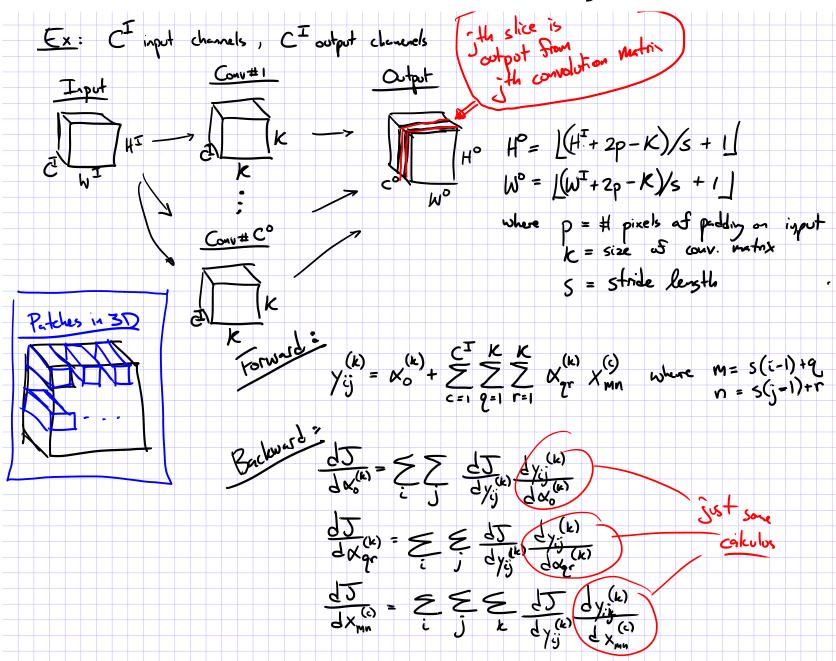
Fully-Connected Layer



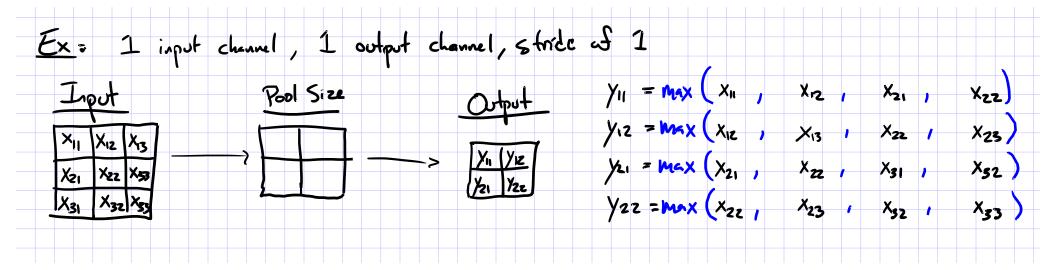
Convolutional Layer



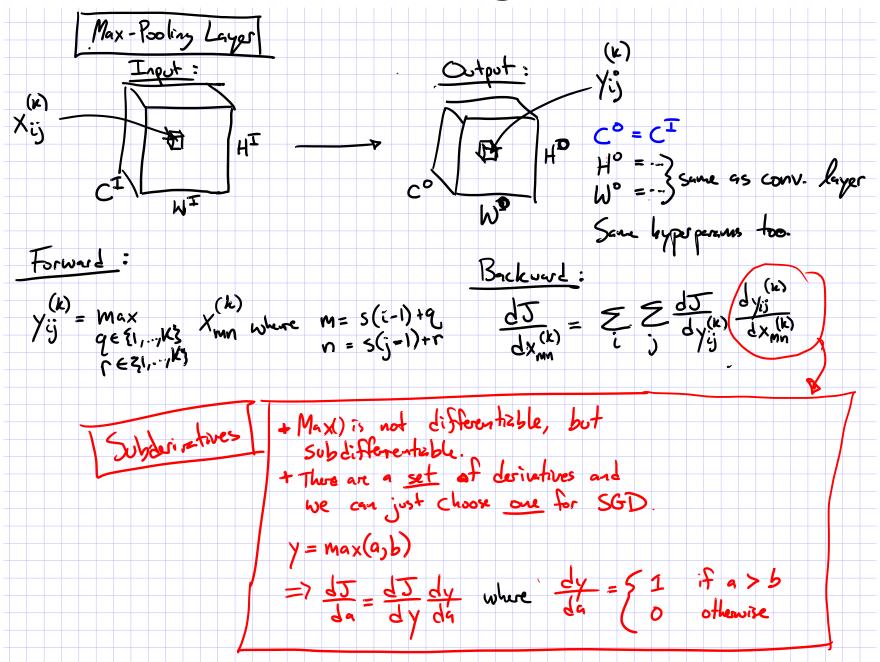
Convolutional Layer



Max-Pooling Layer



Max-Pooling Layer



Convolutional Neural Network (CNN)

- Typical layers include:
 - Convolutional layer
 - Max-pooling layer
 - Fully-connected (Linear) layer
 - ReLU layer (or some other nonlinear activation function)
 - Softmax
- These can be arranged into arbitrarily deep topologies

Architecture #1: LeNet-5

PROC. OF THE IEEE, NOVEMBER 1998

7

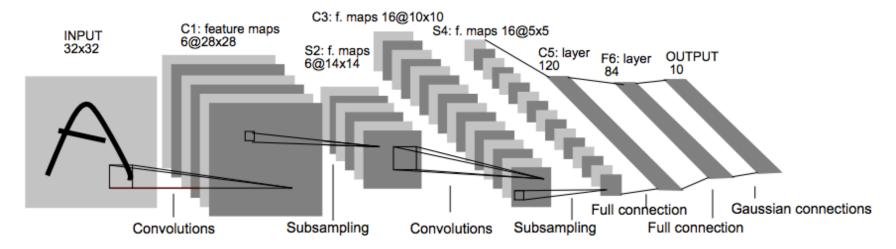


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

Architecture #2: AlexNet

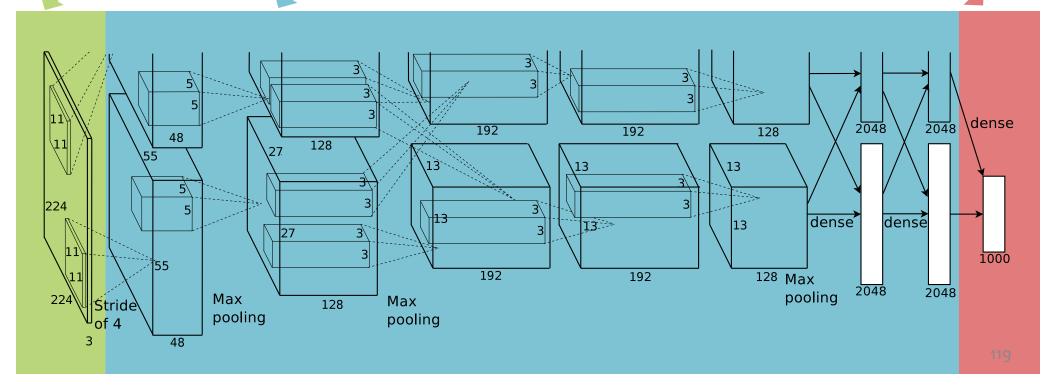
CNN for Image Classification

(Krizhevsky, Sutskever & Hinton, 2012) 15.3% error on ImageNet LSVRC-2012 contest

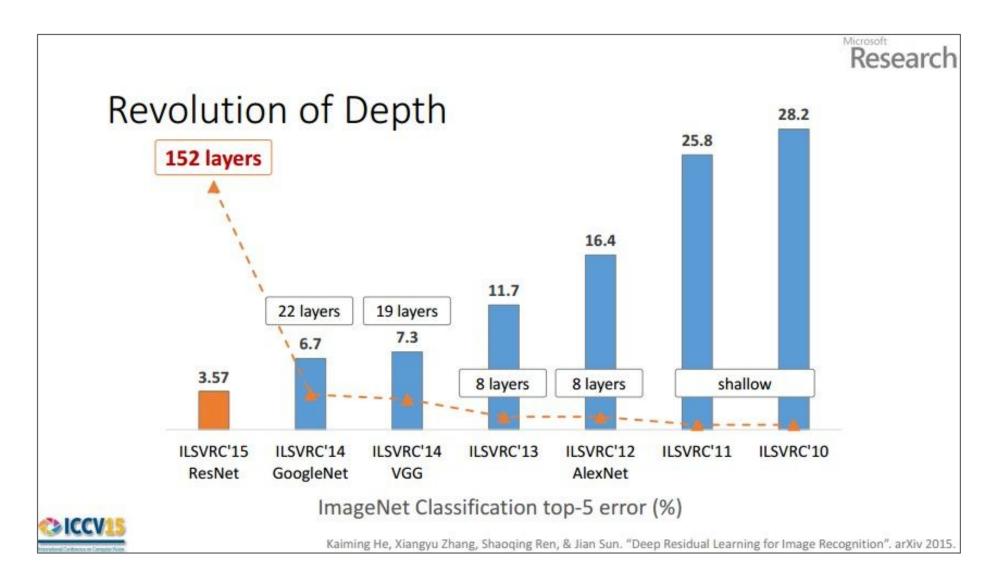
Input image (pixels)

- Five convolutional layers (w/max-pooling)
- Three fully connected layers

1000-way softmax



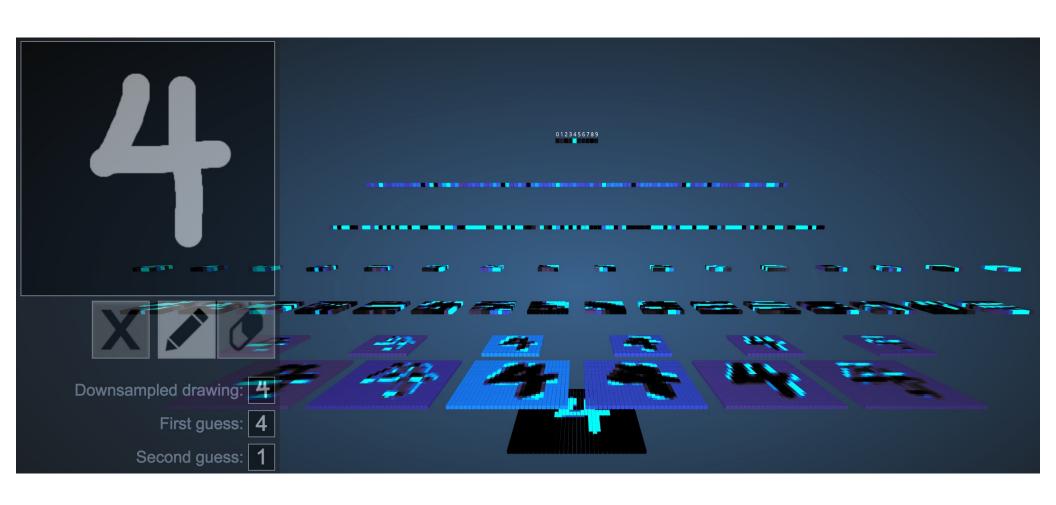
CNNs for Image Recognition



CNN VISUALIZATIONS

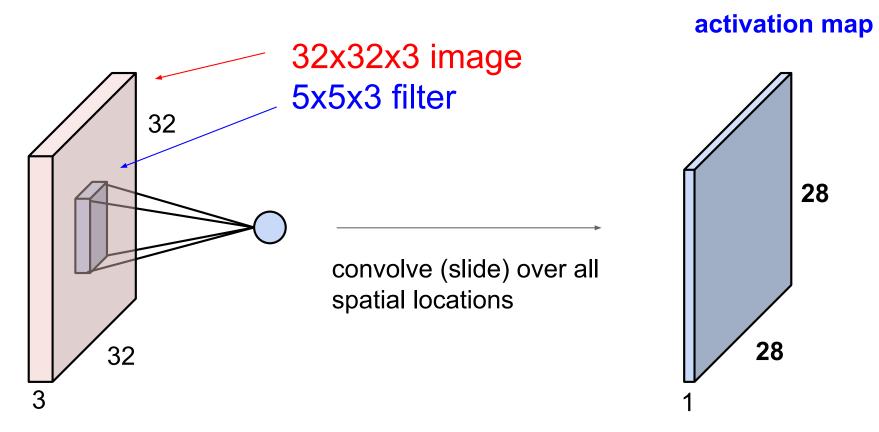
3D Visualization of CNN

http://scs.ryerson.ca/~aharley/vis/conv/



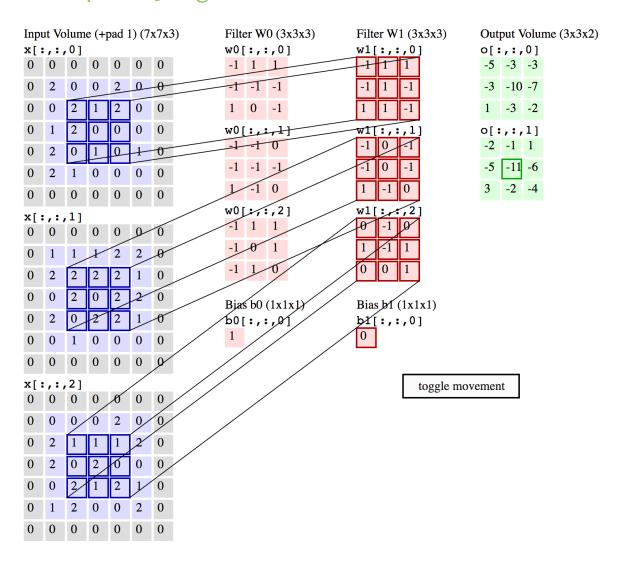
Convolution of a Color Image

- Color images consist of 3 floats per pixel for RGB (red, green blue) color values
- Convolution must also be 3-dimensional



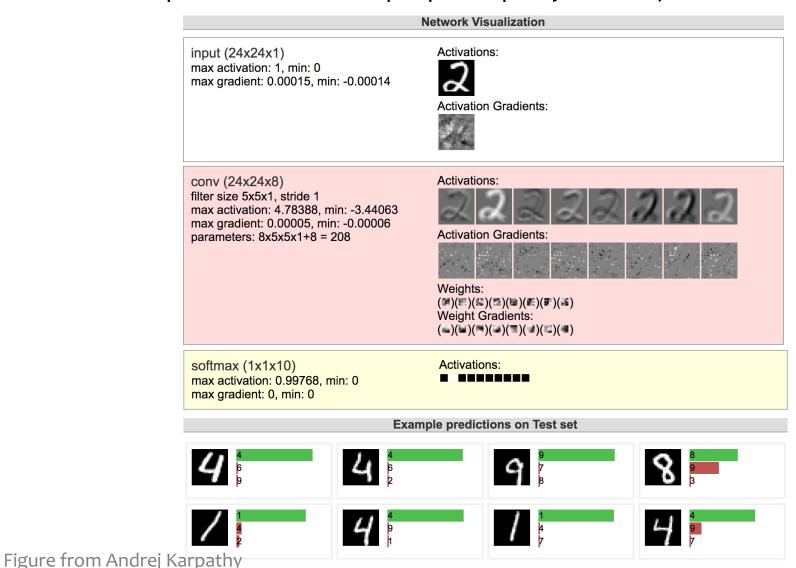
Animation of 3D Convolution

http://cs231n.github.io/convolutional-networks/



MNIST Digit Recognition with CNNs (in your browser)

https://cs.stanford.edu/people/karpathy/convnetjs/demo/mnist.html



CNN Summary

CNNs

- Are used for all aspects of computer vision, and have won numerous pattern recognition competitions
- Able learn interpretable features at different levels of abstraction
- Typically, consist of convolution layers, pooling layers, nonlinearities, and fully connected layers

Other Resources:

- Readings on course website
- Andrej Karpathy, CS231n Notes
 http://cs231n.github.io/convolutional-networks/

Deep Learning Objectives

You should be able to...

- Implement the common layers found in Convolutional Neural Networks (CNNs) such as linear layers, convolution layers, max-pooling layers, and rectified linear units (ReLU)
- Explain how the shared parameters of a convolutional layer could learn to detect spatial patterns in an image
- Describe the backpropagation algorithm for a CNN
- Identify the parameter sharing used in a basic recurrent neural network, e.g. an Elman network
- Apply a recurrent neural network to model sequence data
- Differentiate between an RNN and an RNN-LM

ML Big Picture

Learning Paradigms:

What data is available and when? What form of prediction?

- supervised learning
- unsupervised learning
- semi-supervised learning
- reinforcement learning
- active learning
- imitation learning
- domain adaptation
- online learning
- density estimation
- recommender systems
- feature learning
- manifold learning
- dimensionality reduction
- ensemble learning
- distant supervision
- hyperparameter optimization

Theoretical Foundations:

What principles guide learning?

- probabilistic
- ☐ information theoretic
- evolutionary search
- ☐ ML as optimization

Problem Formulation:

What is the structure of our output prediction?

boolean Binary Classification

categorical Multiclass Classification

ordinal Ordinal Classification

real Regression

ordering Ranking

multiple discrete Structured Prediction

multiple continuous (e.g. dynamical systems)

both discrete & (e.g. mixed graphical models)

cont.

Application Areas

Key challenges?

NLP, Speech, Computer
Vision, Robotics, Medicine

Facets of Building ML Systems:

How to build systems that are robust, efficient, adaptive, effective?

- 1. Data prep
- 2. Model selection
- Training (optimization / search)
- 4. Hyperparameter tuning on validation data
- 5. (Blind) Assessment on test

Big Ideas in ML:

Which are the ideas driving development of the field?

- inductive bias
- generalization / overfitting
- bias-variance decomposition
- generative vs. discriminative
- deep nets, graphical models
- PAC learning
- distant rewards