

MSIS549: ML and AI for Business Applications

Lesson 4



Course Outline

Lesson 1: Introduction to Deep Learning

Lesson 2: Neural Network Fundamentals

Lesson 3: Convolutional Neural Networks

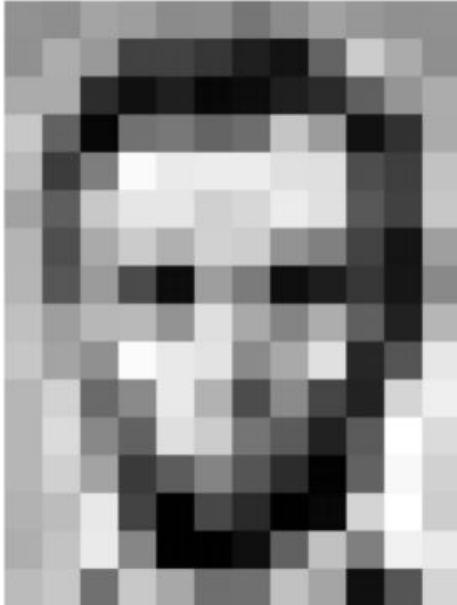
Lesson 4: Recurrent Neural Networks

Lesson 5: Deep Learning Practice

Recap



What Computers “See”



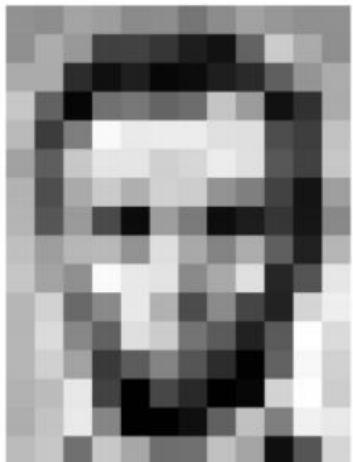
157	153	174	148	150	162	129	151	172	161	165	156
155	182	163	74	75	62	33	17	110	210	180	154
180	180	50	14	34	6	10	33	48	106	159	181
206	109	5	124	131	111	120	204	166	15	56	180
194	68	137	251	237	239	239	228	227	87	71	201
172	105	207	233	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	158	128	75	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	193	158	227	178	143	182	106	36	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	108	227	210	127	103	36	101	255	224
190	214	173	66	103	143	96	50	2	109	249	215
187	196	235	75	1	81	47	0	6	217	255	211
183	202	237	146	0	0	12	108	200	138	243	236
195	206	123	207	177	121	123	200	175	13	96	218

What the computer sees

157	153	174	168	150	152	129	151	172	161	165	156
155	182	163	74	75	62	33	17	110	210	180	154
180	180	50	14	34	6	10	33	48	106	159	181
206	109	5	124	131	111	120	204	166	15	56	180
194	68	137	251	237	239	239	228	227	87	71	201
172	105	207	233	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	158	128	75	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	193	158	227	178	143	182	106	36	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	108	227	210	127	103	36	101	255	224
190	214	173	66	103	143	96	50	2	109	249	215
187	196	235	75	1	81	47	0	6	217	255	211
183	202	237	146	0	0	12	108	200	138	243	236
195	206	123	207	177	121	123	200	175	13	96	218

An image is just a matrix of numbers [0,255]!
i.e., 1080x1080x3 for an RGB image

Tasks in Computer Vision



Input Image



167	153	174	168	150	152	129	161	172	161	155	156
155	182	163	74	75	62	33	17	110	210	180	154
180	180	58	14	34	6	10	33	48	106	159	181
206	109	6	124	131	111	120	204	166	15	56	180
194	68	137	251	237	239	239	228	227	87	71	207
172	105	207	233	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	168	139	75	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	193	198	227	178	143	182	106	36	190
205	174	165	252	236	231	149	178	208	43	95	234
190	216	116	149	236	187	86	160	79	36	218	241
190	224	147	108	227	210	127	102	36	101	255	224
190	214	173	66	103	143	96	50	2	109	249	215
187	196	239	76	1	81	47	0	6	217	255	211
189	202	237	145	0	0	12	108	200	138	243	236
196	206	123	207	177	121	123	200	175	13	96	218

Pixel Representation

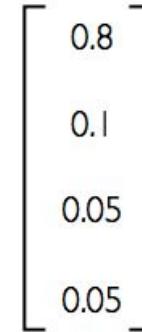
classification

Lincoln

Washington

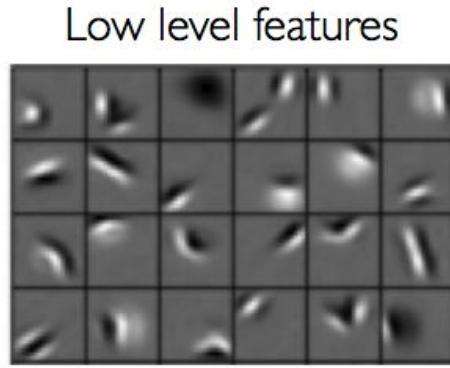
Jefferson

Obama

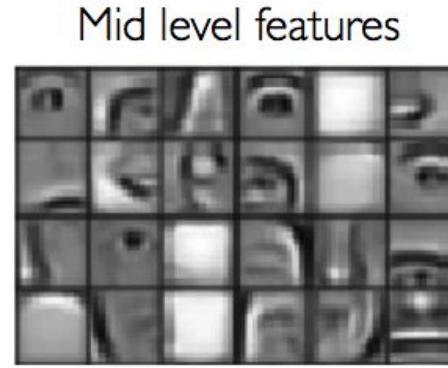


Learning Feature Representations

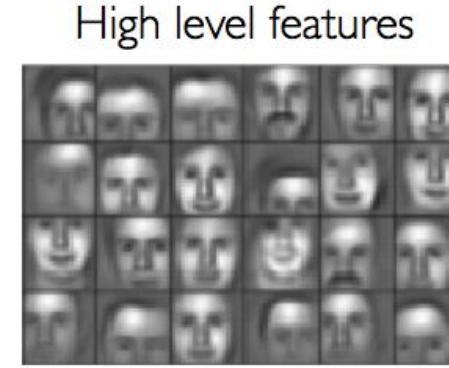
Can we learn a **hierarchy of features** directly from the data instead of hand engineering?



Edges, dark spots

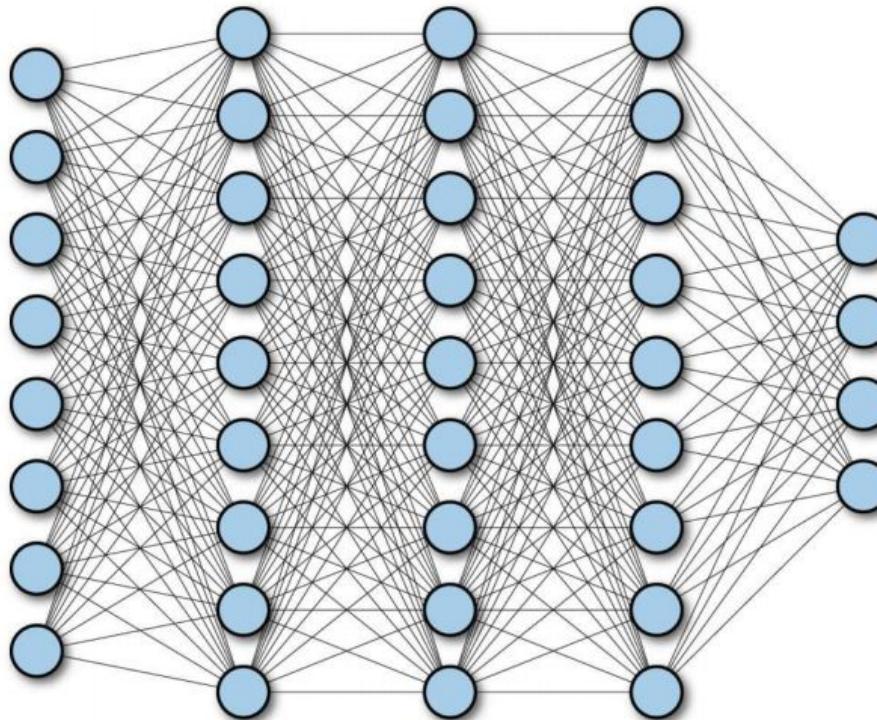


Eyes, ears, nose



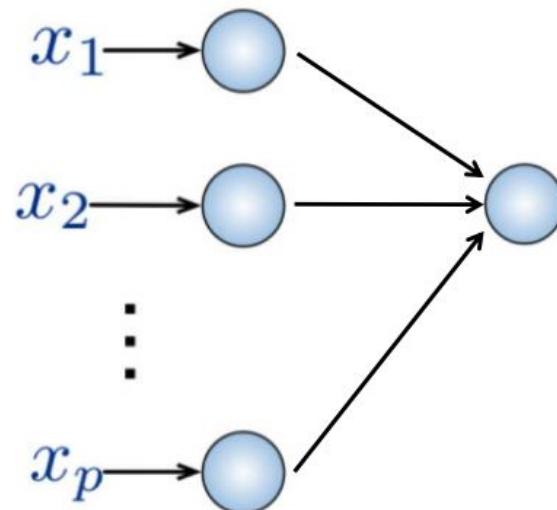
Facial structure

Fully Connected Neural Network



Fully Connected Neural Network

- Input:**
- 2D image
 - Vector of pixel values



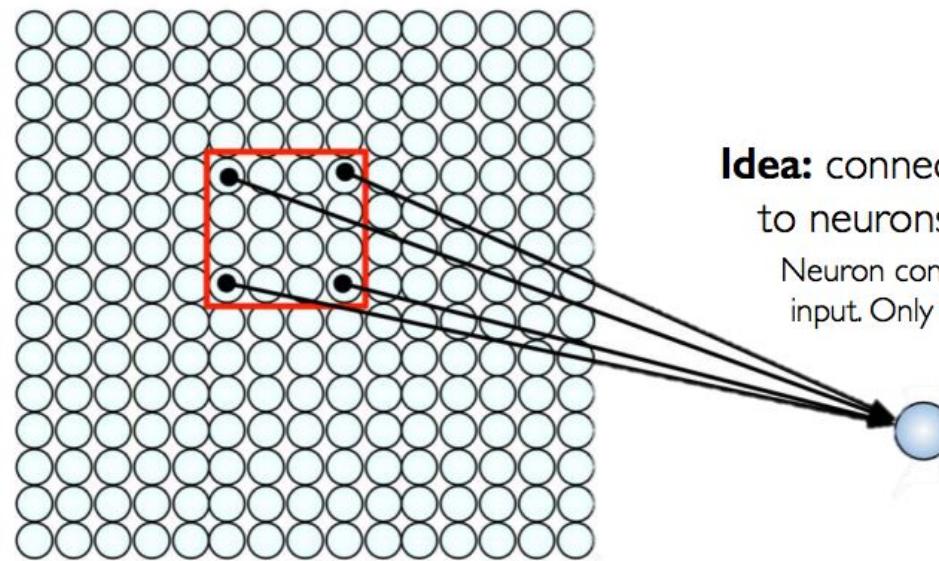
Fully Connected:

- Connect neuron in hidden layer to all neurons in input layer
- No spatial information!
- And many, many parameters!

How can we use **spatial structure** in the input to inform the architecture of the network?

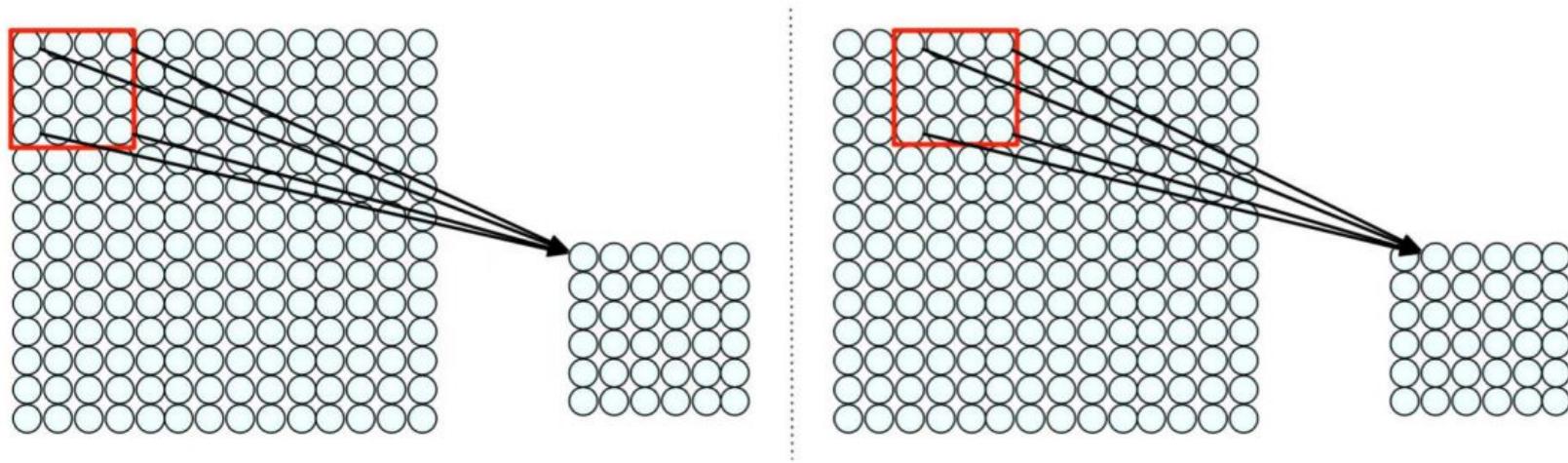
Using Spatial Structure

Input: 2D image.
Array of pixel values



Idea: connect patches of input
to neurons in hidden layer.
Neuron connected to region of
input. Only "sees" these values.

Using Spatial Structure



Connect patch in input layer to a single neuron in subsequent layer.

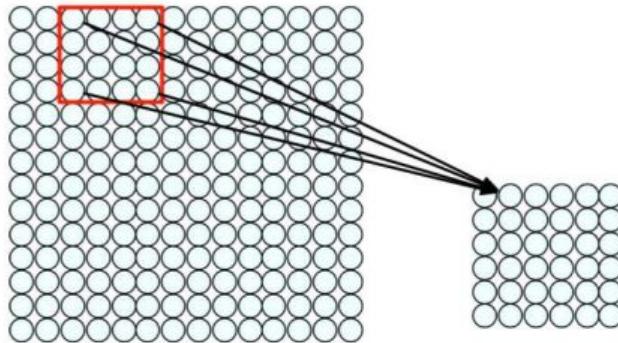
Use a sliding window to define connections.

*How can we **weight** the patch to detect particular features?*

Applying Filters to Extract Features

- 1) Apply a set of weights – a filter – to extract **local features**
- 2) Use **multiple filters** to extract different features
- 3) Spatially **share** parameters of each filter
(features that matter in one part of the input should matter elsewhere)

Feature Extraction with Convolution



- Filter of size 4×4 : 16 different weights
- Apply this same filter to 4×4 patches in input
- Shift by 2 pixels for next patch

This “patchy” operation is **convolution**

- 1) Apply a set of weights – a filter – to extract **local features**
- 2) Use **multiple filters** to extract different features
- 3) **Spatially share** parameters of each filter

Convolution Filters



Original



Sharpen

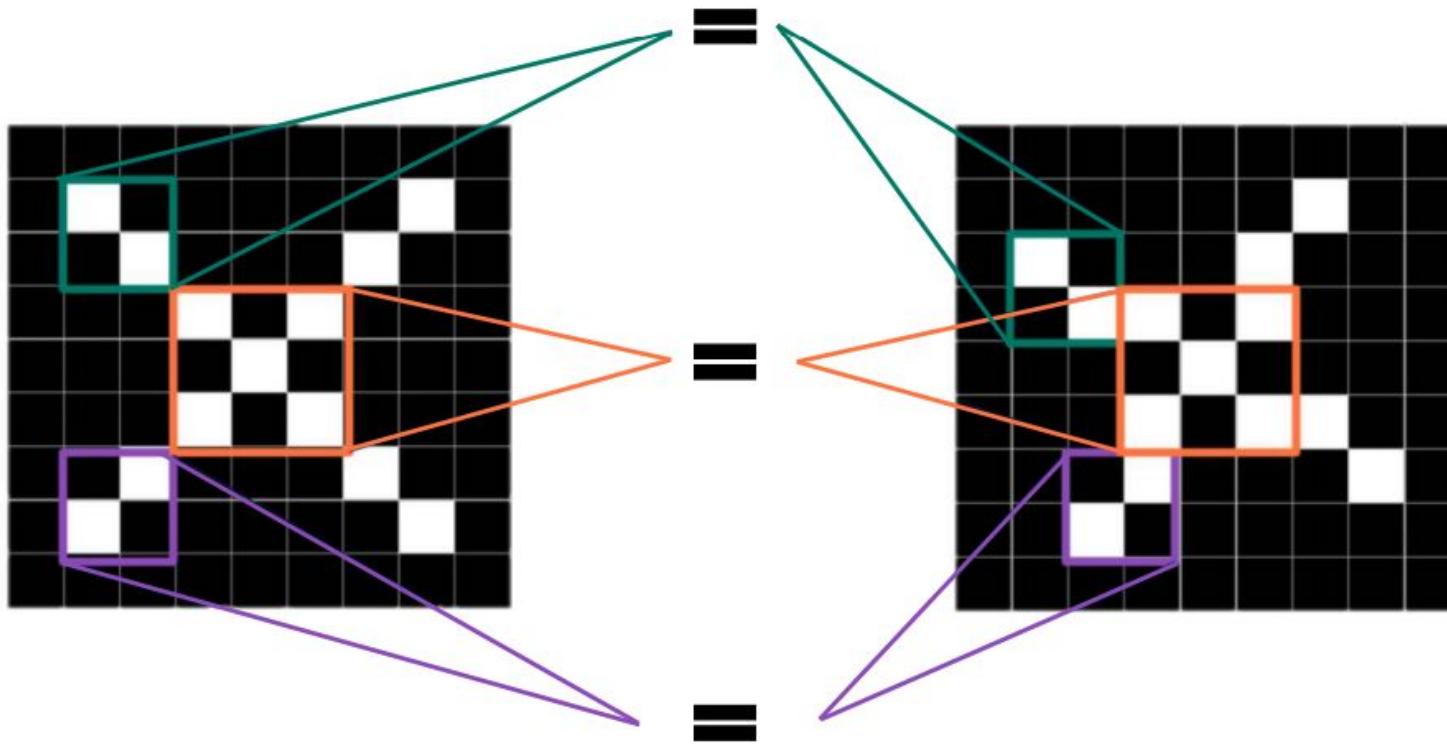


Edge Detect

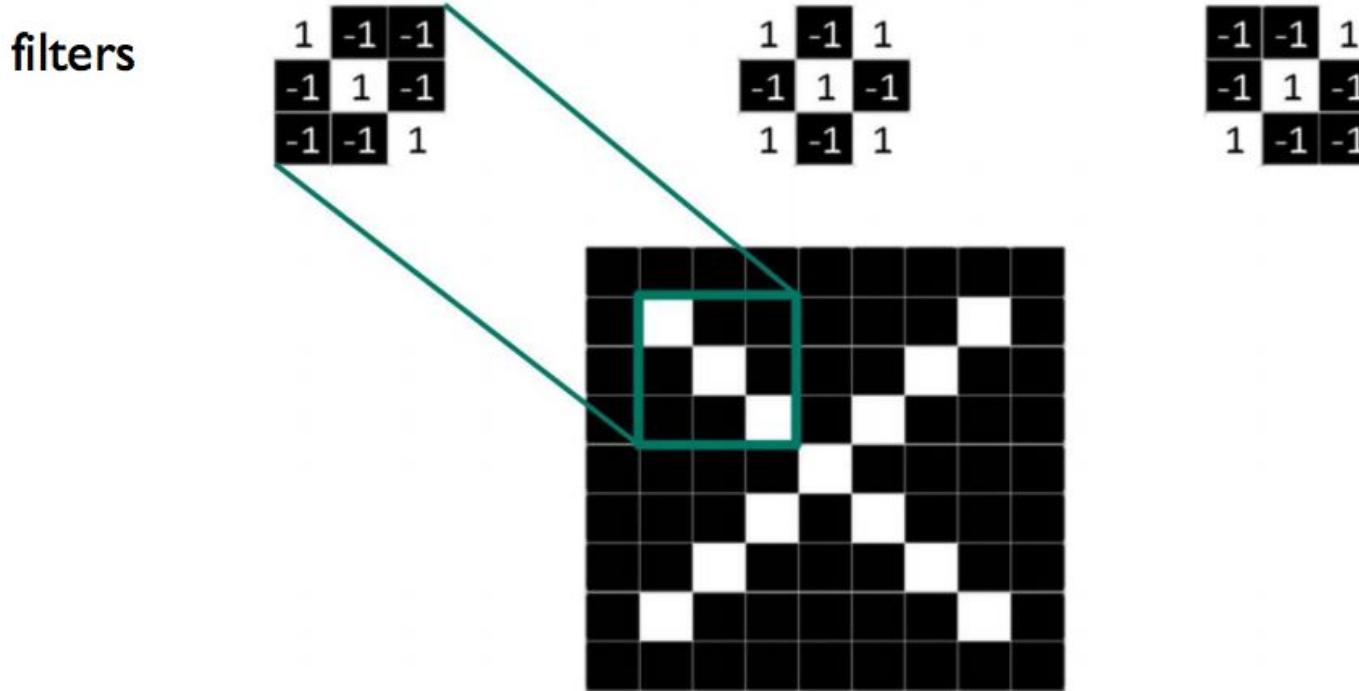


"Strong" Edge
Detect

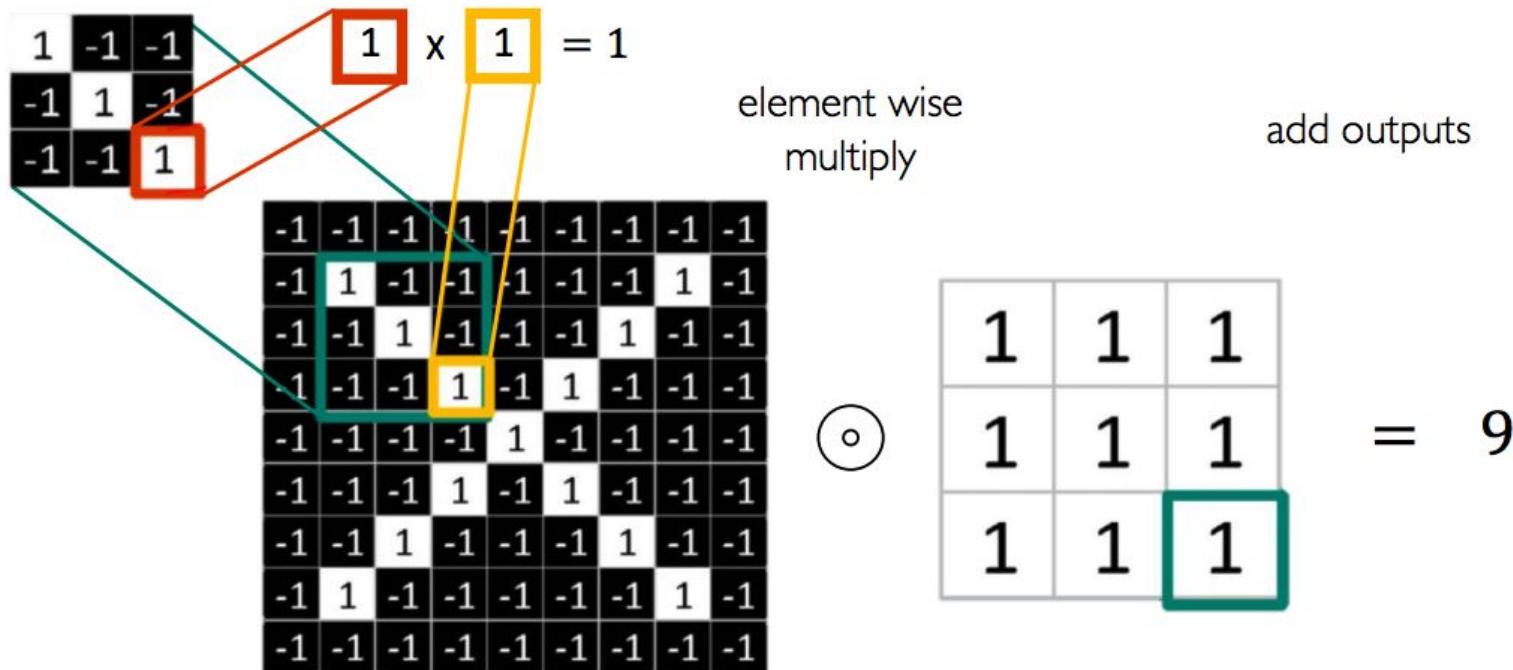
Features of X



Filters to Detect X Features

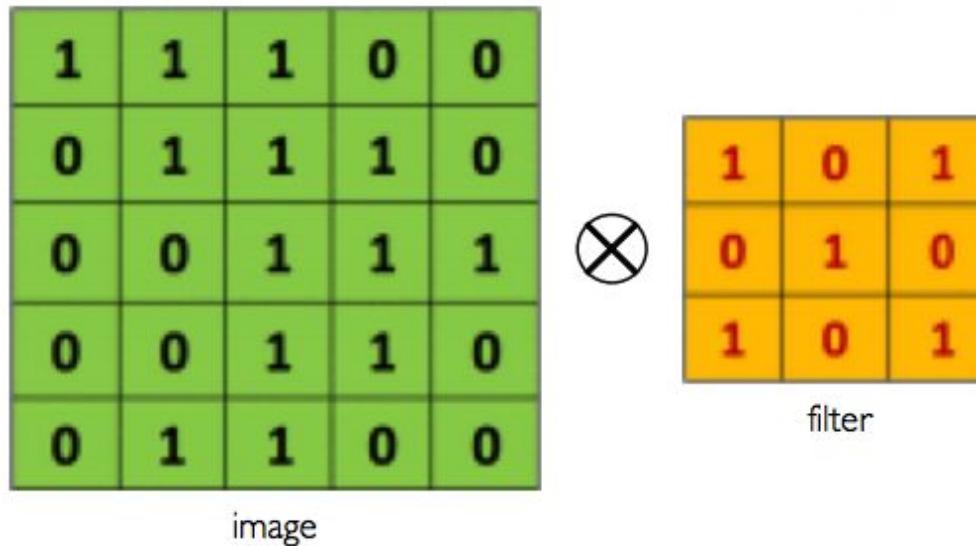


The Convolution Operation



The Convolution Operation

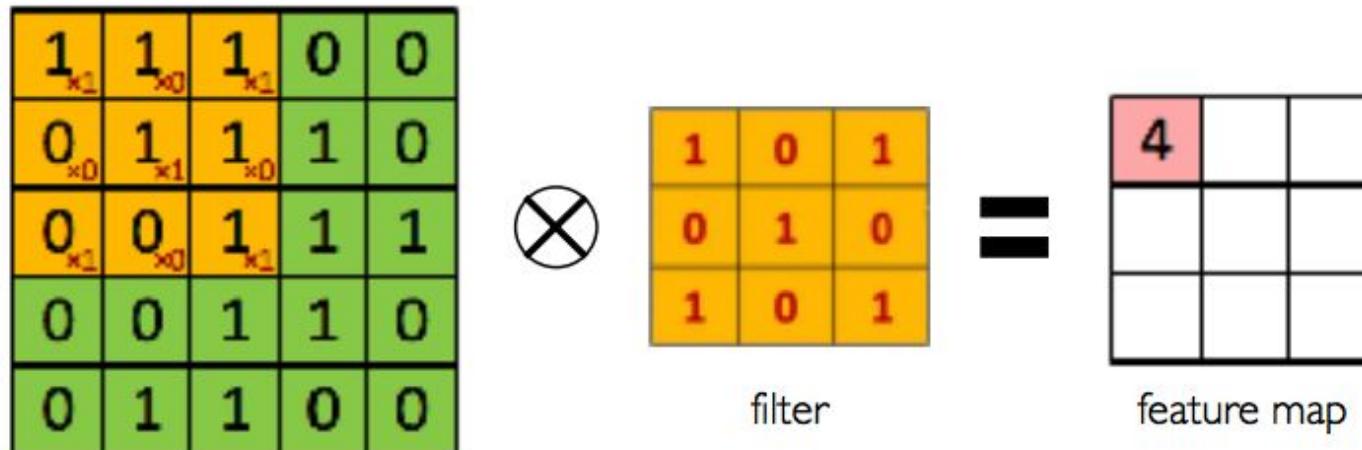
Suppose we want to compute the convolution of a 5x5 image and a 3x3 filter:



We slide the 3x3 filter over the input image, element-wise multiply, and add the outputs...

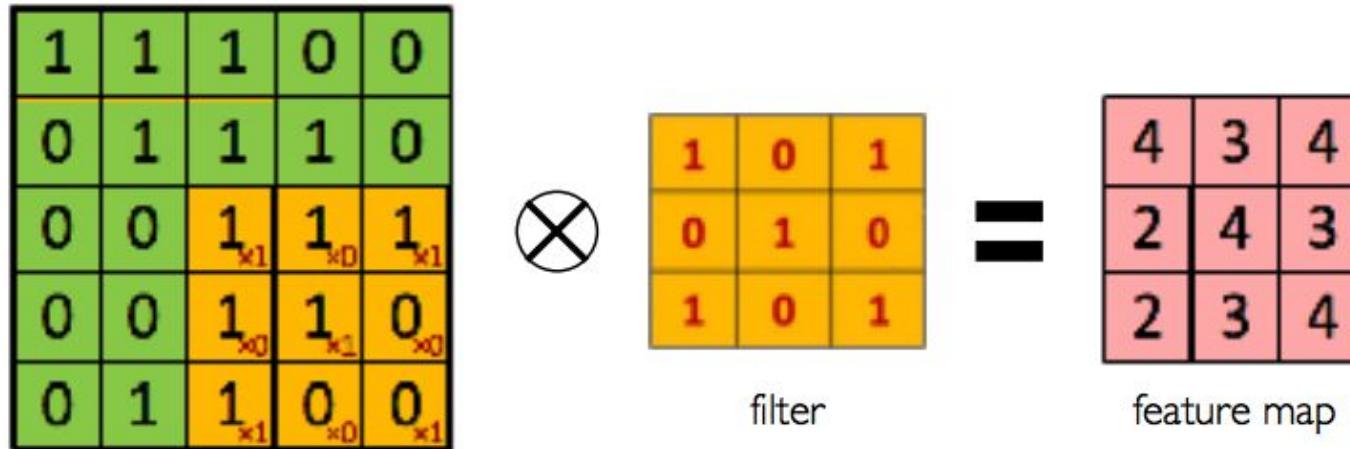
The Convolution Operation

We slide the 3x3 filter over the input image, element-wise multiply, and add the outputs:

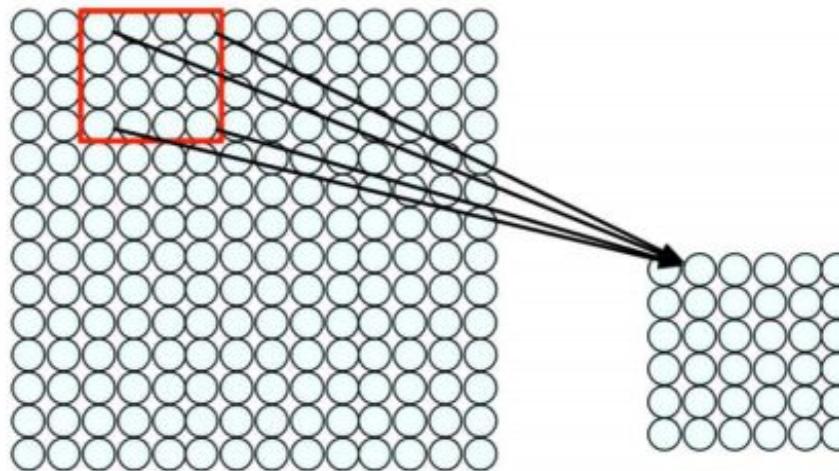


The Convolution Operation

We slide the 3x3 filter over the input image, element-wise multiply, and add the outputs:



Convolution Recap

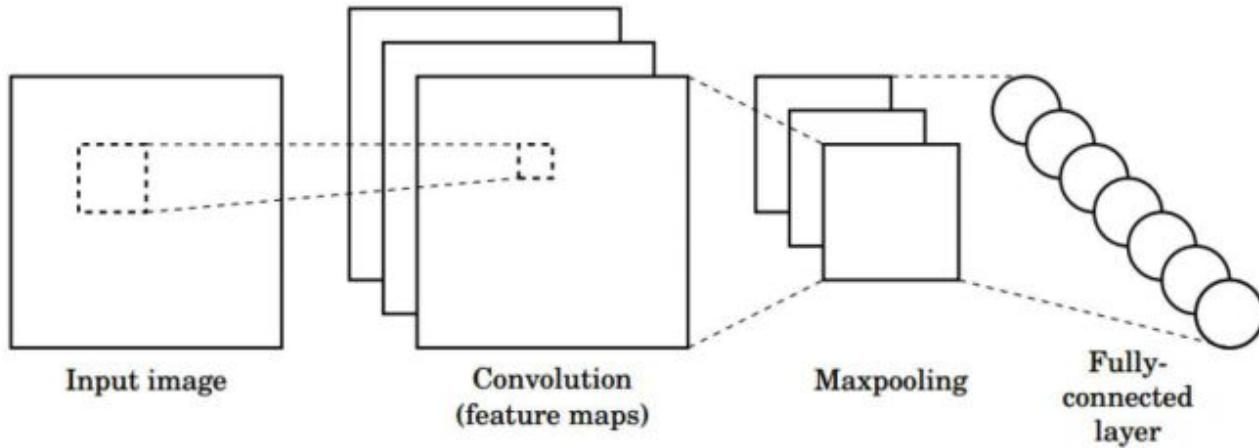


- 1) Apply a set of weights – a filter – to extract **local features**
- 2) Use **multiple filters** to extract different features
- 3) **Spatially share** parameters of each filter

Convolutional Neural Networks (CNNs)



CNNs for Classification

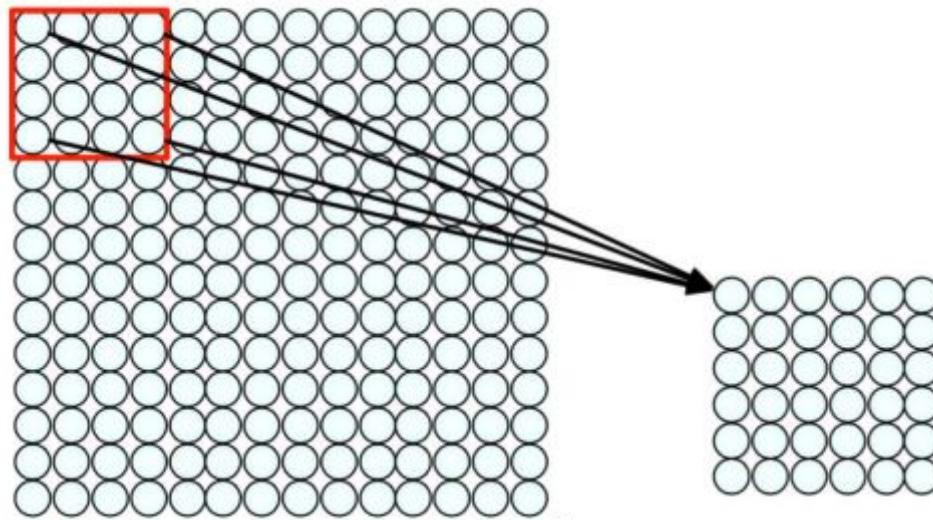


- 1. Convolution:** Apply filters with learned weights to generate feature maps.
- 2. Non-linearity:** Often ReLU.
- 3. Pooling:** Downsampling operation on each feature map.

Train model with image data.

Learn weights of filters in convolutional layers.

Convolutional Layers: Local Connectivity



$$\sum_{i=1}^4 \sum_{j=1}^4 w_{ij} x_{i+p,j+q} + b$$

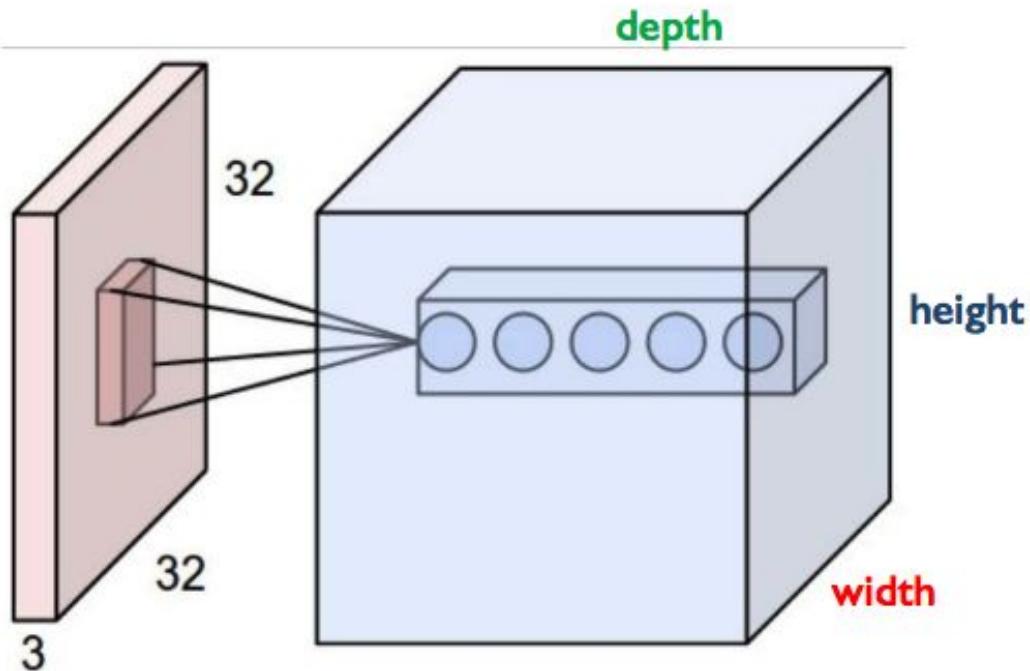
for neuron (p,q) in hidden layer

For a neuron in hidden layer:

- Take inputs from patch
- Compute weighted sum
- Apply bias

- 1) applying a window of weights
- 2) computing linear combinations
- 3) activating with non-linear function

CNNs: Spatial Arrangement of Output Volume



Layer Dimensions:

$$h \times w \times d$$

where h and w are spatial dimensions
d (depth) = number of filters

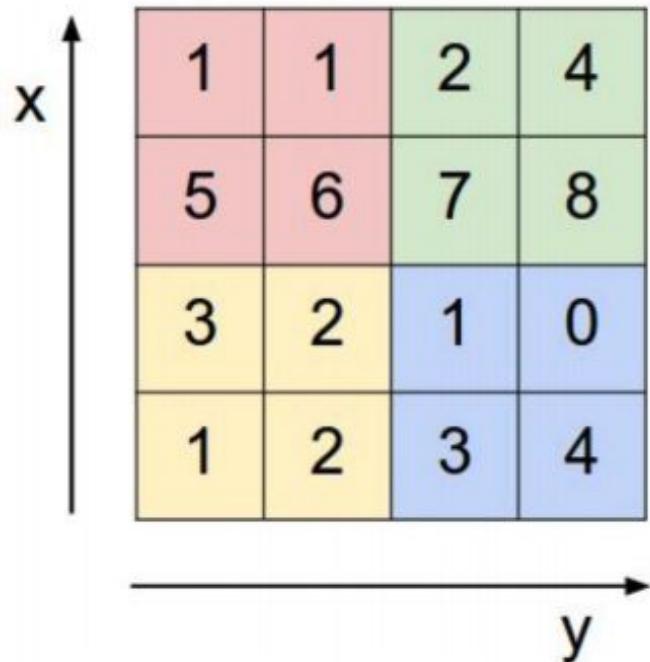
Stride:

Filter step size

Receptive Field:

Locations in input image that
a node is path connected to

Max Pooling



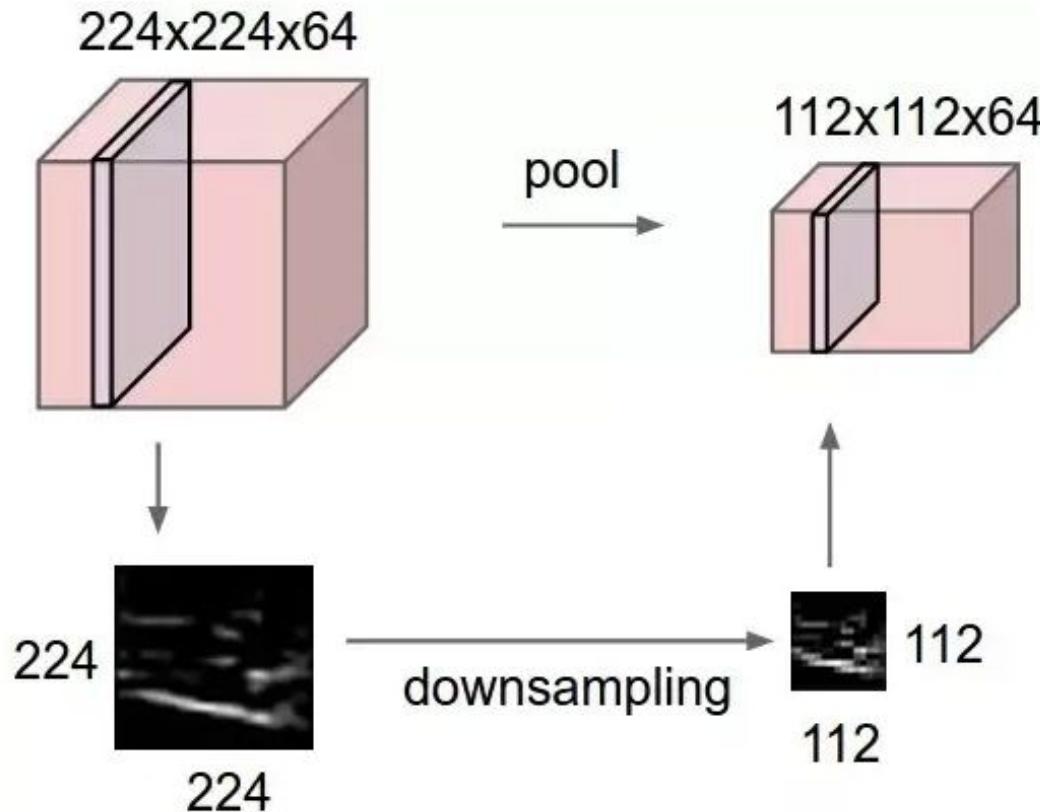
max pool with 2x2 filters
and stride 2



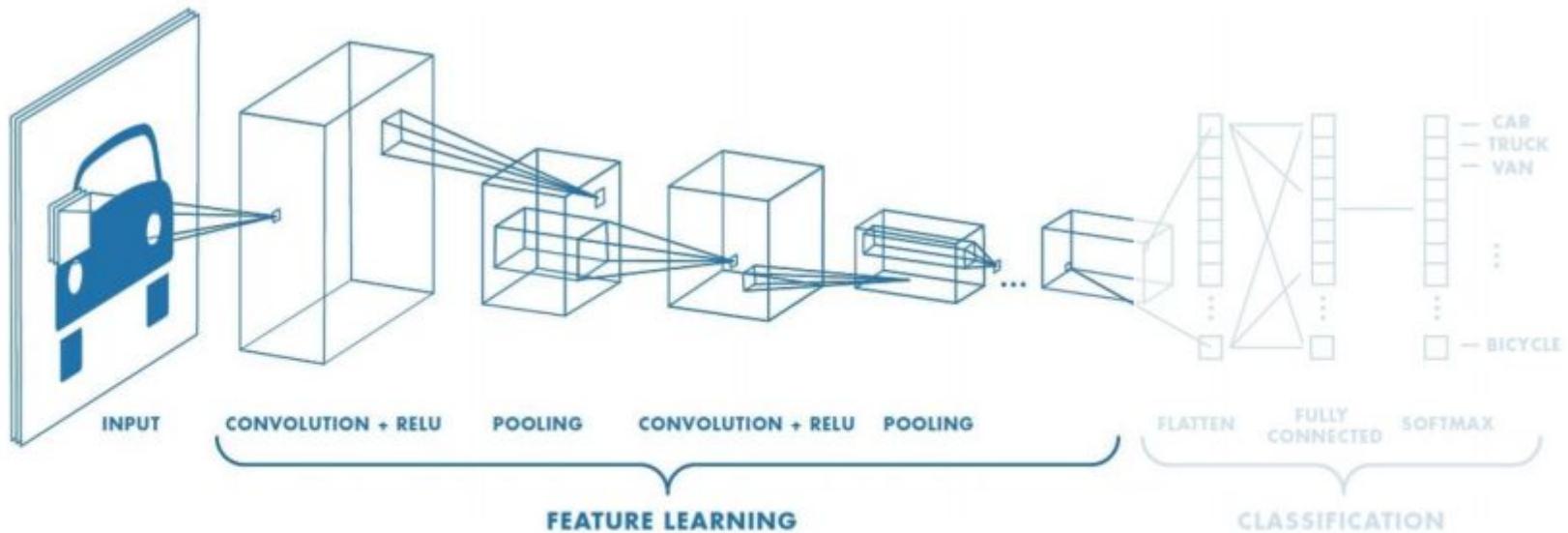
6	8
3	4

- 1) Reduced dimensionality
- 2) Spatial invariance

Max Pooling Example

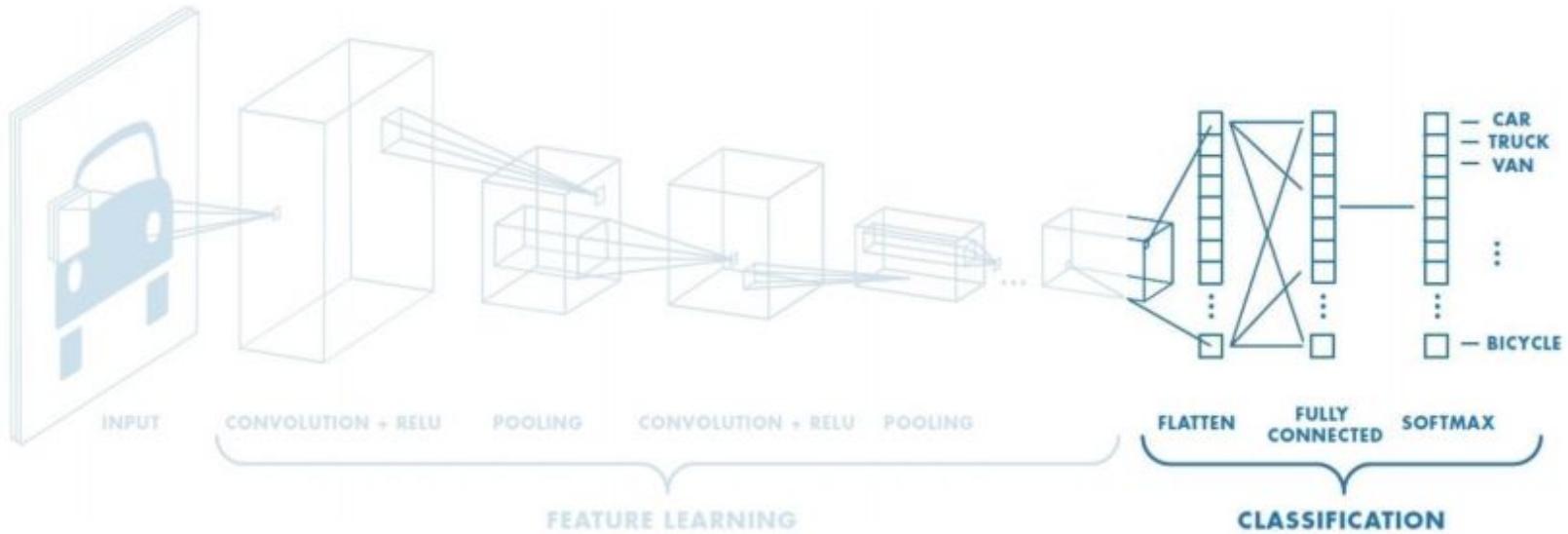


CNNs for Classification: Feature Learning



1. Learn features in input image through **convolution**
2. Introduce **non-linearity** through activation function (real-world data is non-linear!)
3. Reduce dimensionality and preserve spatial invariance with **pooling**

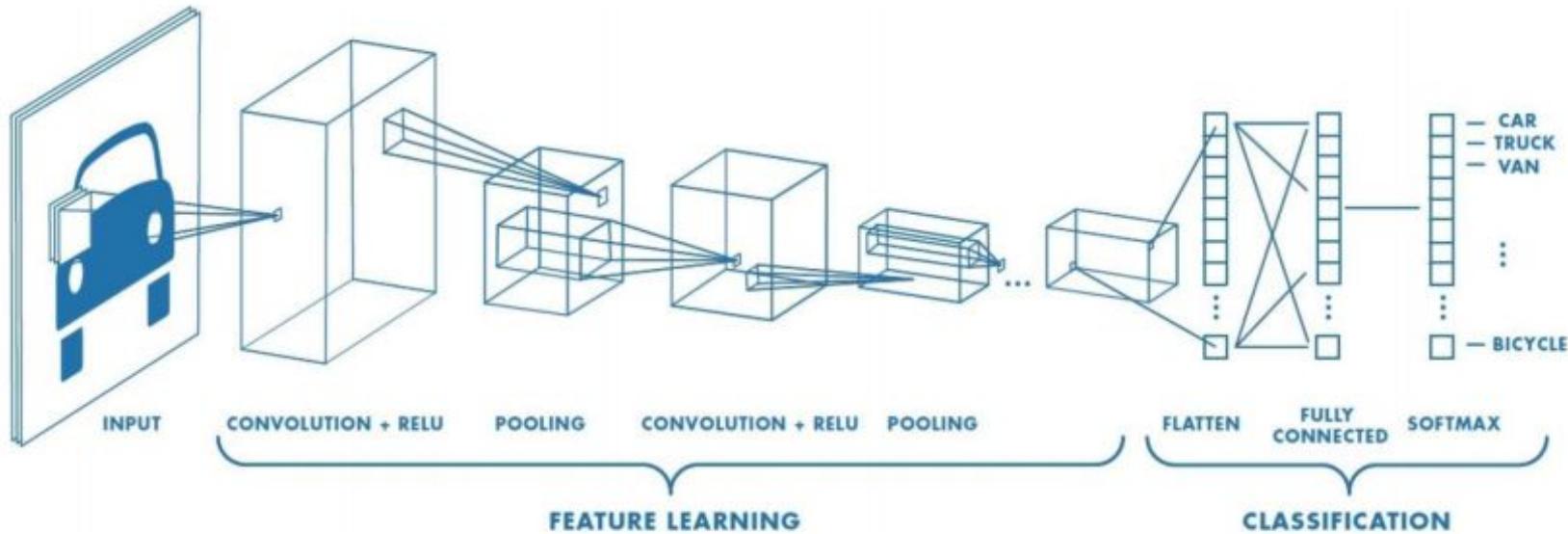
CNNs for Classification: Class Probabilities



- CONV and POOL layers output high-level features of input
- Fully connected layer uses these features for classifying input image
- Express output as **probability** of image belonging to a particular class

$$\text{softmax}(y_i) = \frac{e^{y_i}}{\sum_j e^{y_j}}$$

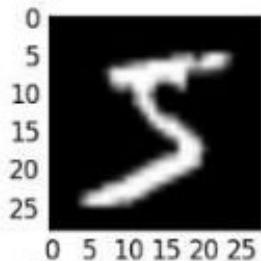
CNNs: Training with Backpropagation



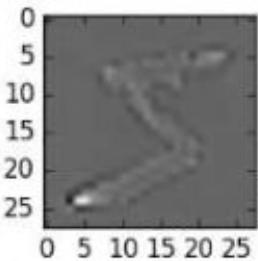
Learn weights for convolutional filters and fully connected layers

Backpropagation: cross-entropy loss

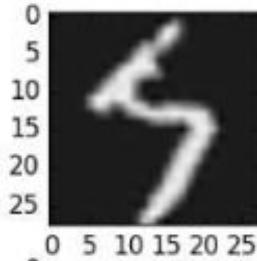
Data Augmentation



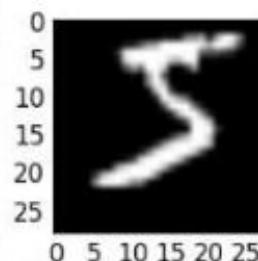
Example MNIST
images



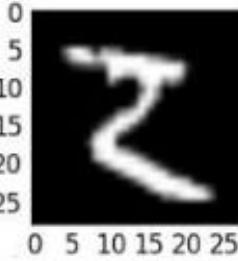
ZCA Whitening



Random Rotations



Random shift



Random flip

- > Approaches that alter the training data in ways that change the array representation while keeping the label the same.
- > They are a way to artificially expand your dataset. Some popular augmentations people use are gray scales , horizontal flips, vertical flips, random crops, color jitters, translations, rotations, and much more .

ImageNet Dataset

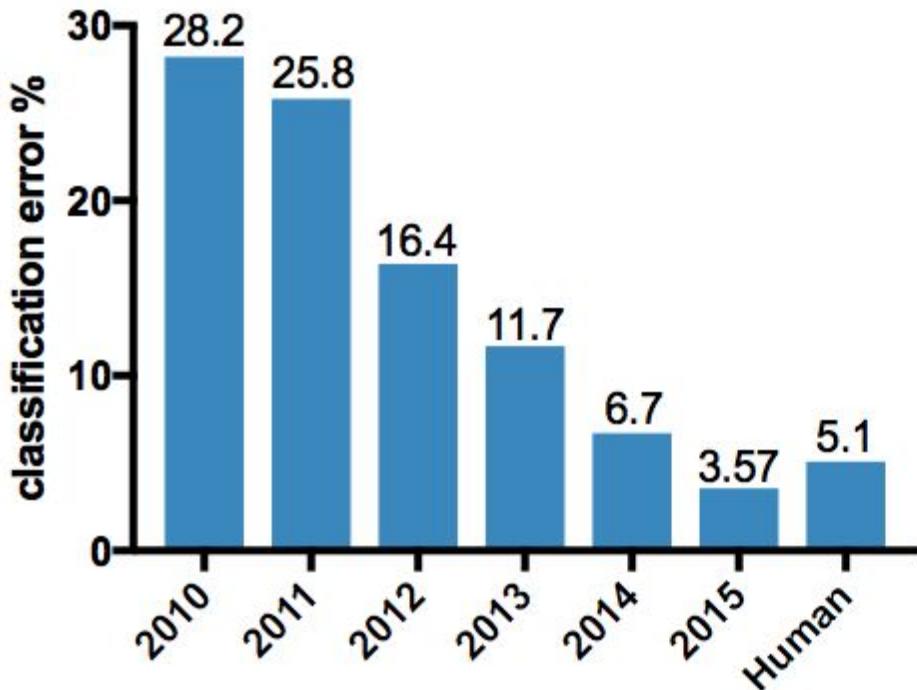
Dataset of over 14 million images across 21,841 categories

"Elongated crescent-shaped yellow fruit with soft sweet flesh"



1409 pictures of bananas.

ImageNet Challenge: Classification Task



2012: AlexNet. First CNN to win.

- 8 layers, 61 million parameters

2013: ZFNet

- 8 layers, more filters

2014: VGG

- 19 layers

2014: GoogLeNet

- "Inception" modules
- 22 layers, 5million parameters

2015: ResNet

- 152 layers

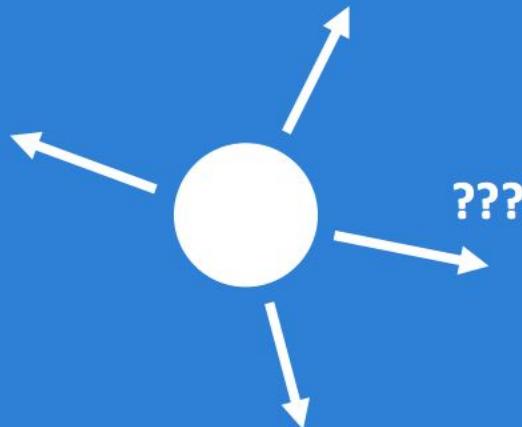
Deep Sequence Modeling



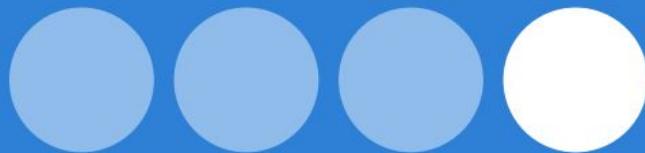
Given an image of a ball,
can you predict where it will go next?



Given an image of a ball,
can you predict where it will go next?



Given an image of a ball,
can you predict where it will go next?



Given an image of a ball,
can you predict where it will go next?



Sequence modeling: predict the next word

“This morning I took my cat for a walk.”

given these words

Sequence modeling: predict the next word

“This morning I took my cat for a walk.”

given these words

predict the
next word

Idea #1: use a fixed window

“This morning I took my cat for a walk.”

given these predict the
two words next word

One-hot feature encoding: tells us what each word is

[1 0 0 0 0 0 1 0 0 0]

for a



prediction

Problem: can't model long-term dependencies

"**France** is where I grew up, but I now live in Boston. I speak fluent ____."

- A) English
- B) French
- C) Spanish
- D) Chinese

We need information from **the distant past** to accurately predict the correct word.

Idea #2: use entire sequence as set of counts

“This morning I took my cat for a”



“bag of words”

[0 1 0 0 1 0 0 ... 0 0 1 1 0 0 0 1]



prediction

Problem: counts don't preserve order



The food was good, not bad at all.

vs.

The food was bad, not good at all.



Idea #3: use a really big fixed window

“This morning I took my cat for a walk.”

given these
words

predict the
next word

[1 0 0 0 0 0 0 0 1 0 0 1 0 0 0 1 0 0 0 0 0 1 0 ...]

morning

|

took

this

cat



prediction

Problem #3: no parameter sharing

[1 0 0 0 0 0 0 0 1 0 0 1 0 0 1 0 0 0 0 0 1 0 ...]
this morning took the cat

Each of these inputs has a **separate parameter**:

[0 0 0 1 0 0 0 1 0 0 0 1 0 0 0 1 0 0 0 0 0 1 ...]
this morning

Problem #3: no parameter sharing

[1 0 0 0 0 0 0 0 1 0 0 1 0 0 1 0 0 0 0 0 1 0 ...]
this morning took the cat

Each of these inputs has a **separate parameter**:

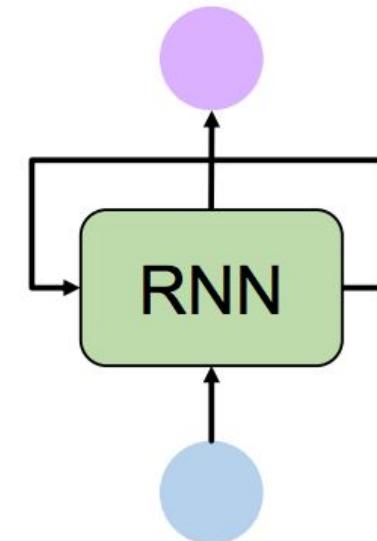
[0 0 0 1 0 0 0 1 0 0 0 1 0 0 0 1 0 0 0 0 0 0 1 ...]
this morning

Things we learn about the sequence **won't transfer** if
they appear **elsewhere** in the sequence.

Sequence modeling: design criteria

To model sequences, we need to:

1. Handle **variable-length** sequences
2. Track **long-term** dependencies
3. Maintain information about **order**
4. **Share parameters** across the sequence



Today: **Recurrent Neural Networks (RNNs)** as
an approach to sequence modeling problems

Recurrent Neural Networks (RNNs)

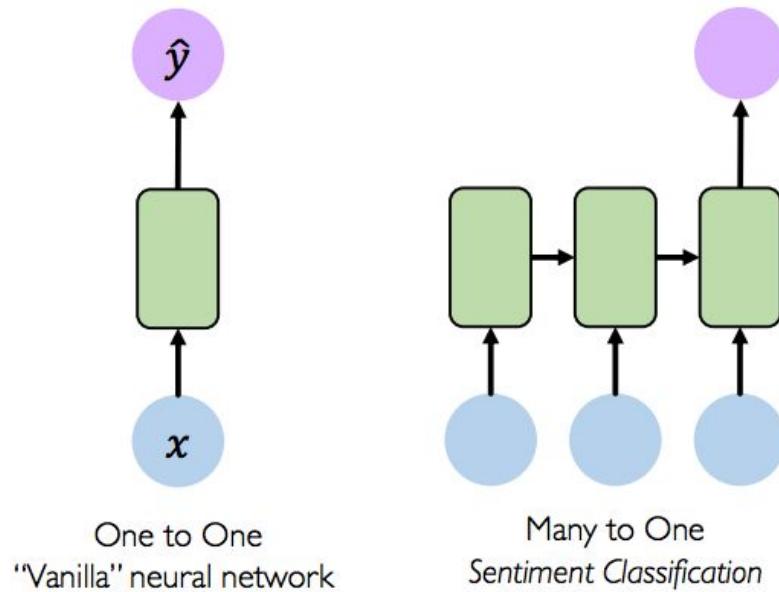


Standard feed-forward neural network



One to One
"Vanilla" neural network

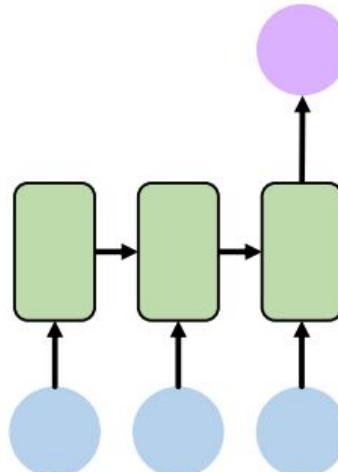
Recurrent neural networks: sequence modeling



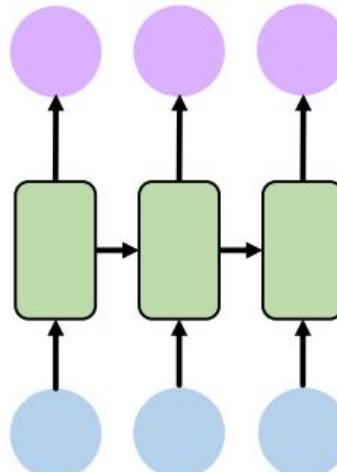
Recurrent neural networks: sequence modeling



One to One
"Vanilla" neural network



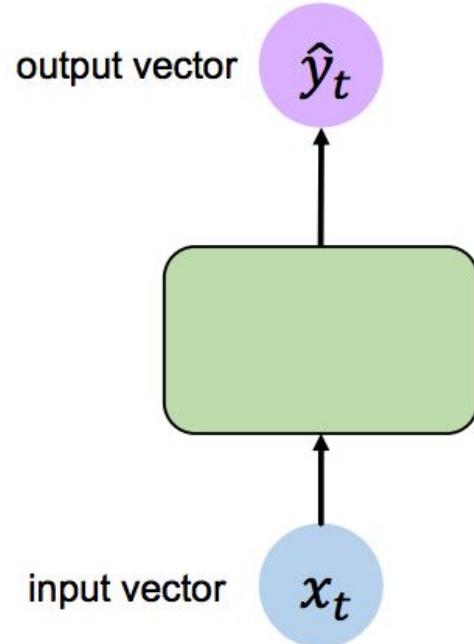
Many to One
Sentiment Classification



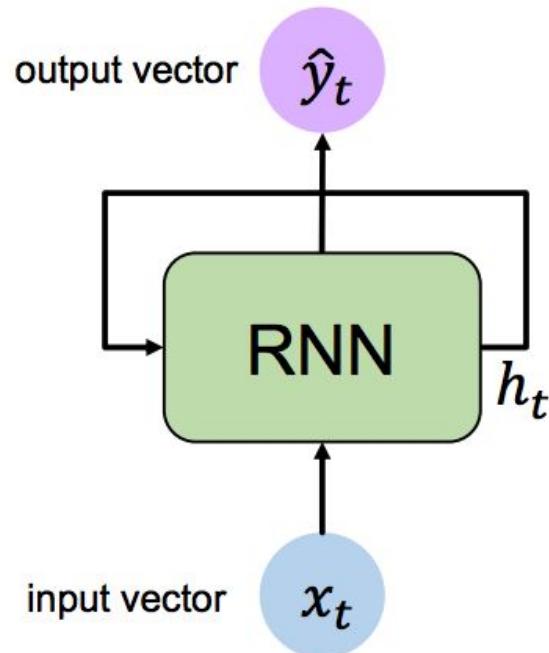
Many to Many
Music Generation

... and many other
architectures and
applications

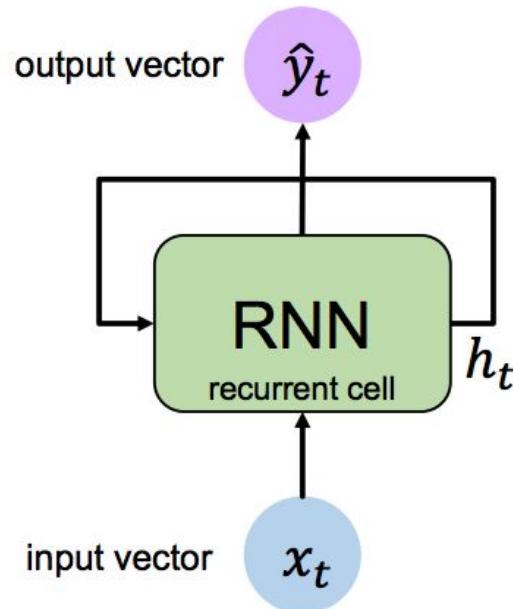
A standard “vanilla” neural network



A recurrent neural network (RNN)



A recurrent neural network (RNN)



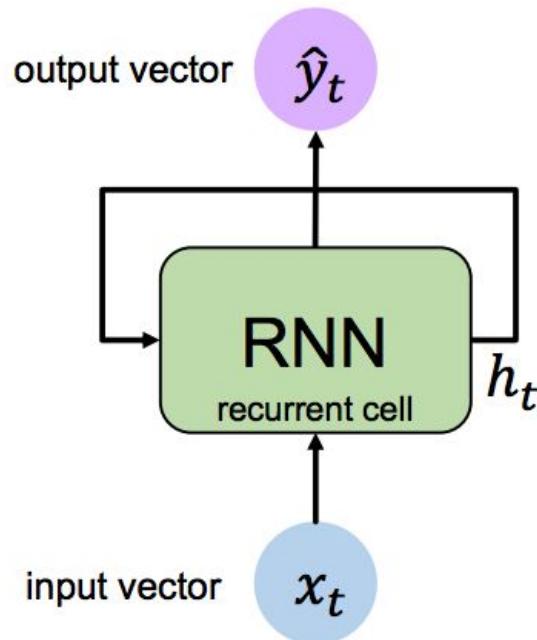
Apply a **recurrence relation** at every time step to process a sequence:

$$h_t = f_W(h_{t-1}, x_t)$$

new state function
 parameterized
 by W old state input vector at
 time step t

Note: the same function and set of parameters are used at every time step

RNN state update and output

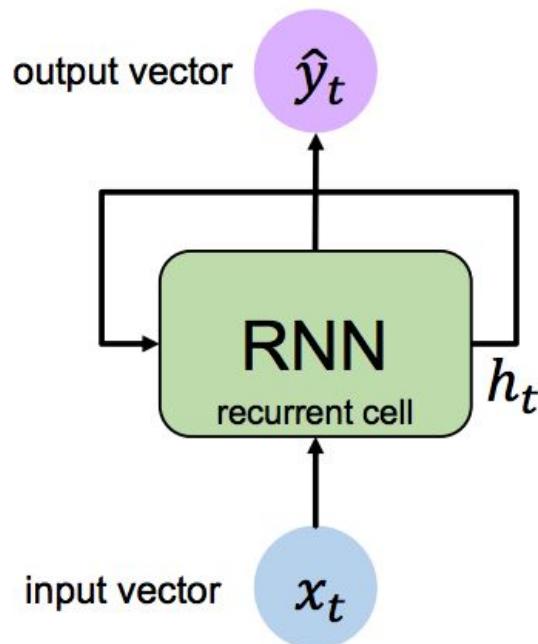


Update Hidden State

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

Input Vector

RNN state update and output



Output Vector

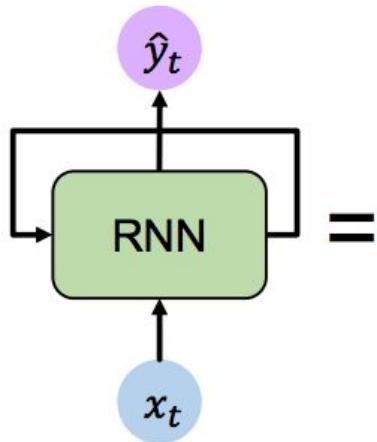
$$\hat{y}_t = \mathbf{W}_{hy} h_t$$

Update Hidden State

$$h_t = \tanh(\mathbf{W}_{hh} h_{t-1} + \mathbf{W}_{xh} x_t)$$

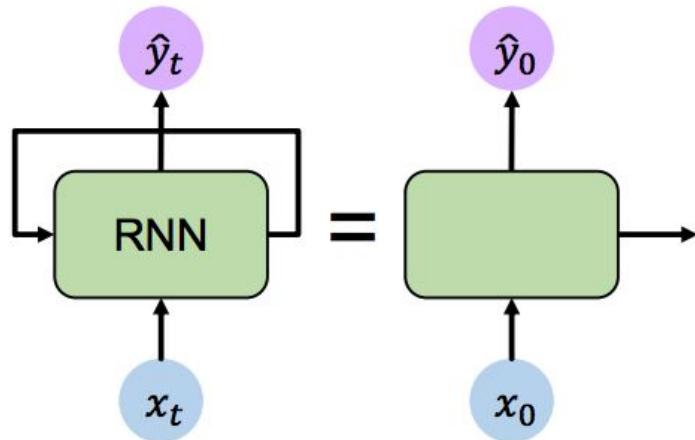
Input Vector

RNNs: computational graph across time

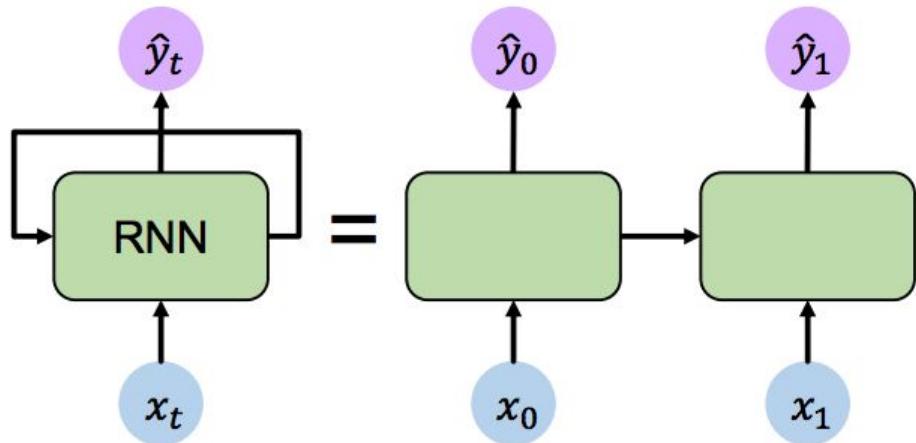


Represent as computational graph unrolled across time

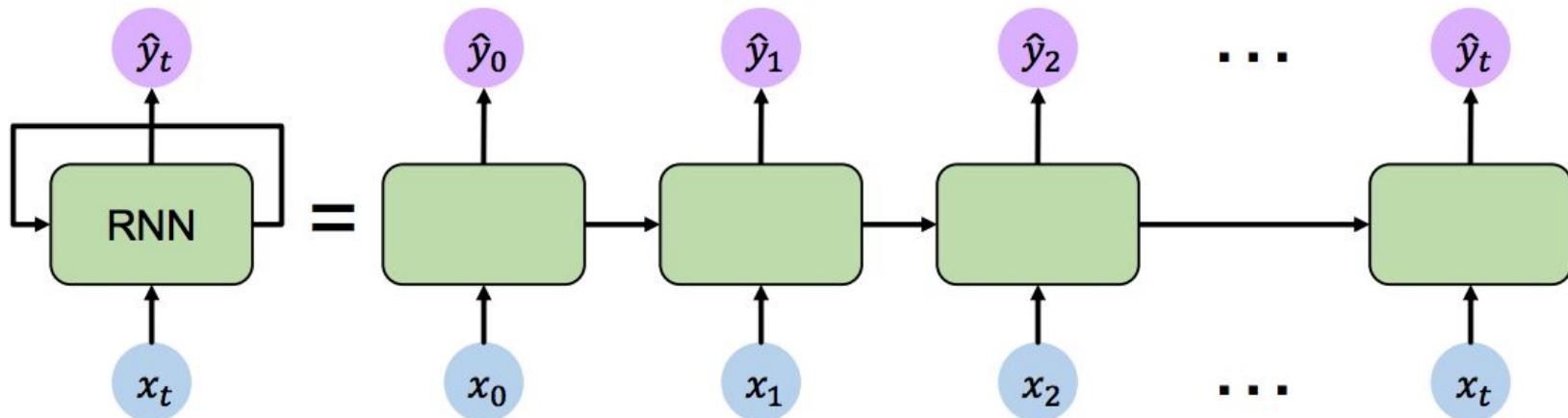
RNNs: computational graph across time



RNNs: computational graph across time

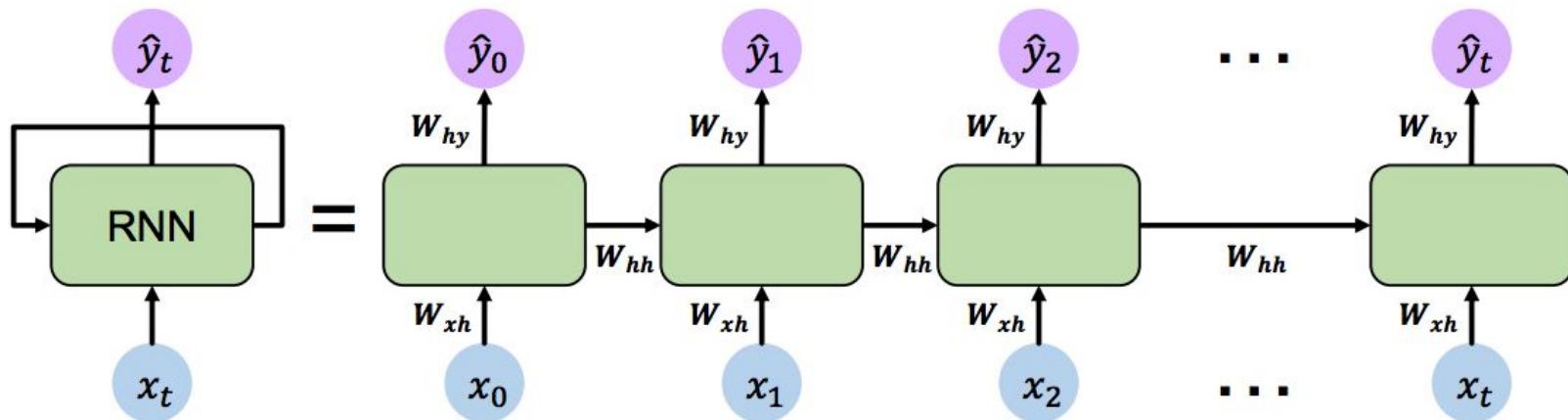


RNNs: computational graph across time

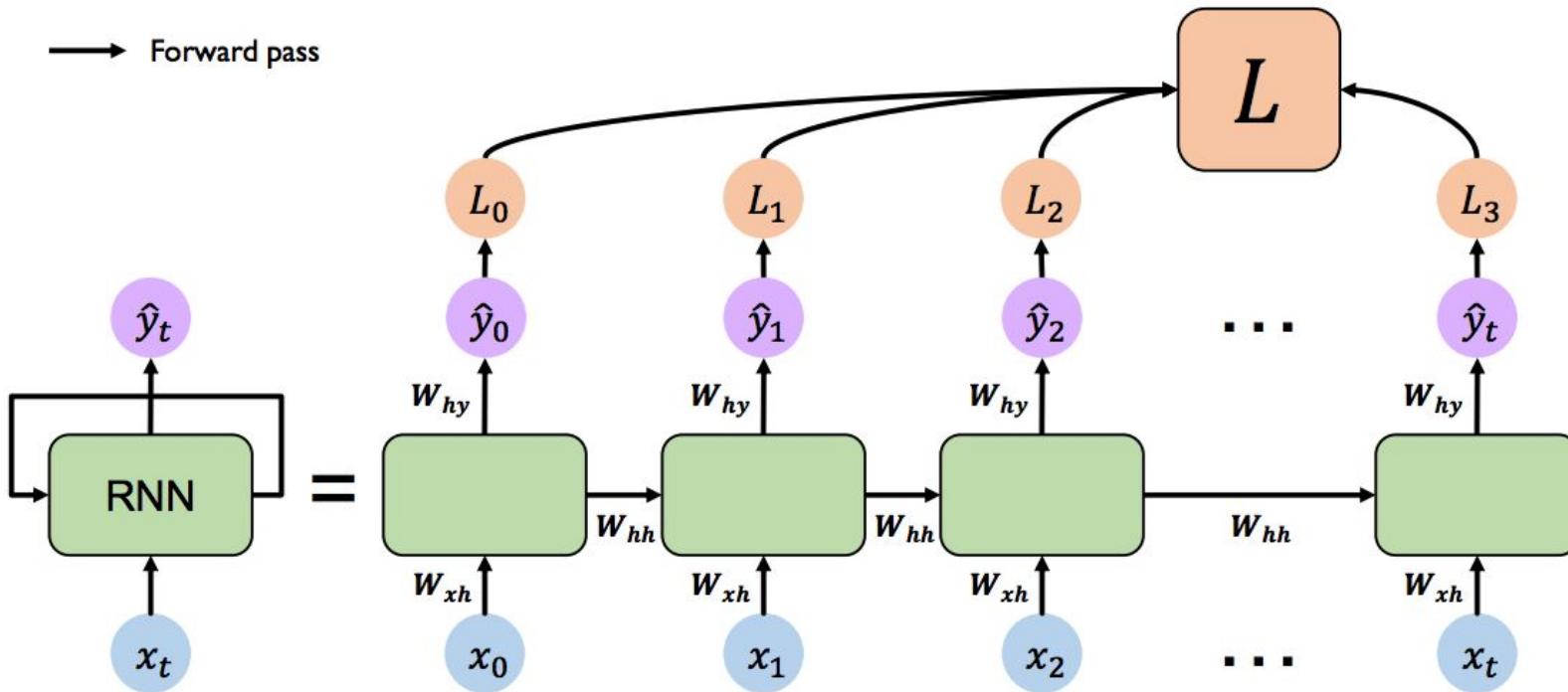


RNNs: computational graph across time

Re-use the **same weight matrices** at every time step

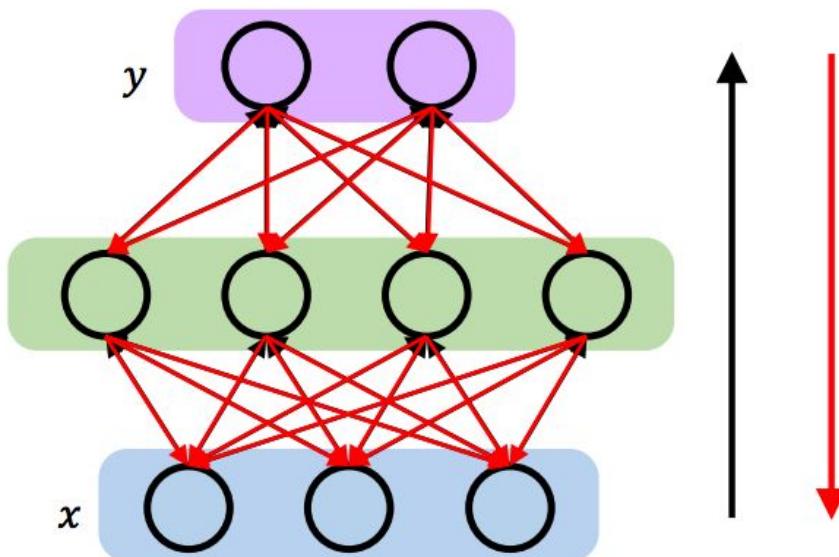


RNNs: computational graph across time



Training RNN with Backpropagation

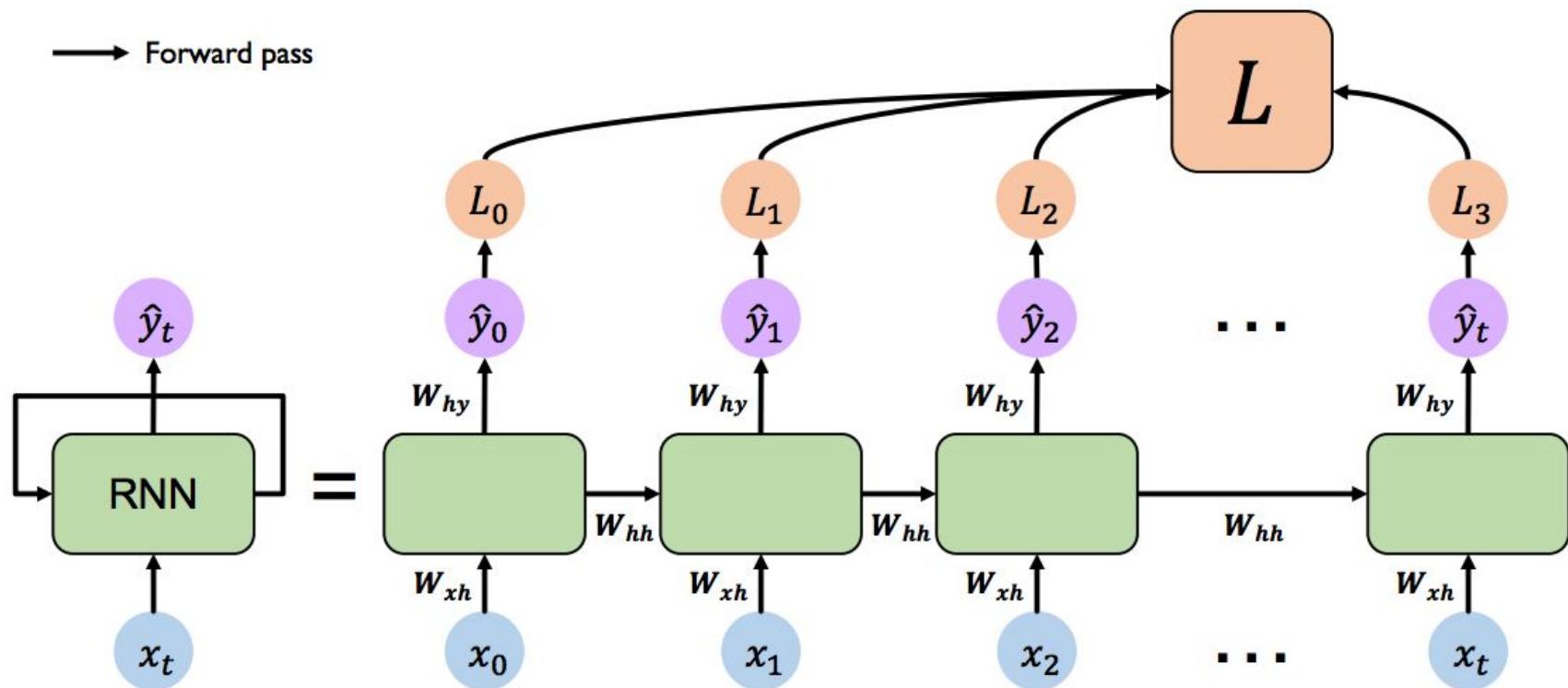
Backpropagation in vanilla neural network



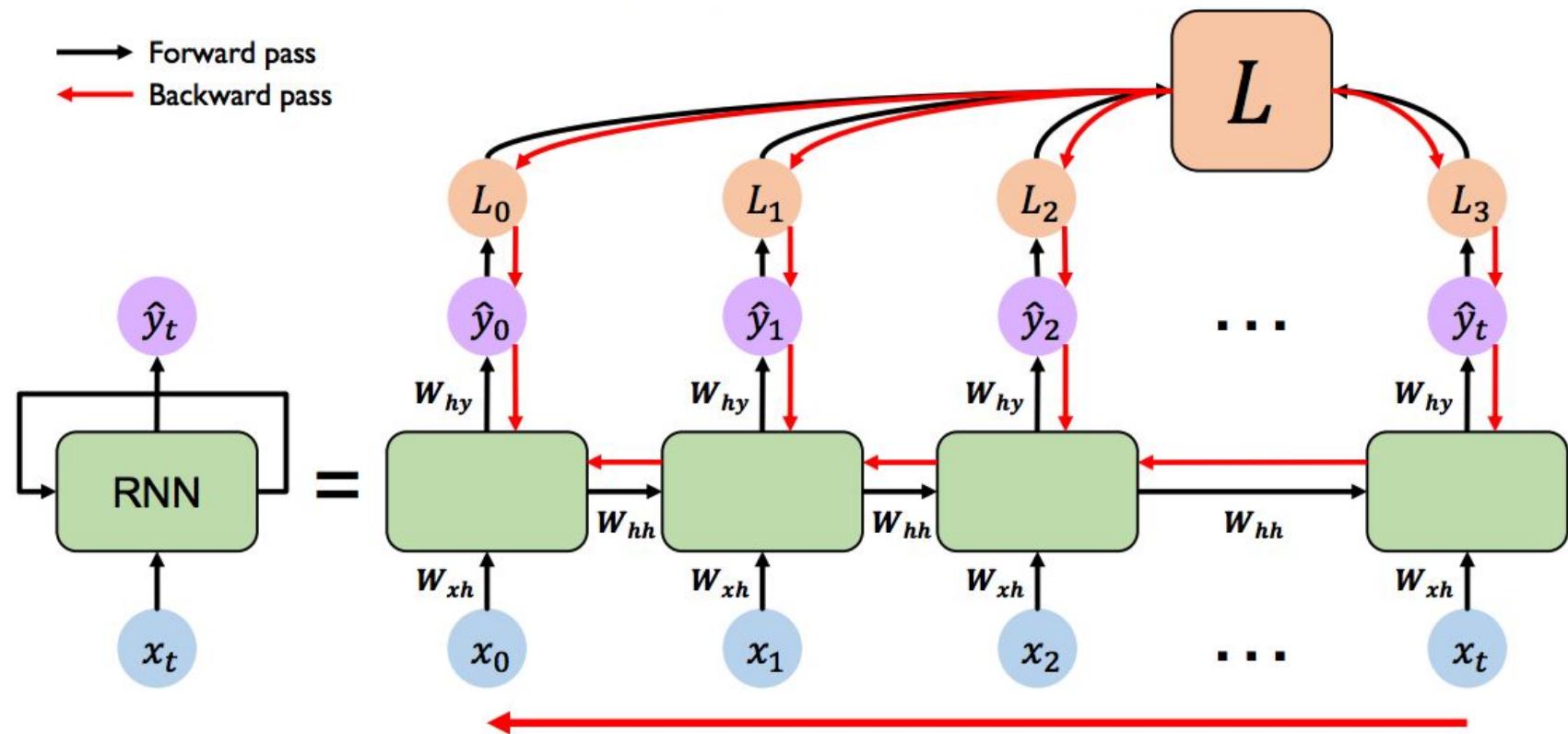
Backpropagation algorithm:

1. Take the derivative (gradient) of the loss with respect to each parameter
2. Shift parameters in order to minimize loss

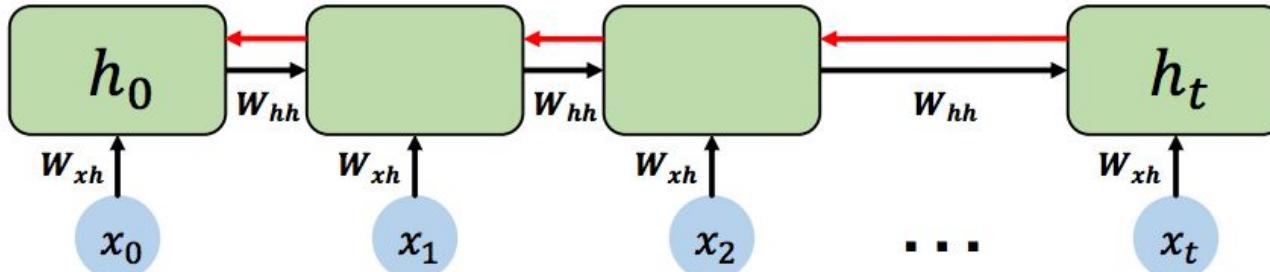
RNNs: backpropagation through time



RNNs: backpropagation through time

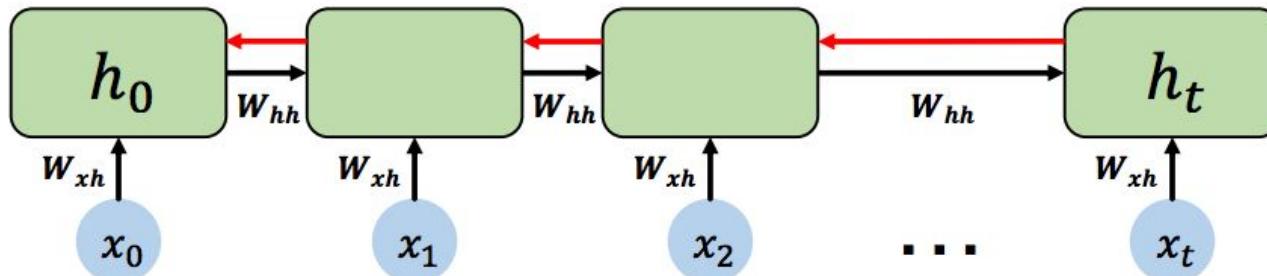


Standard RNN gradient flow



Computing the gradient wrt h_0 involves **many factors of W_{hh}** (and repeated f' !)

RNN gradient flow: exploding gradients

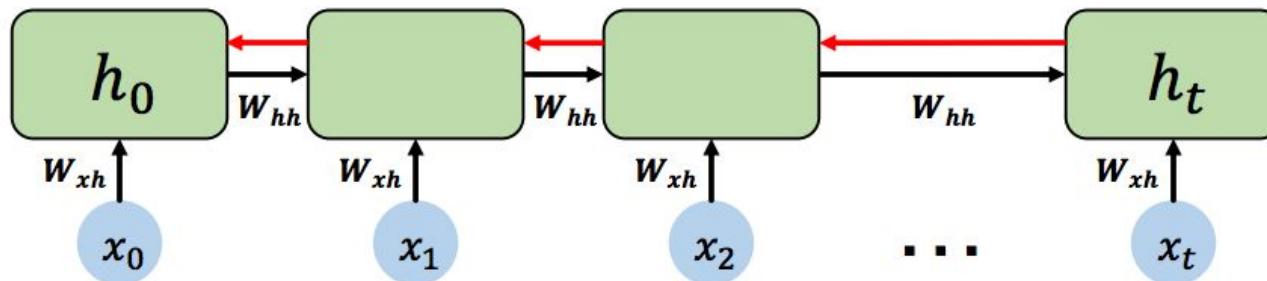


Computing the gradient wrt h_0 involves **many factors of W_{hh}** (and repeated f' !)

Many values > 1 :
exploding gradients

Gradient clipping to
scale big gradients

RNN gradient flow: vanishing gradients



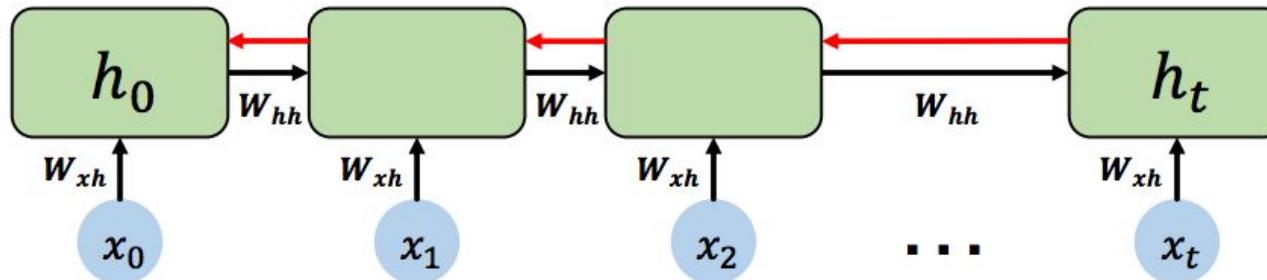
Computing the gradient wrt h_0 involves **many factors of W_{hh}** (and repeated f' !)

Many values > 1 :
exploding gradients

Gradient clipping to
scale big gradients

Many values < 1 :
vanishing gradients

RNN gradient flow: vanishing gradients



Computing the gradient wrt h_0 involves **many factors of W_{hh}** (and repeated f' !)

Largest singular value > 1 :
exploding gradients

Gradient clipping to
scale big gradients

Largest singular value < 1 :
vanishing gradients

1. Activation function
2. Weight initialization
3. Network architecture

The problem of long-term dependencies

Why are vanishing gradients a problem?

Multiply many **small numbers** together



Errors due to further back time steps
have smaller and smaller gradients



Bias network to capture short-term
dependencies

The problem of long-term dependencies

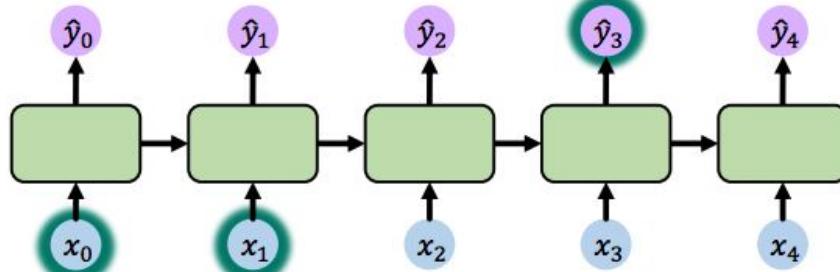
Why are vanishing gradients a problem?

Multiply many **small numbers** together

↓
Errors due to further back time steps
have smaller and smaller gradients

↓
Bias parameters to capture short-term
dependencies

"The clouds are in the ___"



The problem of long-term dependencies

Why are vanishing gradients a problem?

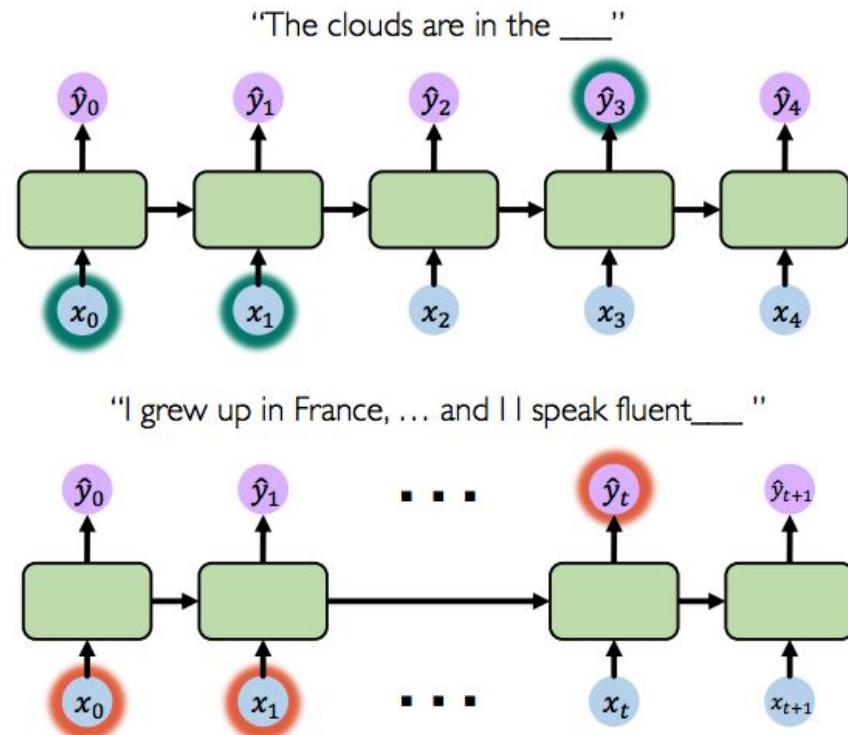
Multiply many **small numbers** together



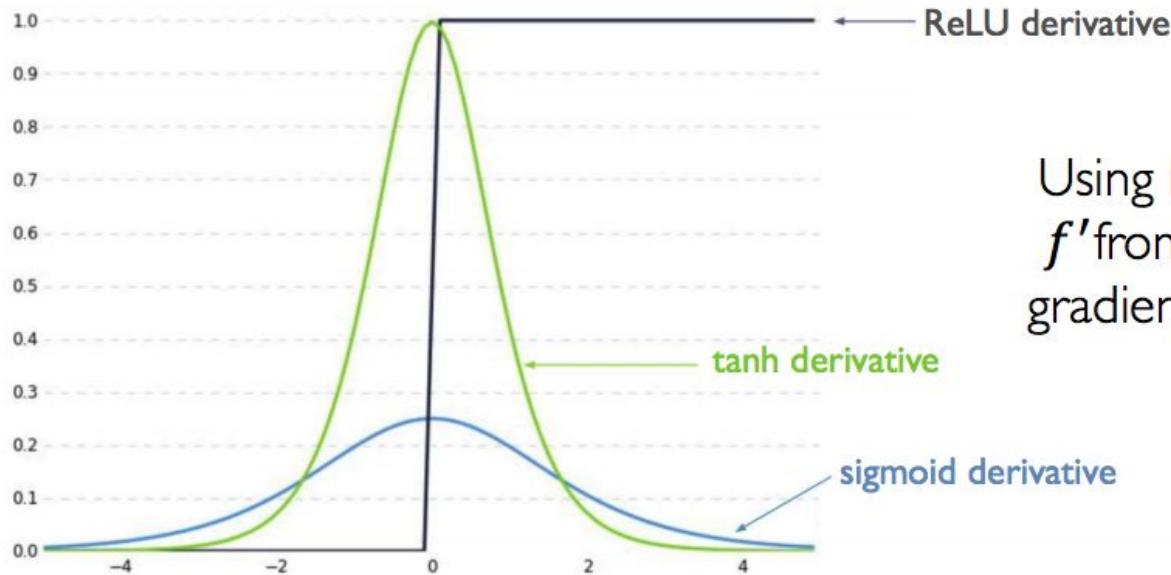
Errors due to further back time steps
have smaller and smaller gradients



Bias parameters to capture short-term
dependencies



Trick #1: activation functions



Using ReLU prevents
 f' from shrinking the
gradients when $x > 0$

Trick #2: parameter initialization

Initialize **weights** to identity matrix

Initialize **biases** to zero

$$I_n = \begin{pmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{pmatrix}$$

This helps prevent the weights from shrinking to zero.

Solution #3: gated cells

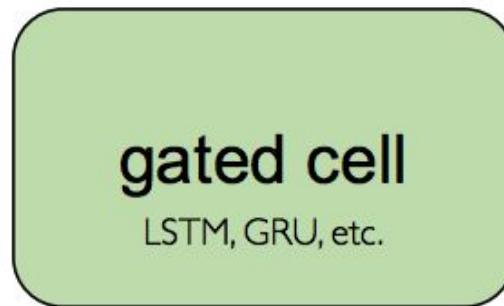
Idea: use a more **complex recurrent unit with gates** to control what information is passed through

gated cell

LSTM, GRU, etc.

Solution #3: gated cells

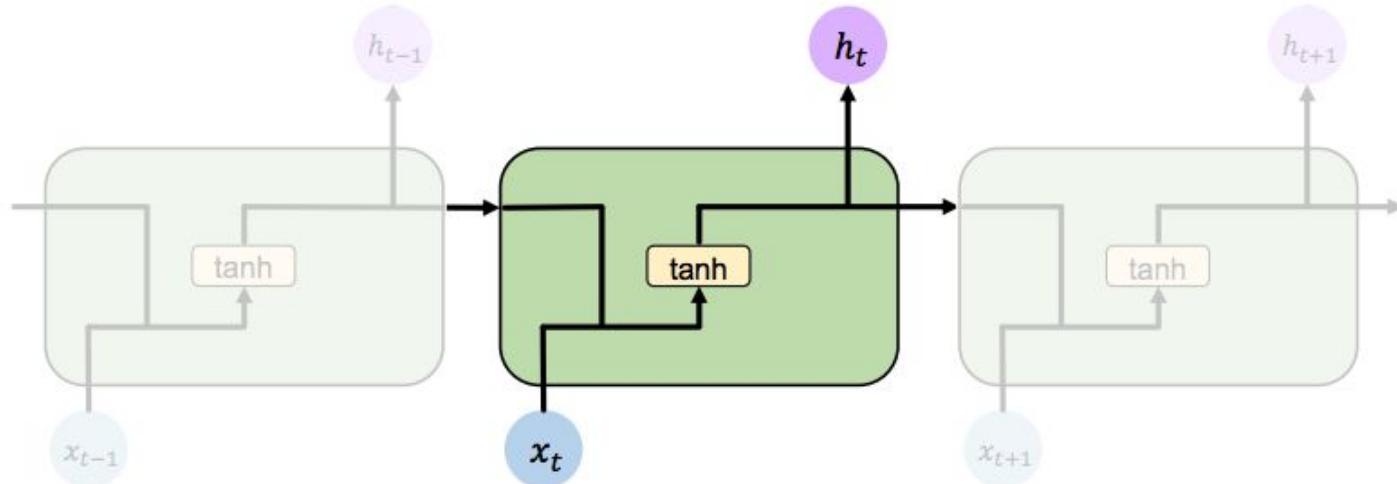
Idea: use a more **complex recurrent unit with gates** to control what information is passed through



Long Short Term Memory (LSTMs) networks rely on a gated cell to track information throughout many time steps.

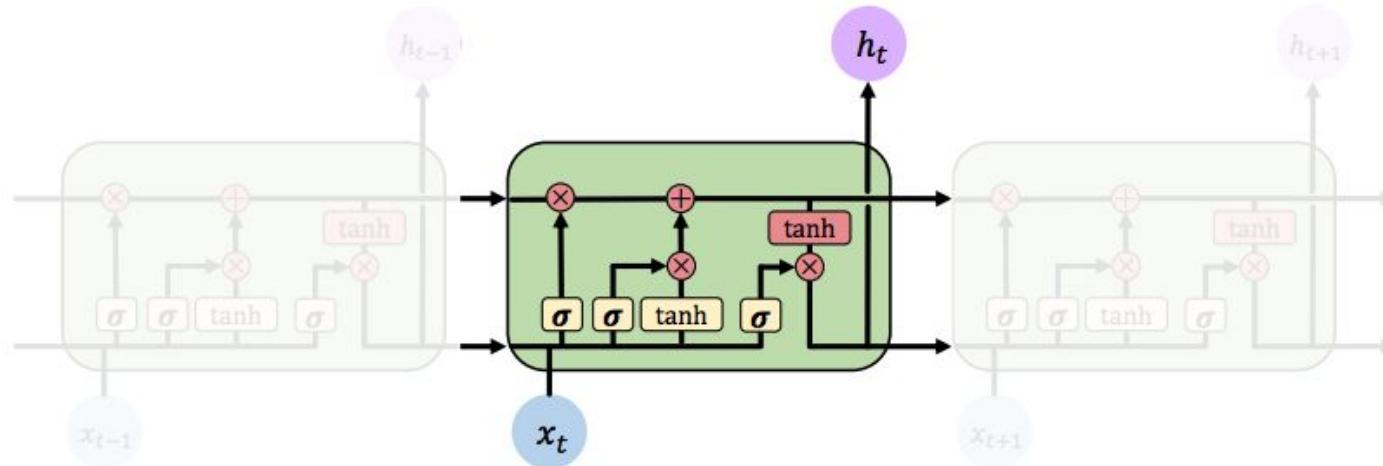
Standard RNN

In a standard RNN, repeating modules contain a **simple computation node**



Long Short Term Memory (LSTMs)

LSTM repeating modules contain **interacting layers** that **control information flow**

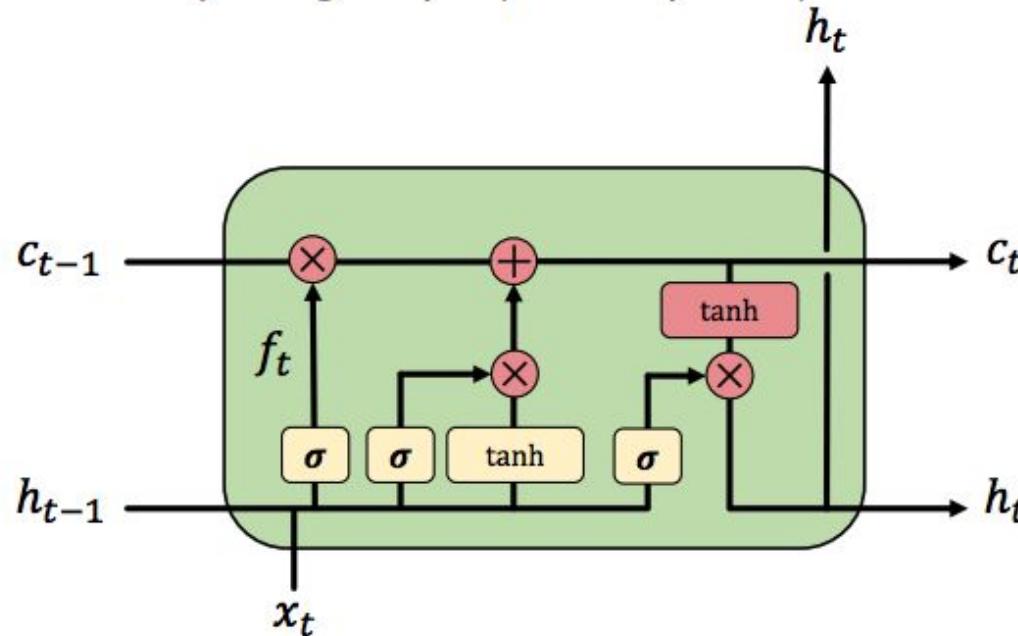


LSTM cells are able to track information throughout many timesteps

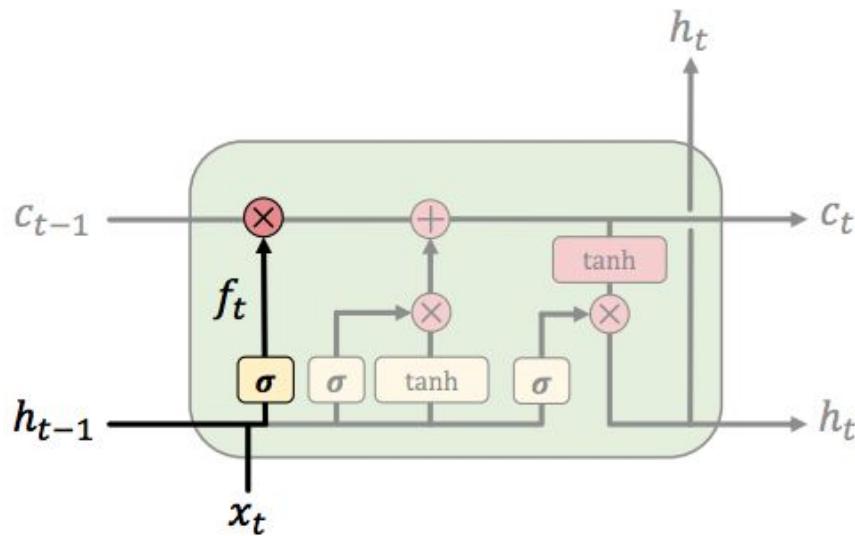
Long Short Term Memory (LSTMs)

How do LSTMs work?

- 1) Forget
- 2) Update
- 3) Output



LSTMs: forget irrelevant information

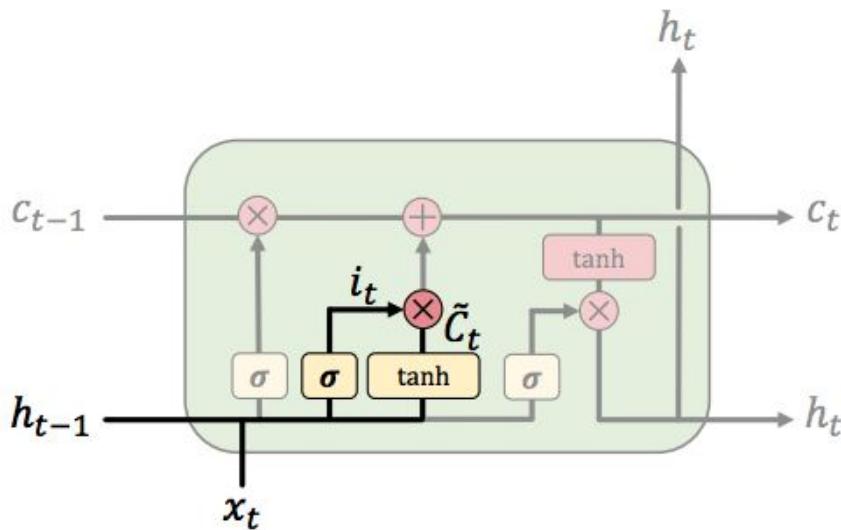


$$f_t = \sigma(\mathbf{W}_f [h_{t-1}, x_t] + b_f)$$

- Use previous cell output and input
- Sigmoid: value 0 and 1 – “completely forget” vs. “completely keep”

ex: Forget the gender pronoun of previous subject in sentence.

LSTMs: identify new info to be stored



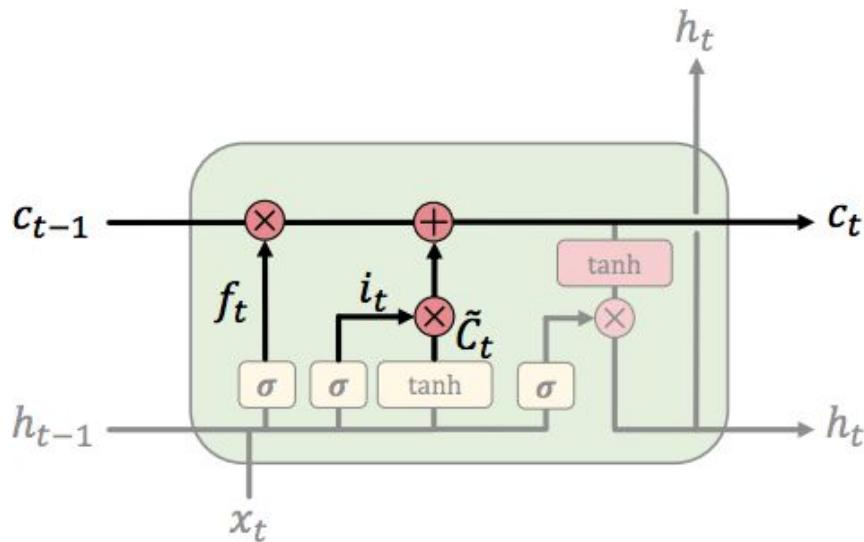
$$i_t = \sigma(\mathbf{W}_i [h_{t-1}, x_t] + b_i)$$

$$\tilde{c}_t = \tanh(\mathbf{W}_c [h_{t-1}, x_t] + b_c)$$

- Sigmoid layer: decide what values to update
- Tanh layer: generate new vector of "candidate values" that could be added to the state

ex: Add gender of new subject to replace that of old subject.

LSTMs: update cell state

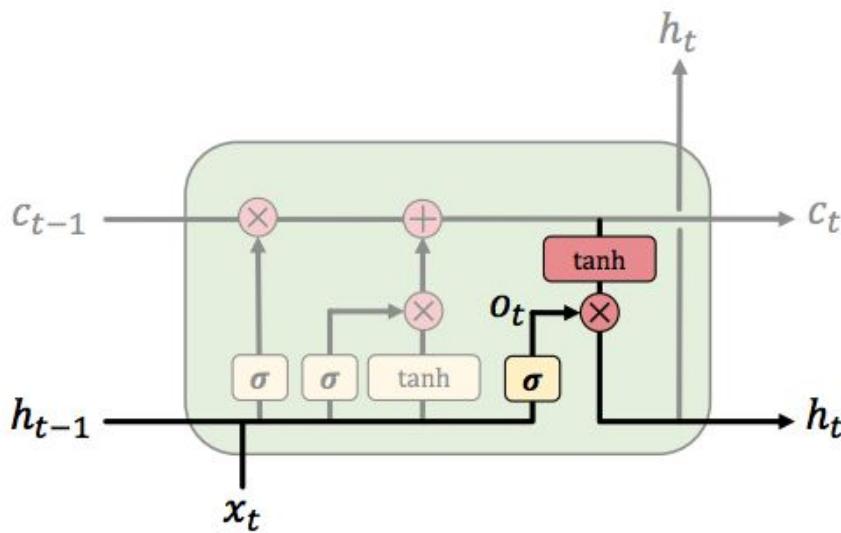


$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

- Apply forget operation to previous internal cell state: $f_t * C_{t-1}$
- Add new candidate values, scaled by how much we decided to update: $i_t * \tilde{C}_t$

ex: Actually drop old information and add new information about subject's gender.

LSTMs: output filtered version of cell state



$$o_t = \sigma(W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$

- Sigmoid layer: decide what parts of state to output
- Tanh layer: squash values between -1 and 1
- $o_t * \tanh(C_t)$: output filtered version of cell state

ex: Having seen a subject, may output information relating to a verb.

LSTMs: key concepts

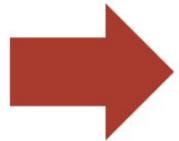
1. Maintain a **separate cell state** from what is outputted
2. Use **gates** to control the **flow of information**
 - Forget gate gets rid of irrelevant information
 - Selectively update cell state
 - Output gate returns a filtered version of the cell state
3. Backpropagation from c_t to c_{t-1} doesn't require matrix multiplication:
uninterrupted gradient flow

Recurrent neural networks (RNNs)

1. RNNs are well suited for **sequence modeling** tasks
2. Model sequences via a **recurrence relation**
3. Training RNNs with **backpropagation through time**
4. Gated cells like **LSTMs** let us model **long-term dependencies**

One-hot encoding of Words

Vocabulary:
Man, woman, boy,
girl, prince,
princess, queen,
king, monarch



	1	2	3	4	5	6	7	8	9
man	1	0	0	0	0	0	0	0	0
woman	0	1	0	0	0	0	0	0	0
boy	0	0	1	0	0	0	0	0	0
girl	0	0	0	1	0	0	0	0	0
prince	0	0	0	0	1	0	0	0	0
princess	0	0	0	0	0	1	0	0	0
queen	0	0	0	0	0	0	1	0	0
king	0	0	0	0	0	0	0	1	0
monarch	0	0	0	0	0	0	0	0	1

Each word gets
a 1×9 vector
representation

Word Embedding

Try to build a lower dimensional embedding

Vocabulary:
Man, woman, boy,
girl, prince,
princess, queen,
king, monarch



	Femininity	Youth	Royalty
Man			
Woman			
Boy			
Girl			
Prince			
Princess			
Queen			
King			
Monarch			

Word Embedding

Try to build a lower dimensional embedding

Vocabulary:
Man, woman, boy,
girl, prince,
princess, queen,
king, monarch



	Femininity	Youth	Royalty
Man	0	0	0
Woman	1	0	0
Boy	0	1	0
Girl	1	1	0
Prince	0	1	1
Princess	1	1	1
Queen	1	0	1
King	0	0	1
Monarch	0.5	0.5	1

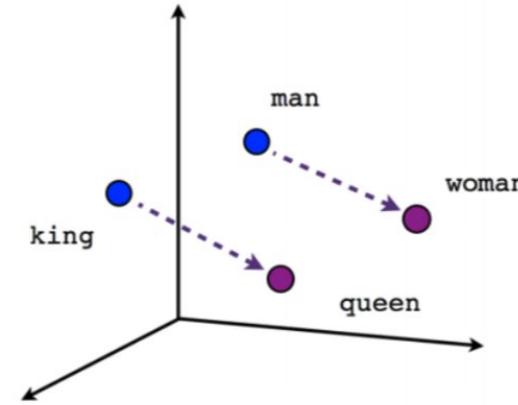
Each word gets a
1x3 vector

Similar words...
similar vectors

Word2Vec Word Embedding

Here's a list of words associated with "Sweden" using Word2vec, in order of proximity:

Word	Cosine distance
norway	0.760124
denmark	0.715460
finland	0.620022
switzerland	0.588132
belgium	0.585835
netherlands	0.574631
iceland	0.562368
estonia	0.547621
slovenia	0.531408



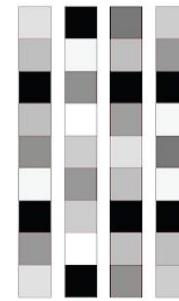
$$[[\text{king}]] - [[\text{man}]] + [[\text{woman}]] = [[\text{queen}]]$$

One-hot encoding vs. Embedding



One-hot word vectors:

- Sparse
- High-dimensional
- Hard-coded



Word embeddings:

- Dense
- Lower-dimensional
- Learned from data

Lab

