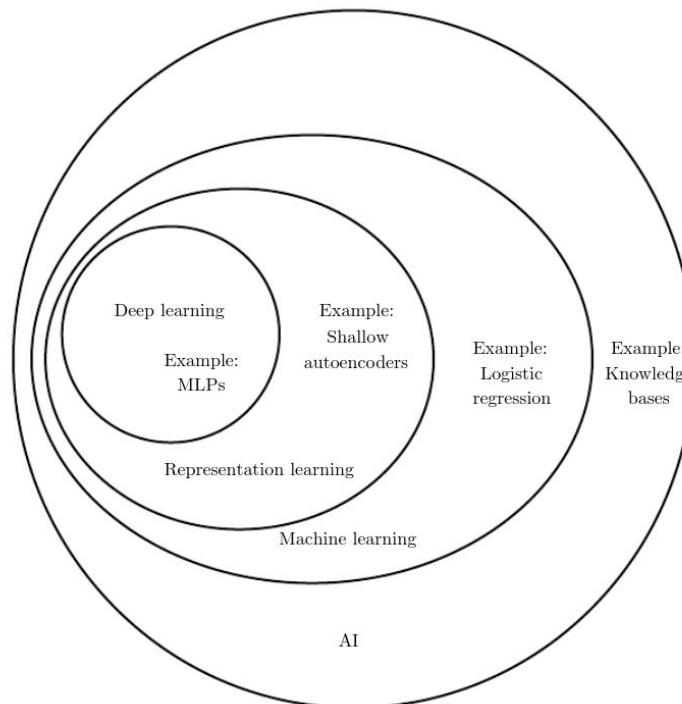


A Selective Introduction to Deep Learning

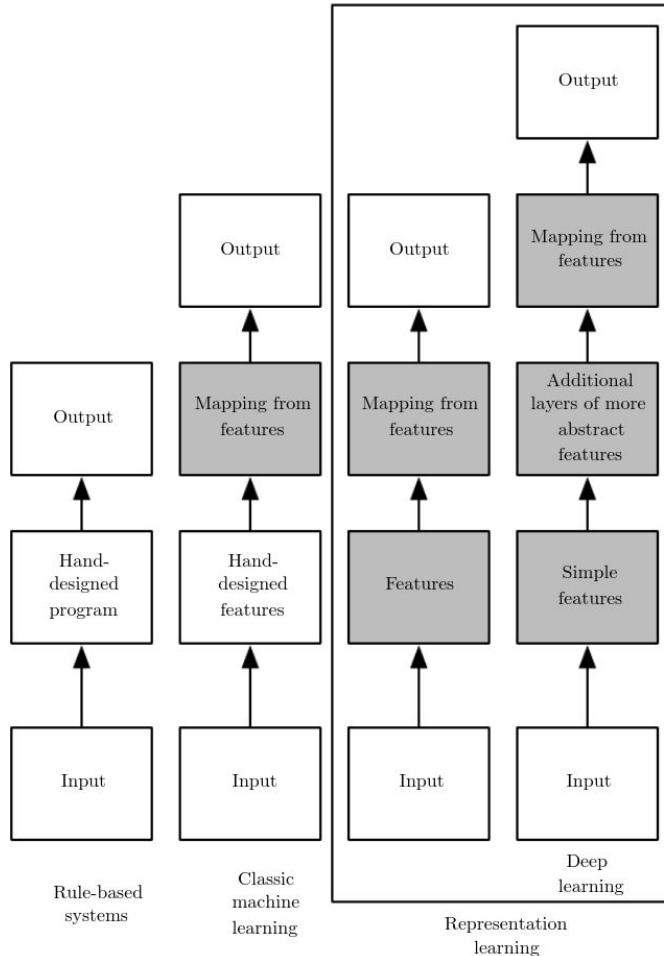
Deep Learning

The goal of Deep Learning is to:

- Allow computers to learn from experience (Algorithms - SGD).
- Understand the world in terms of a hierarchy of concepts, with each concept defined through its relation to simpler concepts (Networks - CNN).
- Gathering knowledge from experience (Data - ImageNet).



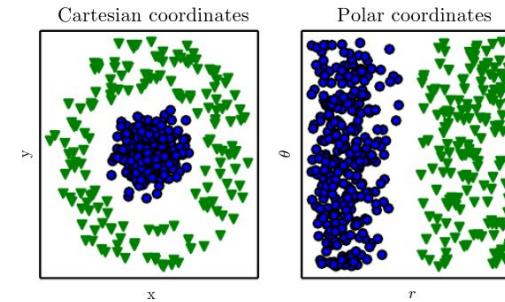
Classic AI and Deep Learning



Systems relying on hard-coded knowledge lack the ability to acquire their own knowledge.

-> Machine learning: enabled computers to tackle problems involving knowledge of the real world and make decisions that appear subjective.

- The performance of ML algorithms depends heavily on the representation of the data they are given.



-> Reuse ML, but this time not only the mapping from representation to output, but also the representation itself.

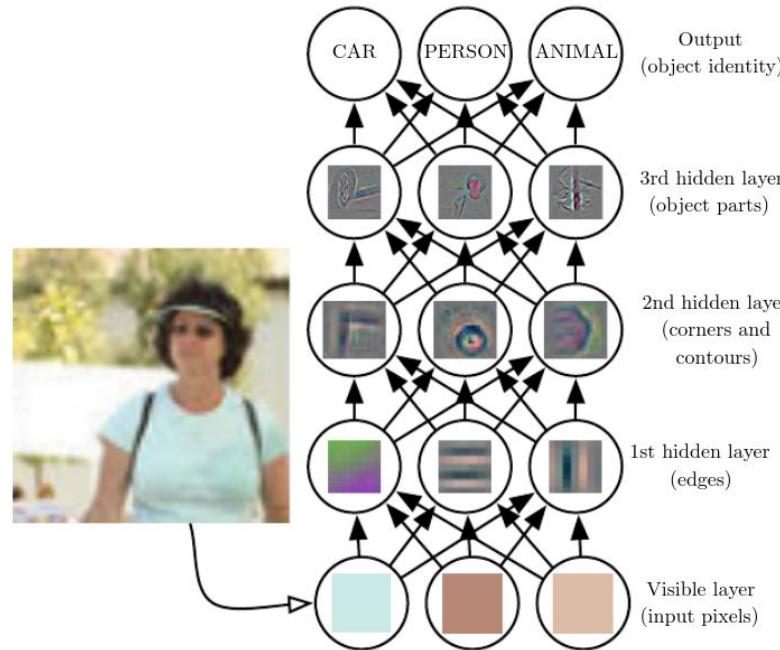
-> Learned representations often result in much better performance than can be obtained with hand-designed representations.

- The factors of variation influence every single piece of data we are able to observe. The individual pixels in an image of a red car might be very close to black at night.

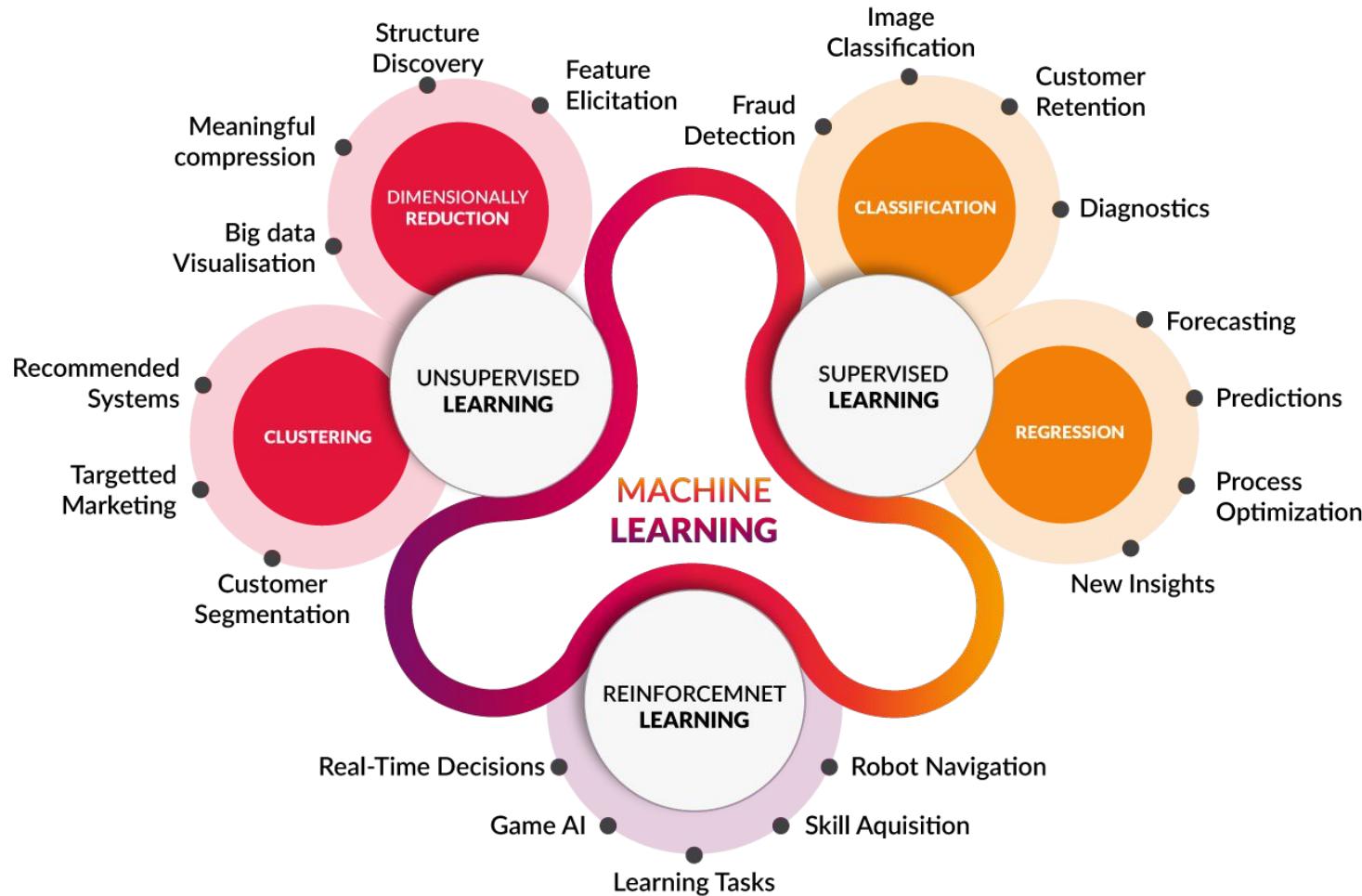
-> **Deep Learning:** disentangles the factors of variation and discard the ones that we do not care about

Deep Learning

Deep learning introduces representations that are expressed in terms of other, simpler representations.
Deep learning enables the computer to build complex concepts out of simpler concepts.



DL: in a nutshell is a composition of many nonlinear functions to model the complex dependency between input features and labels.



Deep Learning

Modern machine learning deals with the problem of *learning from data*.

Given a training dataset $\{(y_i, \mathbf{x}_i)\}_{1 \leq i \leq n}$ Inputs: $\mathbf{x}_i \in \mathbb{R}^d$ Outputs: $y_i \in \mathbb{R}$

-> The objective is to learn a function $f : \mathbb{R}^d \mapsto \mathbb{R}$ from a certain function class \mathcal{F} that has good prediction performance on test data.

Images Classification: \mathbf{x} are raw RGB images and \mathbf{y} are class labels.

Class of functions \mathcal{F}

- Linear classifiers (e.g., linear / logistic regression).
- Kernel methods (e.g., support vector machines).
- Tree-based methods (e.g., decision trees, random forests).
- Nonparametric regression (e.g., nearest neighbors).

In *Deep learning*, we have a compositional function class:

$$\{f(\mathbf{x}; \boldsymbol{\theta}) = \mathbf{W}_L \sigma_L(\mathbf{W}_{L-1} \cdots \sigma_2(\mathbf{W}_2 \sigma_1(\mathbf{W}_1 \mathbf{x}))) \mid \boldsymbol{\theta} = \{\mathbf{W}_1, \dots, \mathbf{W}_L\}\}$$

Deep Learning Example: ImageNet Challenge, aka ILSVRC

A classification task: A dataset of 1.2 million color images with 1000 classes.

The performance of a classifier is then evaluated on a test dataset of 100 thousand images, and in the end the top-5 error is reported.

Model	Year	# Layers	# Params	Top-5 error
Shallow	< 2012	—	—	> 25%
AlexNet	2012	8	61M	16.4%
VGG19	2014	19	144M	7.3%
GoogleNet	2014	22	7M	6.7%
ResNet-152	2015	152	60M	3.6%

DL methods have a clear edge over shallow models that fit linear models / tree-based models on handcrafted features.

But there is a catch. Deep Learning in practice, requires:

- huge datasets that often contain millions of samples.
- immense computing power resulting (GPUs - TPUs).

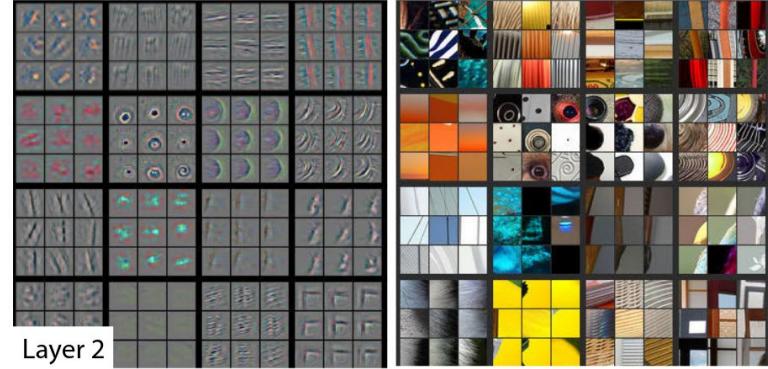
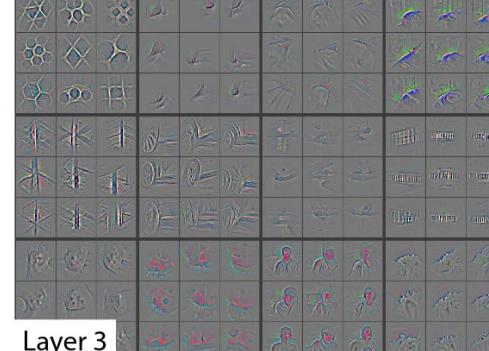
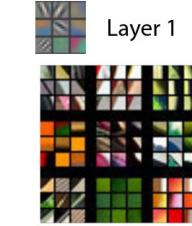
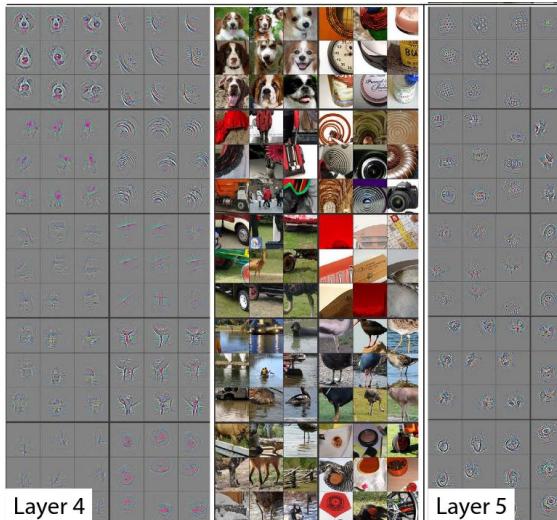
Can be hard to get right due to:

- *over-parameterization*: the number of parameters in state-of-the-art deep learning models is often much larger than the sample size (possible overfitting).
- *nonconvexity*: even with the help of GPUs, training deep learning models is still hard.

Deep in Deep Learning

Deep learning expresses complicated nonlinearity through composing many nonlinear functions;

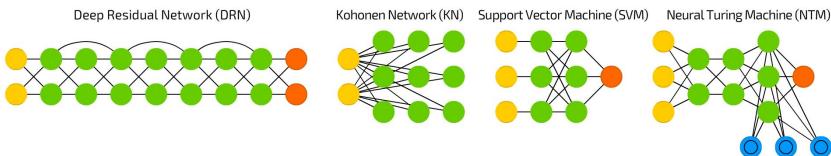
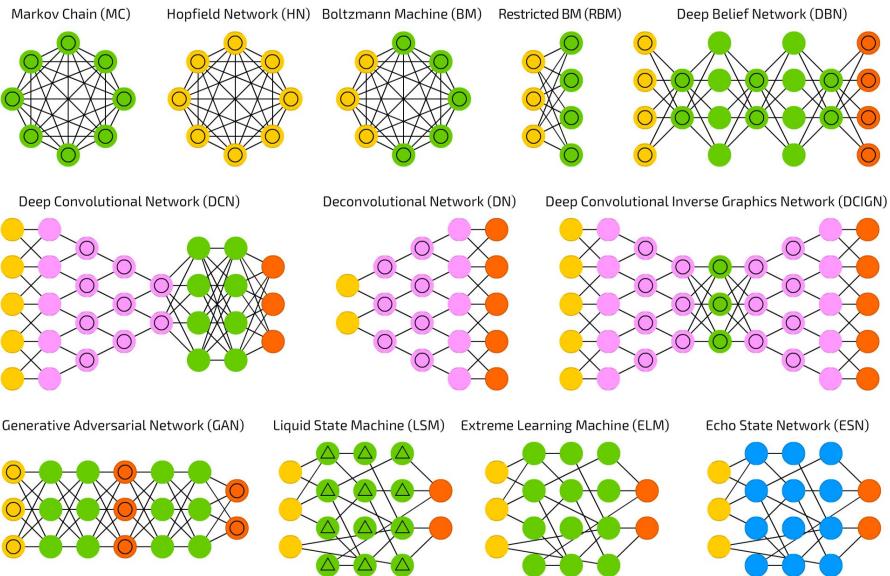
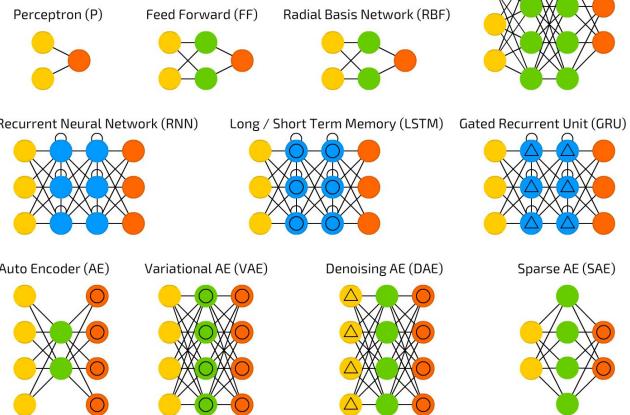
- > Composing lower-level features to create higher-level ones.
- > Represents the data in a different function space suitable for solving difficult tasks.
- > With many layers, we'll have a larger class function \mathcal{F} (small approximation error).
- > A biggest generalization power, with deep net & SGD, we'll often have a small training error.



A mostly complete chart of
Neural Networks

©2016 Fjodor van Veen - asimovinstitute.org

- Backfed Input Cell
- Input Cell
- △ Noisy Input Cell
- Hidden Cell
- Probabilistic Hidden Cell
- △ Spiking Hidden Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- △ Different Memory Cell
- Kernel
- Convolution or Pool



DL basics

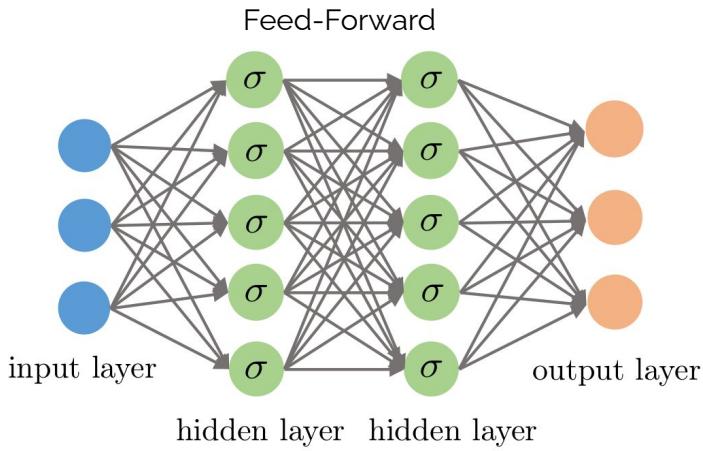
Notations: For a classification problem (regression can be addressed similarly):

The training dataset: $\{(y_i, \mathbf{x}_i)\}_{1 \leq i \leq n}$ where $y_i \in [K] \triangleq \{1, 2, \dots, K\}$ and $\mathbf{x}_i \in \mathbb{R}^d$ (We assume data points are IID)

Goal: find a function $\hat{f}(\mathbf{x})$ that predicts the labels \mathbf{y} for a new inputs \mathbf{x} (for the same distribution).

Terminology:

- inputs x_i : features;
- outputs y_i : labels;
- (x_i, y_i) example
- \hat{f} Classifier



Deep nets use composition of a series of simple nonlinear functions to model nonlinearity:

$$\mathbf{h}^{(L)} = \mathbf{g}^{(L)} \circ \mathbf{g}^{(L-1)} \circ \dots \circ \mathbf{g}^{(1)}(\mathbf{x})$$

$$\text{Recursively: } \mathbf{h}^{(0)} \triangleq \mathbf{x} \quad \mathbf{h}^{(l)} = \mathbf{g}^{(l)}(\mathbf{h}^{(l-1)})$$

Feed-forward nets, aka; *multilayer perceptrons (MLPs)*

$$\mathbf{h}^{(l)} = \mathbf{g}^{(l)}(\mathbf{h}^{(l-1)}) \triangleq \sigma(\mathbf{W}^{(l)} \mathbf{h}^{(l-1)} + \mathbf{b}^{(l)})$$

Examples of non-linearity:

- ReLU $[\sigma(\mathbf{z})]_j = \max\{z_j, 0\}$
- Sigmoid $(1 + e^{-z})^{-1}$

The output:

$$f_k(\mathbf{x}; \boldsymbol{\theta}) \triangleq \frac{\exp(z_k)}{\sum_k \exp(z_k)}, \quad \forall k \in [K]$$

$$\mathbf{z} = \mathbf{W}^{(L+1)} \mathbf{h}^{(L)} + \mathbf{b}^{(L+1)} \in \mathbb{R}^K$$

DL basics

With the output probabilities $f_k(\mathbf{x}; \boldsymbol{\theta}) \triangleq \frac{\exp(z_k)}{\sum_k \exp(z_k)}$, $\forall k \in [K]$, we compute the loss (cross-entropy)

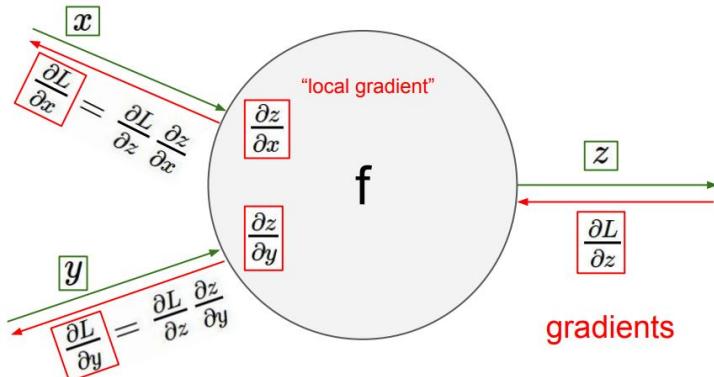
$$\mathcal{L}(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y) = -\sum_{k=1}^K \mathbb{1}\{y = k\} \log p_k, \quad \boldsymbol{\theta} \triangleq \{\mathbf{W}^{(\ell)}, \mathbf{b}^{(\ell)} : 1 \leq \ell \leq L + 1\}.$$

Backpropagation:

Using *stochastic gradient descent*, start from random weights $\boldsymbol{\theta}^0$ then iteratively updates the parameters $\boldsymbol{\theta}^t$ we move in the direction of the negative gradient to reduce the loss.

- Why stochastic? each iteration is based only on a subsample of the data. A batch $\mathcal{B} \subset [n]$
-> Reduces computation, converges to the solutions that can be obtained with GD.

$$\ell_{\mathcal{B}}(\boldsymbol{\theta}) \triangleq |\mathcal{B}|^{-1} \sum_{i \in \mathcal{B}} \mathcal{L}(\mathbf{f}(\mathbf{x}_i; \boldsymbol{\theta}), y_i)$$



Back-prop in computational graphs

$$\frac{\partial \ell_{\mathcal{B}}}{\partial \mathbf{h}^{(\ell-1)}} = \frac{\partial \mathbf{h}^{(\ell)}}{\partial \mathbf{h}^{(\ell-1)}} \cdot \frac{\partial \ell_{\mathcal{B}}}{\partial \mathbf{h}^{(\ell)}} :$$

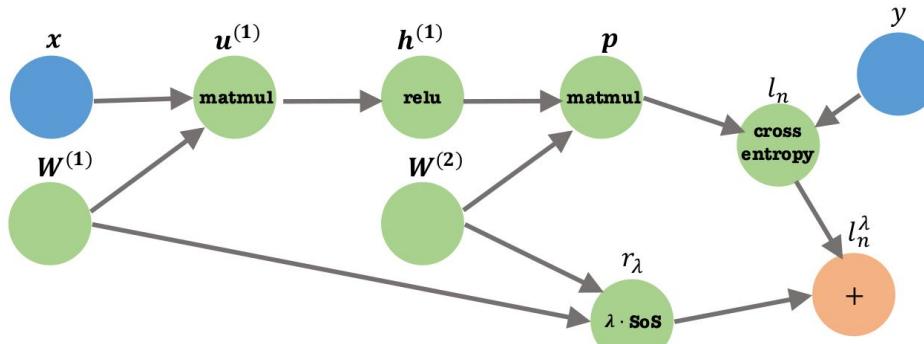
$$\mathbf{W}^{(\ell)} \leftarrow \mathbf{W}^{(\ell)} - \eta \frac{\partial \ell_{\mathcal{B}}}{\partial \mathbf{W}^{(\ell)}}, \quad \text{where} \quad \frac{\partial \ell_{\mathcal{B}}}{\partial W_{jm}^{(\ell)}} = \frac{\partial \ell_{\mathcal{B}}}{\partial h_j^{(\ell)}} \cdot \sigma' \cdot h_m^{(\ell-1)},$$

DL basics

Back-prop in computational graphs: With large models, we can quickly overfit (memorize) the data, without any generalization to the test set -> Penalize large variation in the weights

$$\ell_{\mathcal{B}}^{\lambda}(\boldsymbol{\theta}) = \ell_{\mathcal{B}}(\boldsymbol{\theta}) + r_{\lambda}(\boldsymbol{\theta}) = \ell_{\mathcal{B}}(\boldsymbol{\theta}) + \lambda \left(\sum_{j,j'} (W_{j,j'}^{(1)})^2 + \sum_{j,j'} (W_{j,j'}^{(2)})^2 \right)$$

In practice (PyTorch, TensorFlow), the model is represented as a computation graph: Each node represents a function (inside a circle), which is associated with an output of that function (outside a circle)



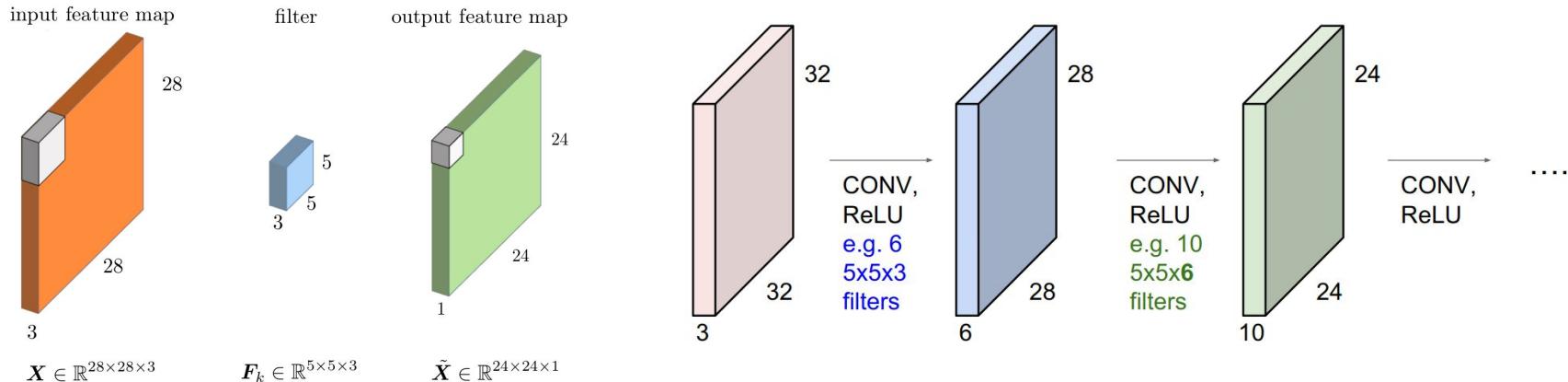
DL Models: Convolutional neural networks

A special type of feed-forward neural networks that is tailored for image processing.

Inputs: 3D tensors; $\mathbf{X} \in \mathbb{R}^{d_1 \times d_2 \times d_3}$

Convolutional layer

$$O_{ij}^k = \langle [\mathbf{X}]_{ij}, \mathbf{F}_k \rangle = \sum_{i'=1}^w \sum_{j'=1}^w \sum_{l=1}^{d_3} [\mathbf{X}]_{i+i'-1, j+j'-1, l} [\mathbf{F}_k]_{i', j', l}.$$

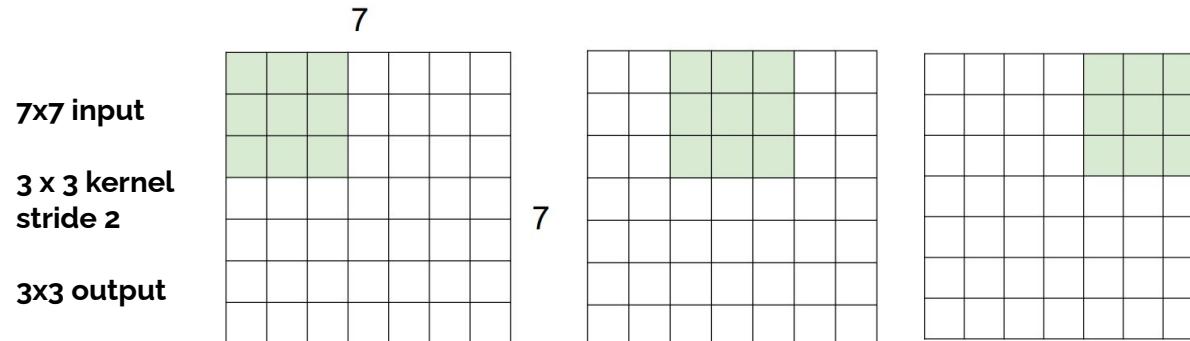
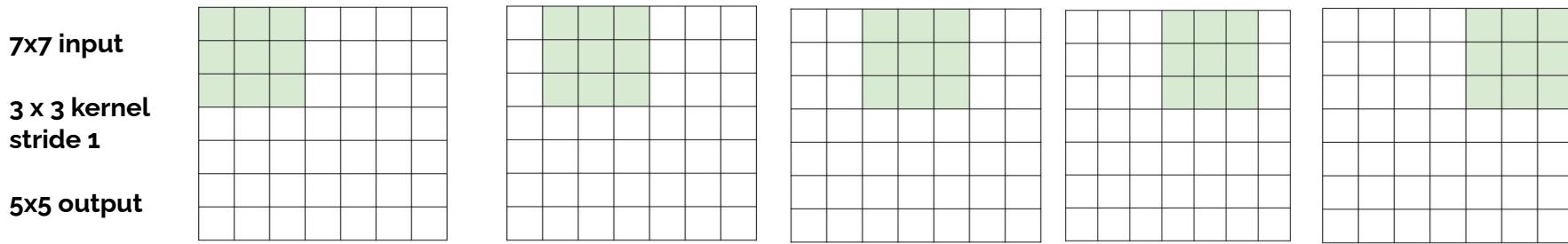
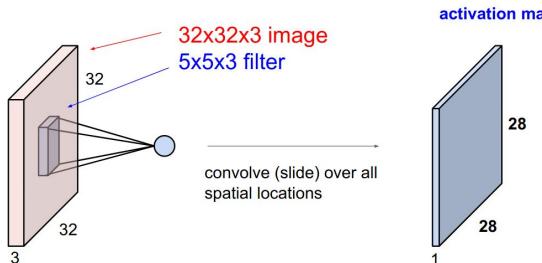


Nonlinearity $\tilde{X}_{ijk} = \sigma(O_{ijk}), \quad \forall i \in [d_1 - w + 1], j \in [d_2 - w + 1], k \in [\tilde{d}_3].$

$$[\sigma(\mathbf{z})]_j = \max\{z_j, 0\}$$

DL Models: Convolutional neural networks

Convolutional layer: A closer look



Output size:
(N - F) / stride + 1

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

DL Models: Convolutional neural networks

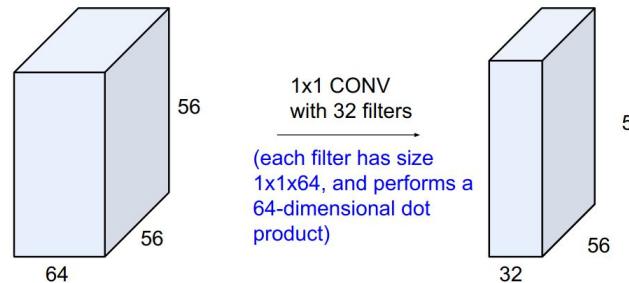
Convolutional layer: Keeping the same spatial dimensions

7x7 input

3 x 3 kernel | stride 1 | padding 1

7x7 output

0	0	0	0	0	0				
0									
0									
0									
0									

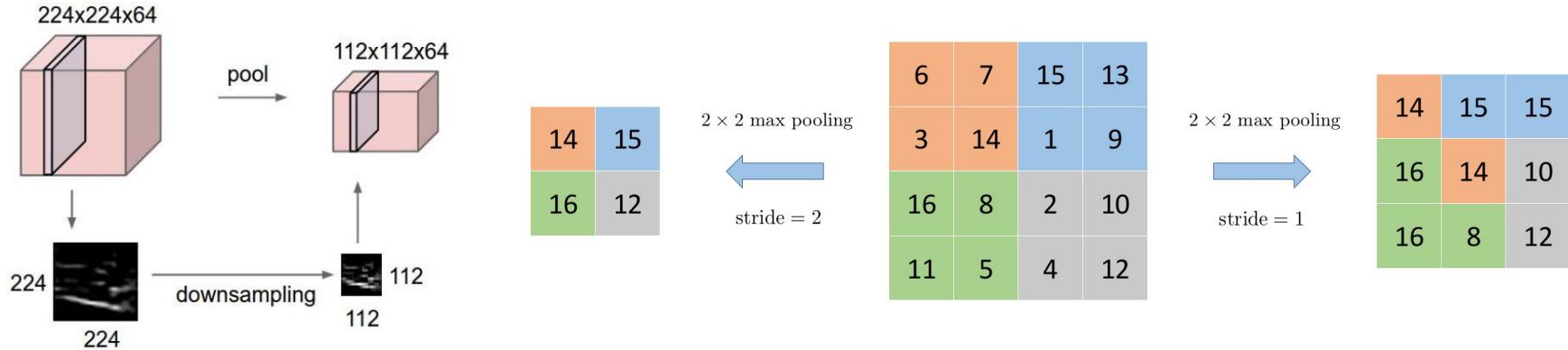


- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires four hyperparameters:
 - Number of filters K ,
 - their spatial extent F ,
 - the stride S ,
 - the amount of zero padding P .
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
 - $W_2 = (W_1 - F + 2P)/S + 1$
 - $H_2 = (H_1 - F + 2P)/S + 1$ (i.e. width and height are computed equally by symmetry)
 - $D_2 = K$
- With parameter sharing, it introduces $F \cdot F \cdot D_1$ weights per filter, for a total of $(F \cdot F \cdot D_1) \cdot K$ weights and K biases.
- In the output volume, the d -th depth slice (of size $W_2 \times H_2$) is the result of performing a valid convolution of the d -th filter over the input volume with a stride of S , and then offset by d -th bias.

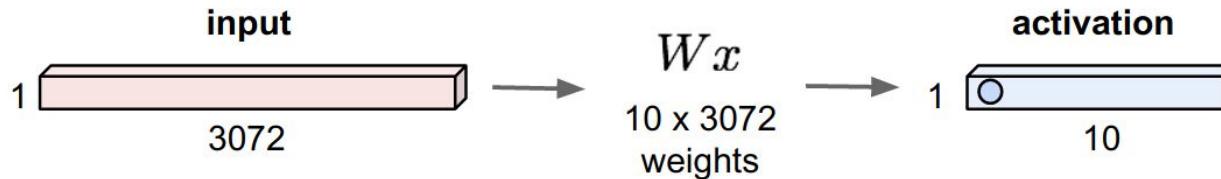
DL Models: Convolutional neural networks

3- **Pooling layers:** aggregates the information of nearby features into a single one. A 2×2 max-pooling filter computes

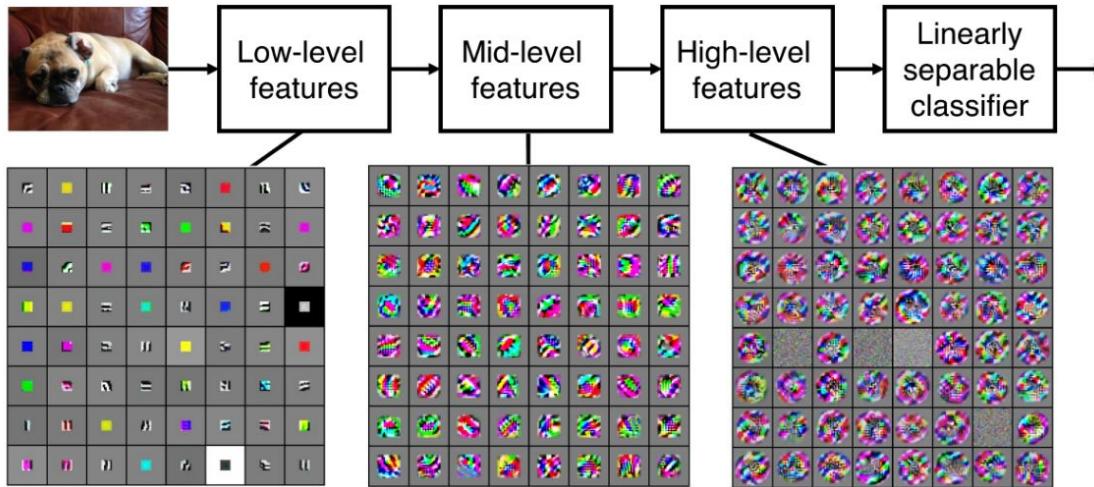
$$\max\{X_{i,j,k}, X_{i+1,j,k}, X_{i,j+1,k}, X_{i+1,j+1,k}\}$$



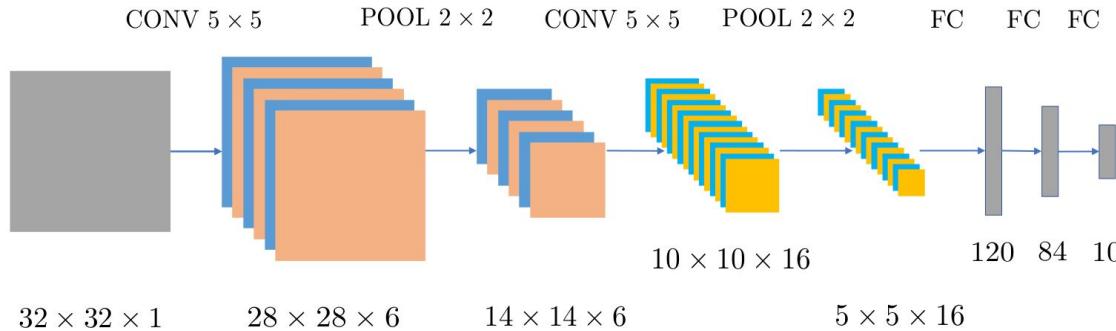
4- **Fully connected layers:** Treating 3D tensors as 2D vectors to compute the output probabilities $\tilde{\mathbf{X}} = \sigma(\mathbf{W}\text{Vec}(\mathbf{X}))$



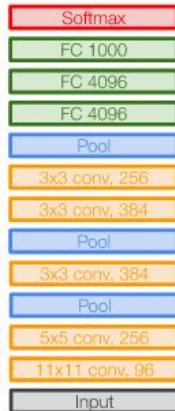
CNN architectures



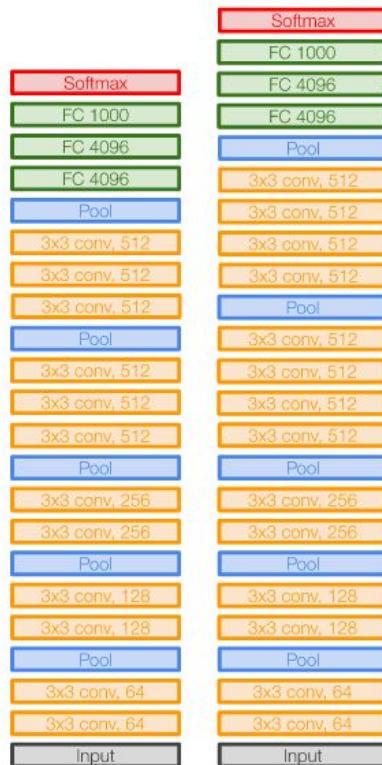
LeNet



CNN architectures

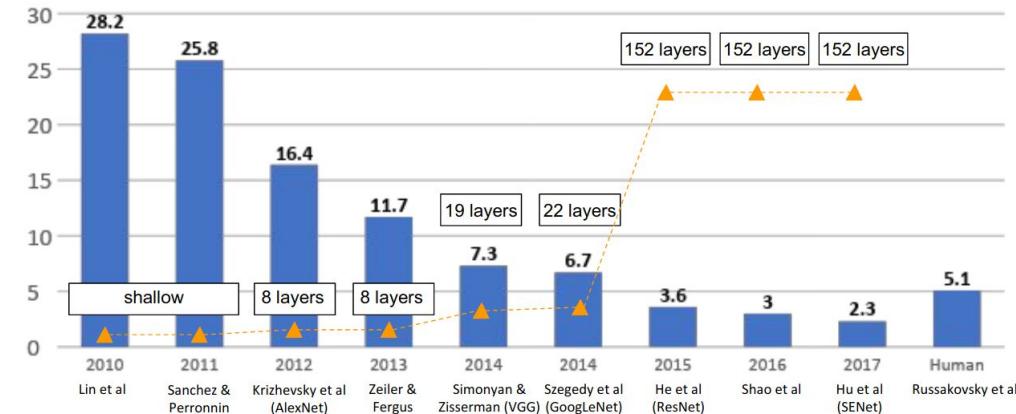


AlexNet



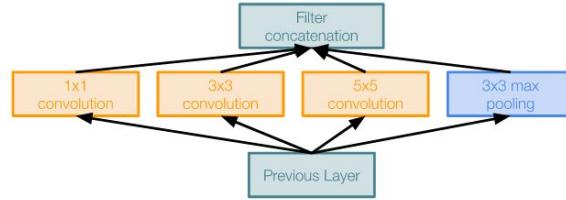
VGG16

VGG19

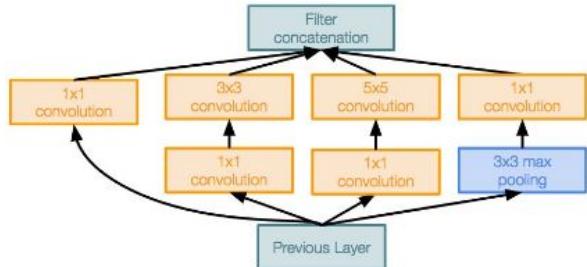


CNN architectures: GoogleNet

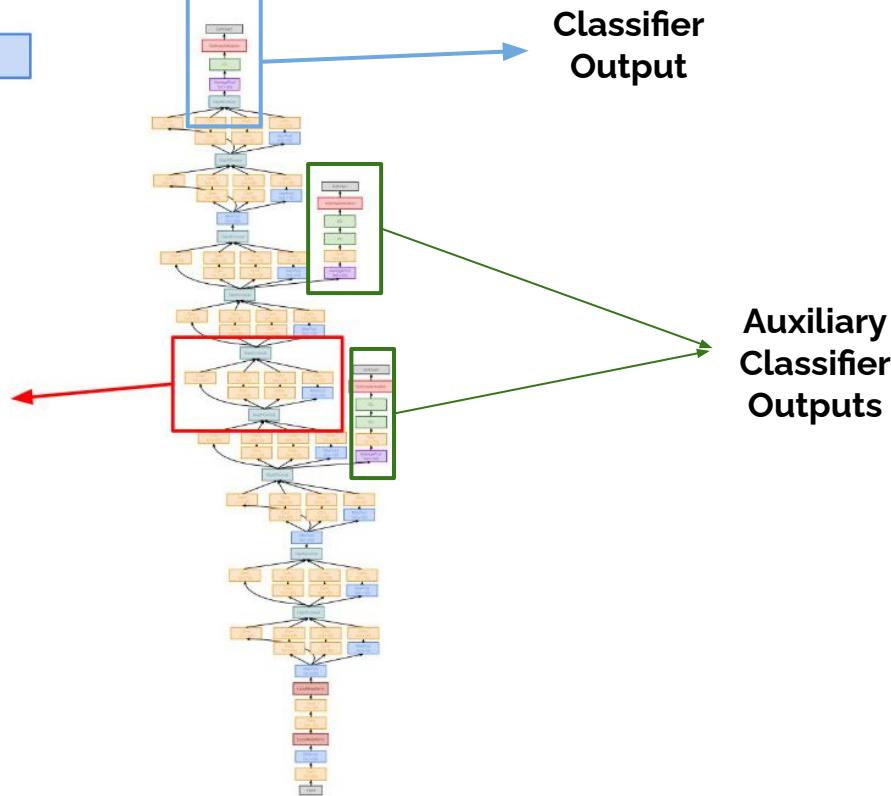
"Inception module": design a good local network topology (network within a network) and then stack these modules on top of each other



Naive Inception module



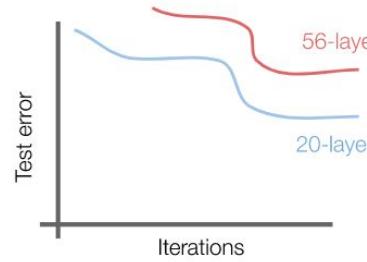
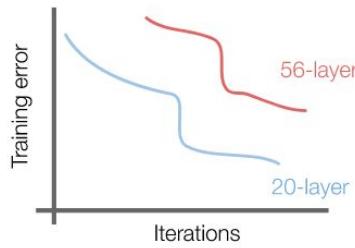
Inception module



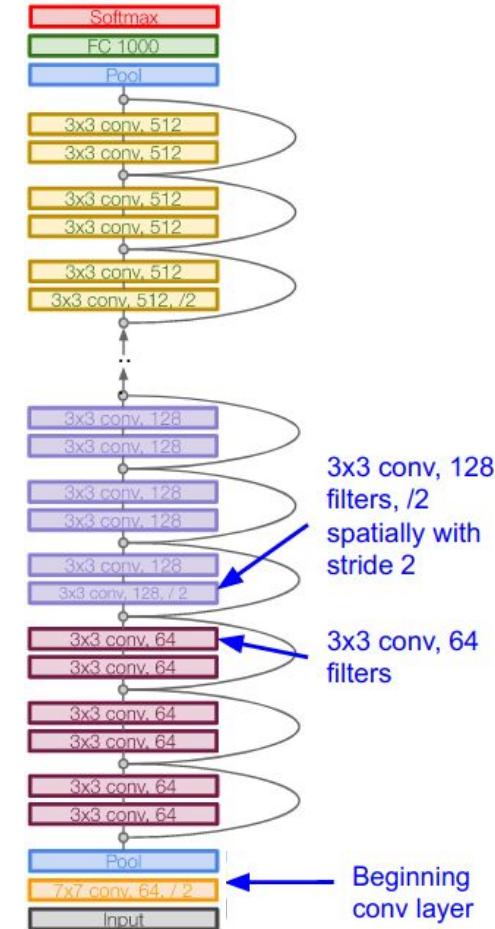
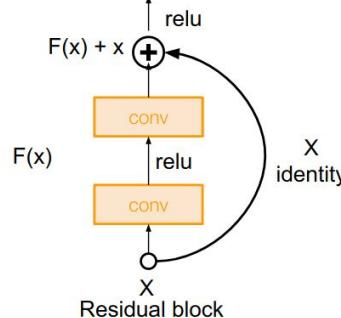
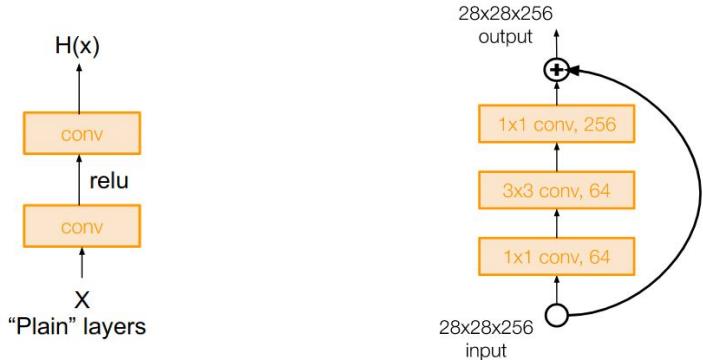
CNN architectures: ResNet

Very deep networks with residual connections: What happens when we continue stacking deeper layers on a "plain" convolutional neural network?

- the problem is an optimization problem, deeper models are harder to optimize



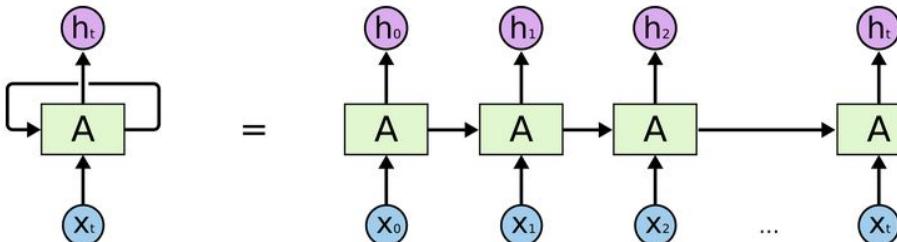
- + ACopying the learned layers from the shallower model and setting additional layers to identity mapping.



DL Models: Recurrent neural networks

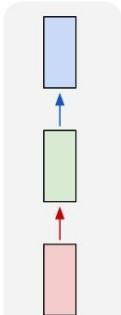
Recurrent nets:

A way to process sequences

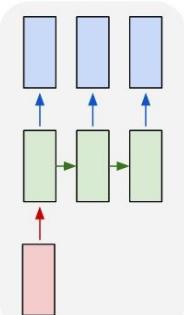


Different use cases

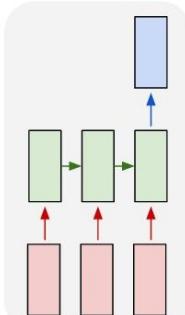
one to one



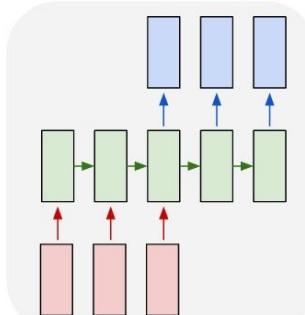
one to many



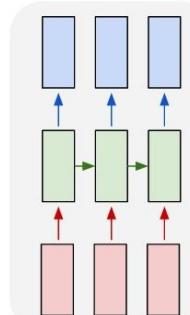
many to one



many to many



many to many



(1) Normal feedforward nets (e.g. image classification).

(2) Sequence output (e.g. image captioning - music generation).

(3) Séquence input (e.g. sentiment analysis).

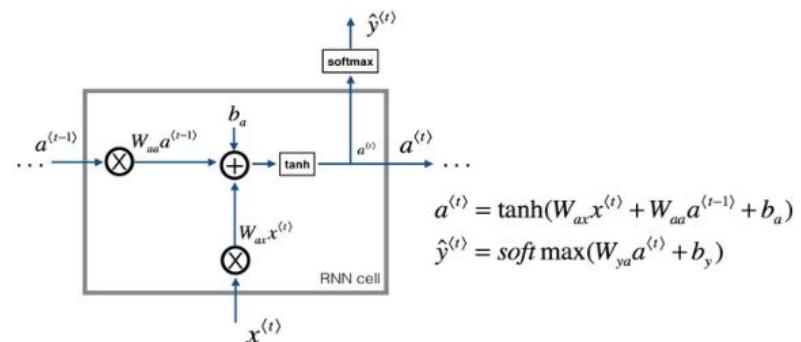
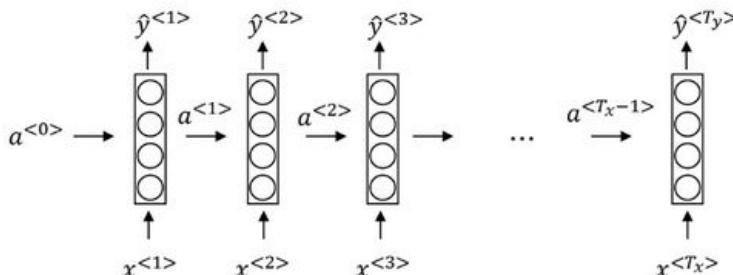
(4) Sequence input & output (e.g. machine translation).

(5) Synced sequence input and output (e.g. video classification).

Can be also used with inputs/outputs are fixed vectors

RNNs ~ a repetition of a single rnn cell

Machine translation	Voulez-vous chanter avec moi?	→	Do you want to sing with me?
Video activity recognition		→	Running
Name entity recognition	Yesterday, Harry Potter met Hermione Granger.	→	Yesterday, Harry Potter met Hermione Granger.



$$cache = (a^{(t)}, a^{(t-1)}, x^{(t)}, parameters)$$

- Propagate the input through the network to get the output.
- Calculate the error using the labels.
- Calculate the derivatives of the error with respect to the network weights.
- Adjust to minimise the error

$$\frac{\partial J}{\partial a^{(t-1)}} = \frac{\partial J}{\partial a^{(t)}} \frac{\partial a^{(t)}}{\partial a^{(t-1)}}$$

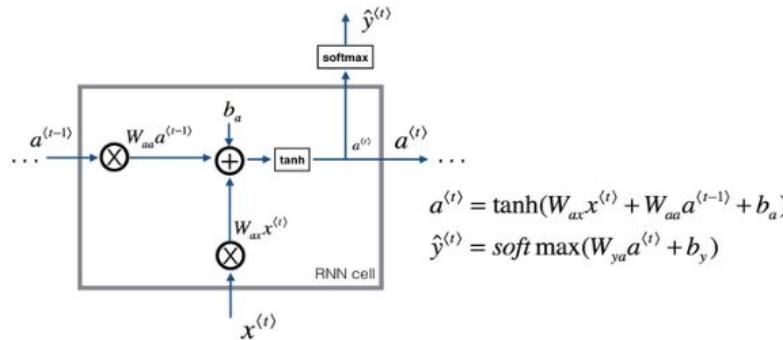
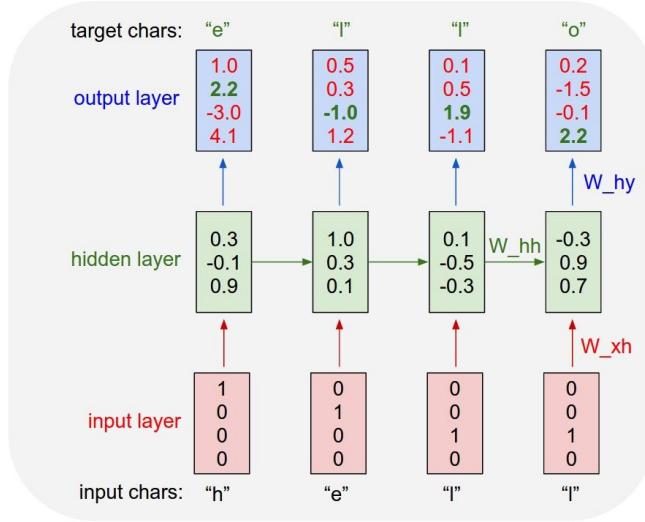
parameters gradients:

$$\frac{\partial a^{(t)}}{\partial W_x}, \frac{\partial a^{(t)}}{\partial W_a}, \frac{\partial a^{(t)}}{\partial b}$$

$$\frac{\partial J}{\partial x^{(t)}} = \frac{\partial J}{\partial a^{(t)}} \frac{\partial a^{(t)}}{\partial x^{(t)}}$$

$$\frac{\partial J}{\partial a^{(t)}}$$

Use case: Character-Level RNN



Pytorch example

```
import torch.nn as nn

class RNN(nn.Module):
    def __init__(self, input_size, hidden_size, output_size):
        super(RNN, self).__init__()

        self.hidden_size = hidden_size

        self.i2h = nn.Linear(input_size + hidden_size, hidden_size)
        self.i2o = nn.Linear(input_size + hidden_size, output_size)
        self.softmax = nn.LogSoftmax(dim=1)

    def forward(self, input, hidden):
        combined = torch.cat((input, hidden), 1)
        hidden = self.i2h(combined)
        output = self.i2o(combined)
        output = self.softmax(output)
        return output, hidden

    def initHidden(self):
        return torch.zeros(1, self.hidden_size)

n_hidden = 128
rnn = RNN(n_letters, n_hidden, n_categories)
```

Visualizing the predictions

t t o : / / w w w . y n e t n e w s . c o m /] E n g l i s h - l a n g u a g e w e b s i t e o f I s r a e l ' s l a r
l p : / / w w w . b a c a h e t s . c o m / - x g l i s h - l i n g u a g e s a i r s i t e o f I s r a e l i s s i n g
d : x n e . w a e a . a w a t o a . s & n t i a c a - s a r d e e l h o a n t b i s a n f a n r e i f ' a a t d
m w - 2 ♦ p i i s o e s s i . e r n . c] (d c e e n e p e s a a k i i e e l e d h , i r t h r a o n s e , c o s e
d r . < : a h b - n p t w t . x i g h / m a) T v d r y z i c o u e d l s u : t h a - o o t u , s t u i f l v e p e r y
s t p , t c o a z 2 d r u l w o c l e n s r] p . l l v a o d . , e y t c - n d m - o i b u v s] b b i m s u l t a t l y b n
g e s i t n e w s p a p e r ' ' [[Y e d i g h A h r o n o t h]] ' ' H e b r e w - l a n g u a g e p e r i o d e t a a w s p a p e r s o ' [[T e l t i (f e a n e m t i) : , : ' e r r e w s l e n g u a g e : a r o s o d i i r s c o e e n a i T T h o a i n n S r m u w] e y s [' i n e i a ' s i w d d e ' h s o l r i f r : u s . s e t l g o r s . a s a t C a r e e g ' a C l r i s z i e ' : , # : T A a a a t B a s e e i l o ' i a n f v l - t u a e v r t i d , B A M S u s y u t] A s a o i g s] , . . . : s M B o l o u s : T o u a - n : d w o a p n u a , d , i i u i t i c p .] (I S V H v t u s u i e D n o e g a n o . . .) : { C C u i b o h e C y b k s l s : r - e p c n t s
i c a l s : ' ' ' ' ' [G l o b e s] ' ' [h t t p : / / w w w . g l o b e s . c o . i l /] b u s i n e s s d a c a l : ' ' * ' ' [T a a b a] / / ([l i p : / / w w w . b u o b a l . c o m u n / s A - y t i n e s s a e t s t l ' [h A e o v e l t s a h a d : x g e . w a o i r . r t o a . e l . i T & a i e g e o o y t t ' ' & [& m C o e r o n e ' : , i ' o d w . : n i i s a a u r e . e n i / o m l c C . (e f t g i r i i u a ' n : , C : & # * : a f D r u s u l , o m e l p < , d h a ; d e u o o t / i h n c s i f S , u r h o s t , t u n n k i <] : & 1 s T G u i t r s i , : b a c m r - x t p o b - g r e s i s l e r l n a f a D] l o s p t a d , i f r m
i l y * ' ' [[H a a r e t z H ' A r e t z i]] ' ' [h t t p : / / w w w . h a a r e t z . c o . i l /] R e l a t i v l y * ' [[T e r r d n F e r a n t a h]] / ([l i p : / / w w w . b o n m d s t . c o m u n / s - e s a t o i r e ' ' h A i l n n t e H a l s r c n o l ' s a h a d : x n e . w a a m r t d h e o h . o l . c & o p i n i v e k i . : * s C O S a n t h i T i m ' l i e : , i m c d w - 2 ♦ p h i i s e r d i t . i n a / c m f i . (a f l c a n a d s - ! [t B T C o m m g d]] W o n a a e , : . b a e r r . < t a i b - d u l c n c / a r n e s i] l i c e y s t o n d s # & : G l D u v c c s a o S u c l t e l] z , : o ' o m t] , : e o a 2 n i v f s r o o e i u n a l a) u v v r o

Cell sensitive to position in line:

The sole importance of the crossing of the Berezina lies in the fact that it plainly and indubitably proved the fallacy of all the plans for cutting off the enemy's retreat and the soundness of the only possible line of action--the one Kutuzov and the general mass of the army demanded--namely, simply to follow the enemy up. The French crowd fled at a continually increasing speed and all its energy was directed to reaching its goal. It fled like a wounded animal and it was impossible to block its path. This was shown not so much by the arrangements it made for crossing as by what took place at the bridges. When the bridges broke down, unarmed soldiers, people from Moscow and women with children who were with the French transport, all--carried on by vis inertiae--pressed forward into boats and into the ice-covered water and did not, surrender.

Cell that turns on inside quotes:

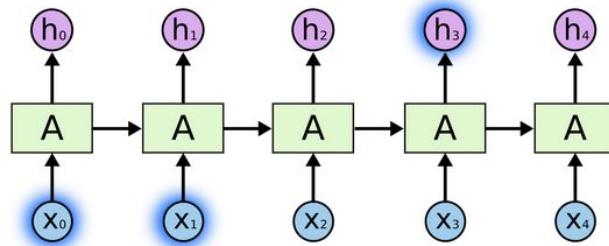
"You mean to imply that I have nothing to eat out of.... On the contrary, I can supply you with everything even if you want to give dinner parties," warmly replied Chichagov, who tried by every word he spoke to prove his own rectitude and therefore imagined Kutuzov to be animated by the same desire.

Kutuzov, shrugging his shoulders, replied with his subtle penetrating smile: "I meant merely to say what I said."

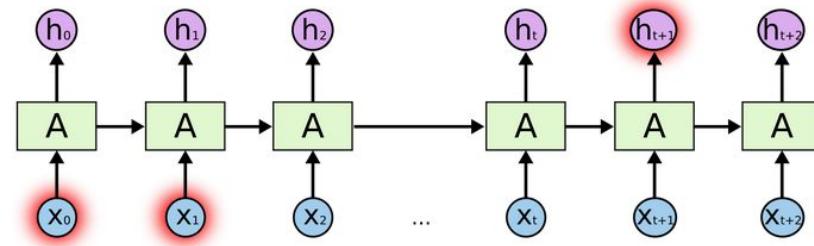
A large portion of cells are not easily interpretable. Here is a typical example:

```
/* Unpack a filter field's string representation from user-space
 * buffer. */
char* audit_unpack_string(void *bufp, size_t *remain, size_t len)
{
    char *str;
    if (!bufp || (len == 0) || (len > *remain))
        return ERR_PTR(-EINVAL);
    /* Of the currently implemented string fields, PATH_MAX
     * defines the longest valid length.
    */
```

The Problem of Long-Term Dependencies

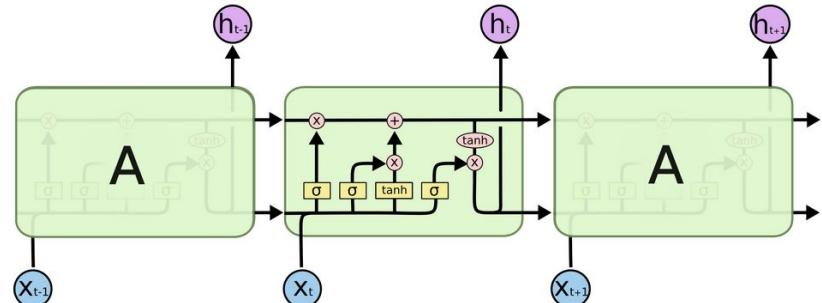
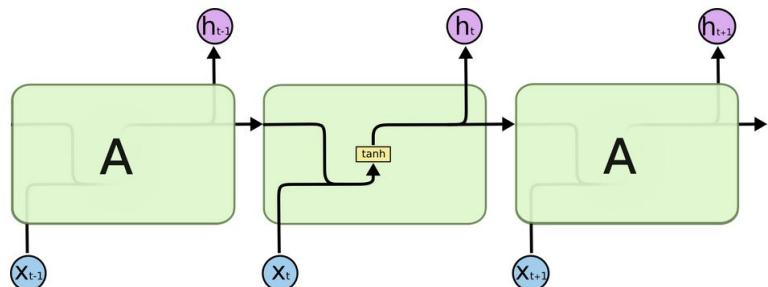


The **clouds** are in the **sky**

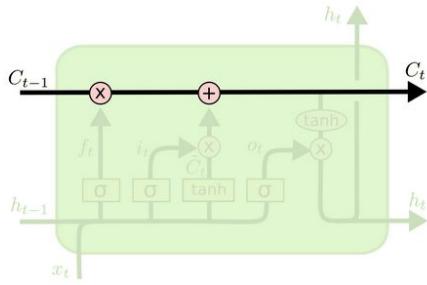


I grew up in **France**... I speak fluent **French**

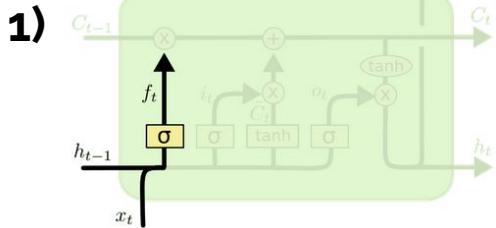
LSTM to the rescue



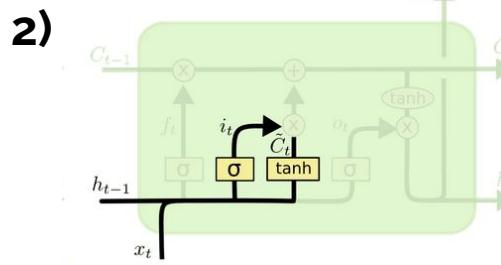
LSTMs Step-by-Step



The ability to remove or add information to the main cell state

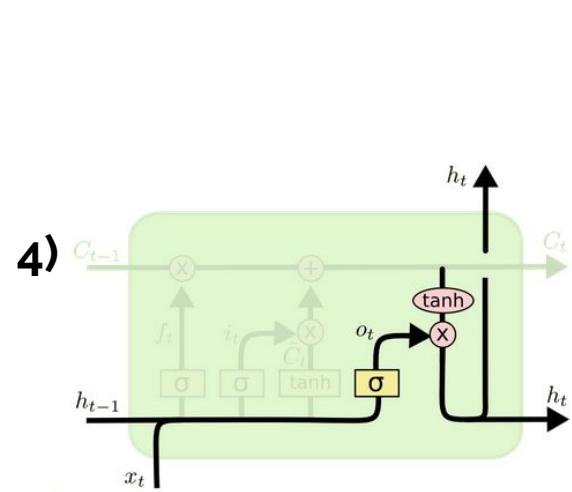


Forget gate : what information we're going to forget



Forget & Update

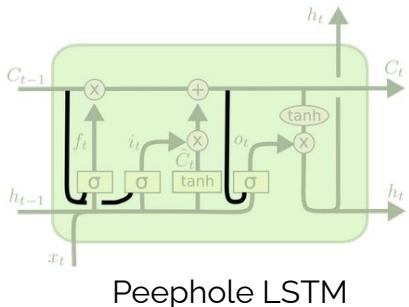
the cat/cats with .. ate was/were



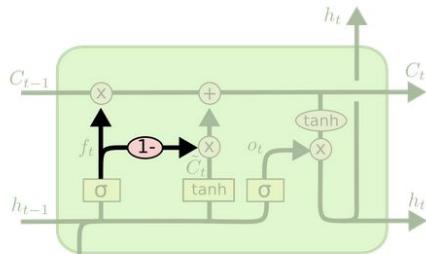
Output gate: From the new state of the main cell, output the relevant informations

Recurrent nets

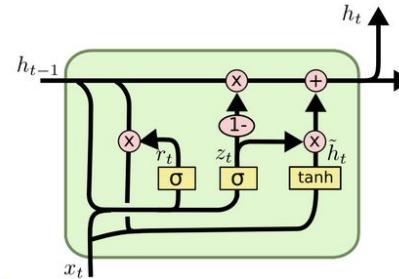
LSTMs variants



Peephole LSTM

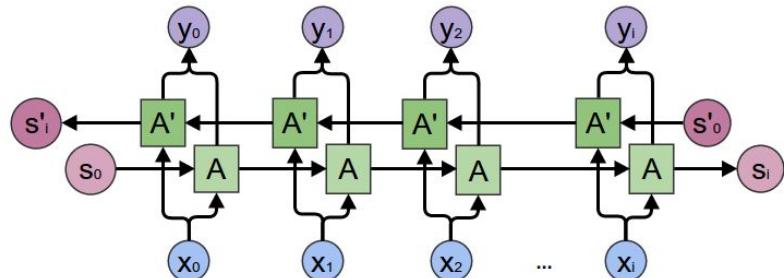


LSTM with coupled forget & update

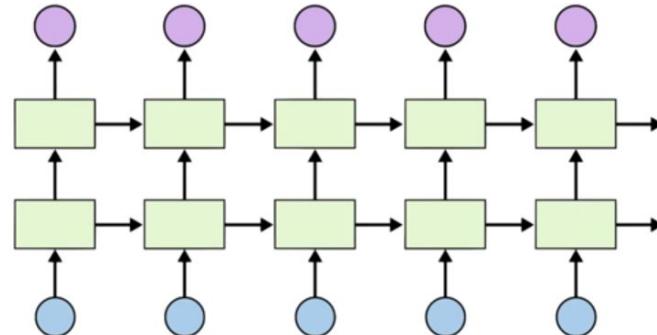


GRU: Gated recurrent unit

Bi-LSTMs



Stacked layers of LSTM cells

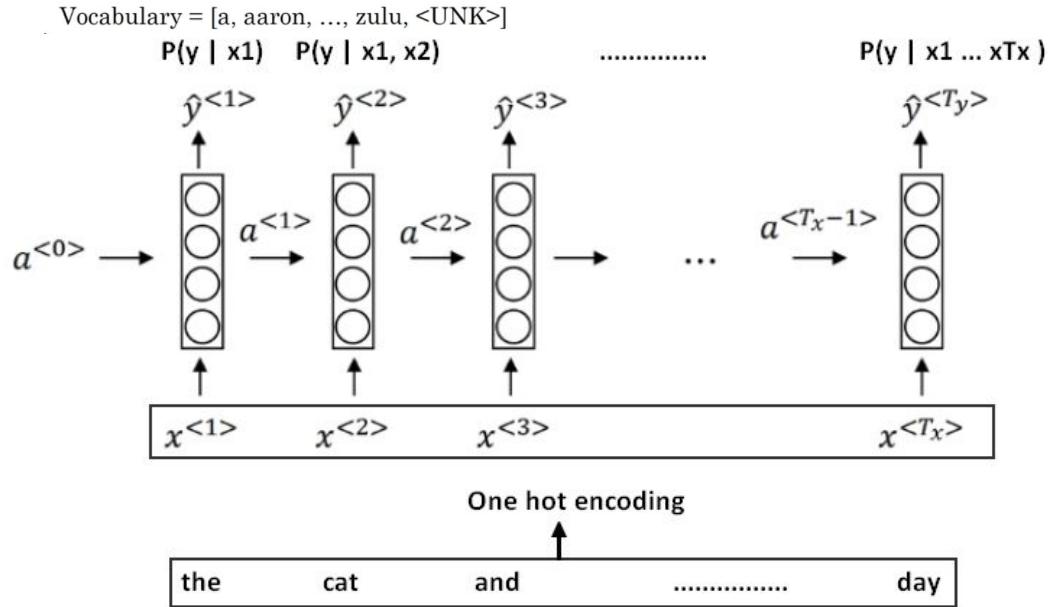


RNNs for Language Modeling

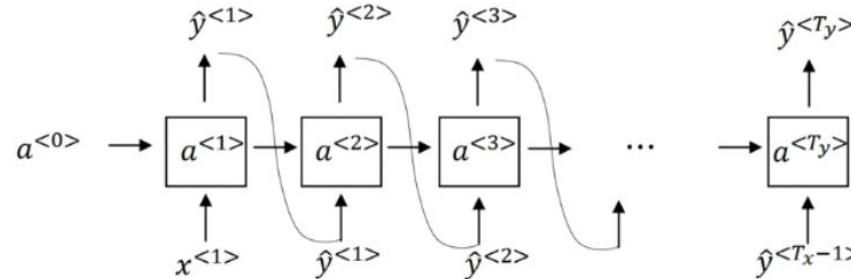
- Preprocess the inputs.
- One hot encoding.
- Batch the dataset.

For each example / batch:

- Calculate the CE loss at each time step.
- Back propagate the loss.



Use the trained RNN for language generation



Unsupervised learning

In supervised learning, we have a labeled dataset $\{(y_i, \mathbf{x}_i)\}$ and we focus on discriminative model, representing $\mathbb{P}(y | \mathbf{x})$ by a neural network $f(\mathbf{x}; \theta)$ with parameters θ .

Supervised Learning

Data: (\mathbf{x}, y)

\mathbf{x} is data, y is label

Goal: Learn a *function* to map $\mathbf{x} \rightarrow y$

Examples: Classification,
regression, object detection,
semantic segmentation, image
captioning, etc.

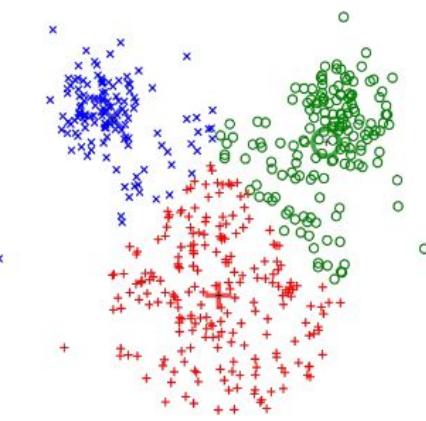
Unsupervised Learning

Data: \mathbf{x}

Just data, no labels!

Goal: Learn some underlying
hidden *structure* of the data

Examples: Clustering,
dimensionality reduction, feature
learning, density estimation, etc.



K-means clustering

Deep unsupervised learning

In unsupervised learning: We are interested in extracting information from the unlabeled data \mathbf{x} (without using labels \mathbf{y}).

- Can be low-dimensional embeddings of the data. -> **Autoencoders**
- A generative model capable of approximating the data distribution $\mathbb{P}_{\mathbf{X}}(\mathbf{x})$. -> **GaNs**

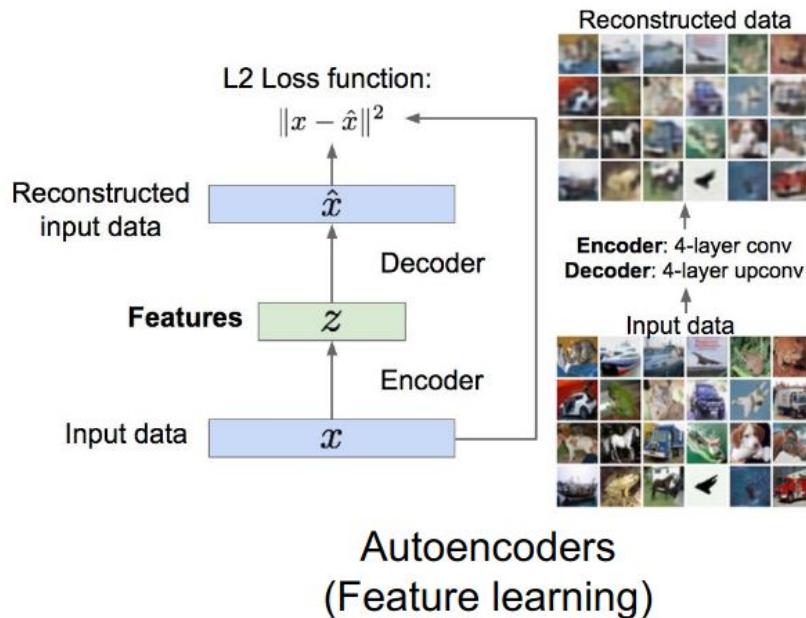
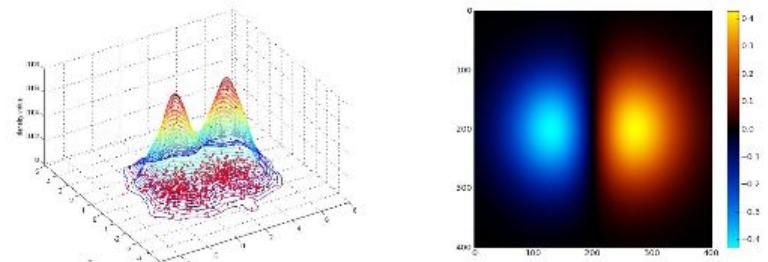


Figure copyright Ian Goodfellow, 2016. Reproduced with permission.

1-d density estimation



2-d density estimation

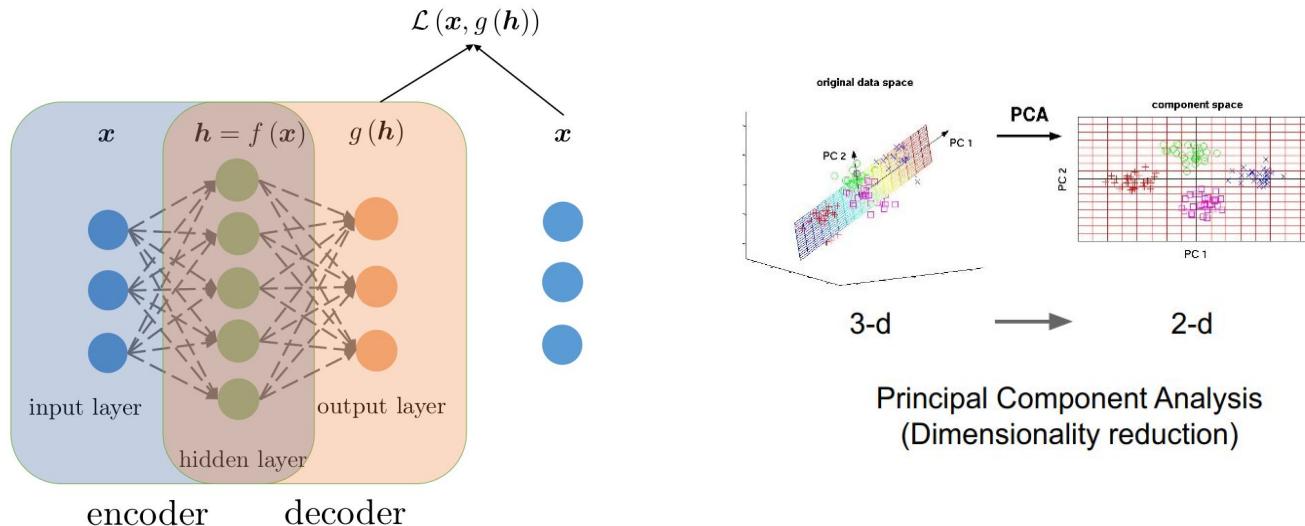
Autoencoders

in dimension reduction, the goal is to reduce the dimensionality of the data and at the same time preserve its salient features:

- PCA: embedding the data \mathbf{x}_i into a low-dimensional space using a linear function \mathbf{f} such that maximum variance can be explained (PCA is a special case of autoencoders).
- Equivalently; with autoencoders we want to learn two functions \mathbf{f} and \mathbf{g} so that the difference between $\mathbf{g}(\mathbf{f}(\mathbf{x}_i))$ and \mathbf{x}_i is minimal.

$$\mathbf{f}(\mathbf{x}) = \mathbf{W}_f \mathbf{x} \triangleq \mathbf{h} \quad \text{and} \quad \mathbf{g}(\mathbf{h}) = \mathbf{W}_g \mathbf{h}, \quad \text{where} \quad \mathbf{W}_f \in \mathbb{R}^{k \times d} \text{ and } \mathbf{W}_g \in \mathbb{R}^{d \times k}.$$

$$\underset{\mathbf{W}_f, \mathbf{W}_g}{\text{minimize}} \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{W}_f \mathbf{W}_g \mathbf{x}_i\|_2^2. \quad \underset{\mathbf{f}, \mathbf{g}}{\text{minimize}} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(\mathbf{x}_i, \mathbf{g}(\mathbf{h}_i)) \quad \text{with} \quad \mathbf{h}_i = \mathbf{f}(\mathbf{x}_i), \quad \text{for all } i \in [n].$$

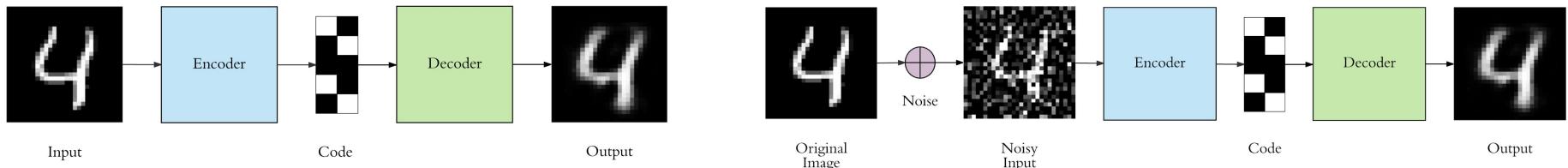


Autoencoders

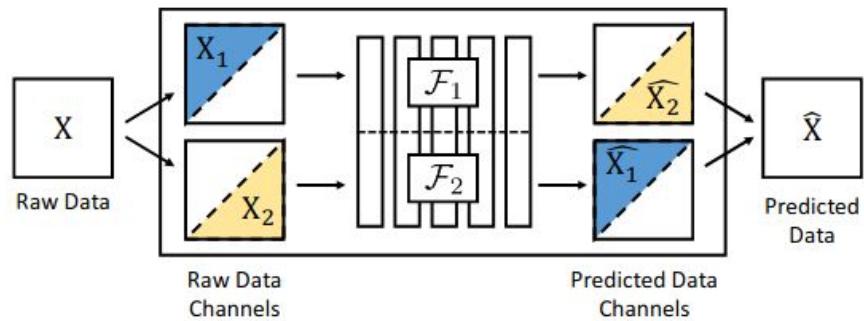
Sparse autoencoders: Find a sparse representation of the input data.

$$\min_{\mathbf{f}, \mathbf{g}} \underbrace{\frac{1}{n} \sum_{i=1}^n \mathcal{L}(\mathbf{x}_i, \mathbf{g}(\mathbf{h}_i))}_{\text{loss}} + \lambda \underbrace{\|\mathbf{h}_i\|_1}_{\text{regularizer}} \quad \text{with } \mathbf{h}_i = \mathbf{f}(\mathbf{x}_i), \text{ for all } i \in [n].$$

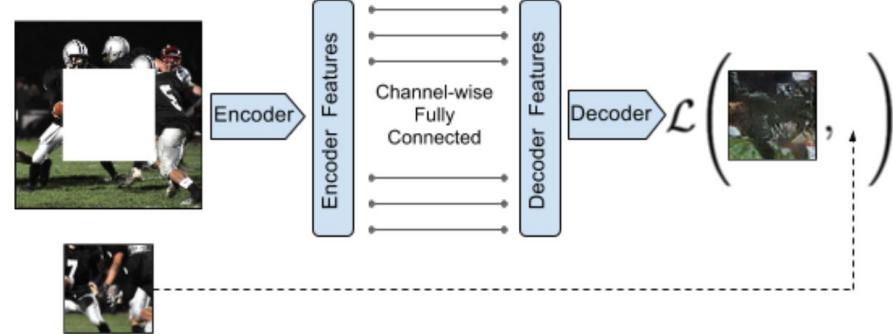
Denoising autoencoders: The model must be robust to noise



Split-Brain autoencoder



Context autoencoder



Generative Models



Training data $\sim p_{\text{data}}(x)$



Generated samples $\sim p_{\text{model}}(x)$

Want to learn $p_{\text{model}}(x)$ similar to $p_{\text{data}}(x)$

Taxonomy of Generative Models

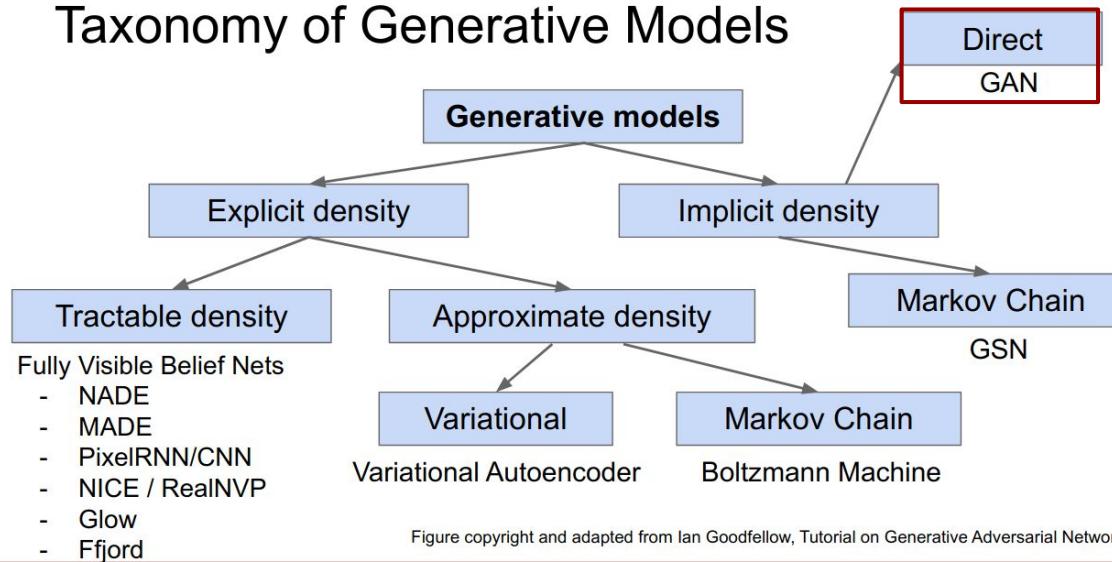


Figure copyright and adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.

Generative Adversarial Networks

PixelCNNs define tractable density function, optimize likelihood of training data:

$$p_{\theta}(x) = \prod_{i=1}^n p_{\theta}(x_i|x_1, \dots, x_{i-1})$$

VAEs define intractable density function with latent \mathbf{z} :

$$p_{\theta}(x) = \int p_{\theta}(z)p_{\theta}(x|z)dz$$

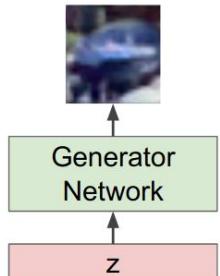
Cannot optimize directly, derive and optimize lower bound on likelihood instead

What if we give up on explicitly modeling density, and just want ability to sample?

GANs: don't work with any explicit density function!

Instead, take game-theoretic approach: learn to generate from training distribution through 2-player game

Output: Sample from training distribution

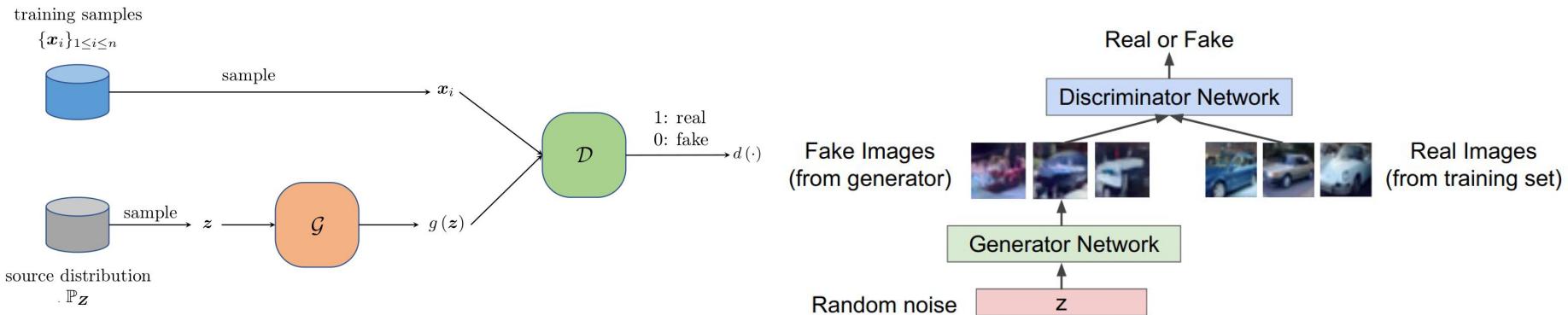


Generative Adversarial Networks

Aimed at density estimation: estimating the underlying probability density function \mathbb{P}_X from which the data is generated.

GANs are capable of estimating the density in high dimensional spaces (Images):

- GANs put more emphasis on sampling from the distribution \mathbb{P}_x than estimation.
- GANs define the density estimation implicitly through a source distribution \mathbb{P}_z and a generator function $g(\cdot)$.



Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log \underbrace{D_{\theta_d}(x)}_{\text{Discriminator output for real data } x} + \mathbb{E}_{z \sim p(z)} \log(1 - \underbrace{D_{\theta_d}(G_{\theta_g}(z))}_{\text{Discriminator output for generated fake data } G(z)}) \right]$$

Discriminator outputs likelihood in $(0,1)$ of real image

Discriminator output for generated fake data $G(z)$

Generative Adversarial Networks

for number of training iterations **do**

for k steps **do**

- Sample minibatch of m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ from noise prior $p_g(\mathbf{z})$.
- Sample minibatch of m examples $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ from data generating distribution $p_{\text{data}}(\mathbf{x})$.
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[\log D_{\theta_d}(x^{(i)}) + \log(1 - D_{\theta_d}(G_{\theta_g}(z^{(i)}))) \right]$$

end for

- Sample minibatch of m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ from noise prior $p_g(\mathbf{z})$.
- Update the generator by ascending its stochastic gradient (improved objective):

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log(D_{\theta_d}(G_{\theta_g}(z^{(i)})))$$

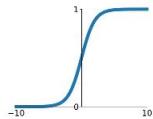
end for



Training Deep Nets

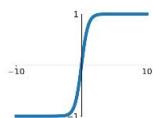
Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



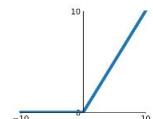
tanh

$$\tanh(x)$$



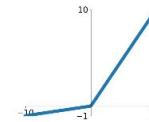
ReLU

$$\max(0, x)$$



Leaky ReLU

$$\max(0.1x, x)$$

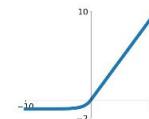


Maxout

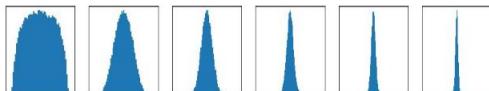
$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$

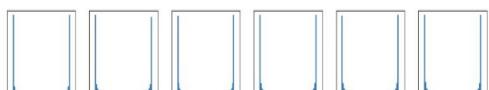


Activations



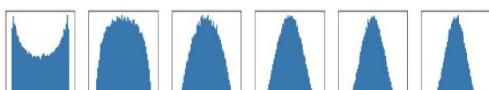
Initialization too small:

Activations go to zero, gradients also zero,
No learning =()



Initialization too big:

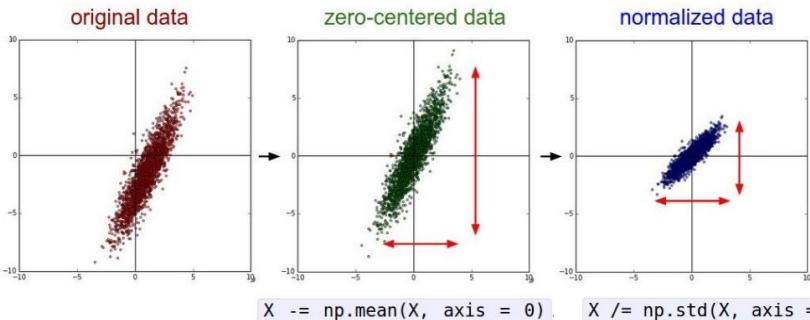
Activations saturate (for tanh),
Gradients zero, no learning =()



Initialization just right:

Nice distribution of activations at all layers,
Learning proceeds nicely

Initializing Weights



Data preprocessing

Input: $x : N \times D$

$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j} \quad \text{Per-channel mean, shape is D}$$

Learnable scale and shift parameters:

$$\gamma, \beta : D$$

Learning $\gamma = \sigma$,
 $\beta = \mu$ will recover the identity function!

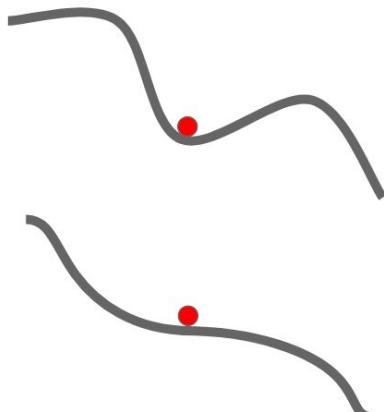
$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2 \quad \text{Per-channel var, shape is D}$$

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}} \quad \text{Normalized x, Shape is N x D}$$

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j \quad \text{Output, Shape is N x D}$$

Batch Normalization

Training Deep Nets: Optimization



SGD

$$x_{t+1} = x_t - \alpha \nabla f(x_t)$$

SGD+Momentum

$$v_{t+1} = \rho v_t + \nabla f(x_t)$$

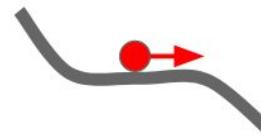
$$x_{t+1} = x_t - \alpha v_{t+1}$$

Nesterov Momentum

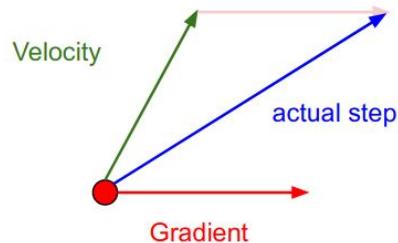
$$v_{t+1} = \rho v_t - \alpha \nabla f(x_t + \rho v_t)$$

$$x_{t+1} = x_t + v_{t+1}$$

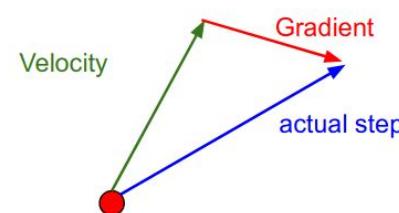
Saddle points



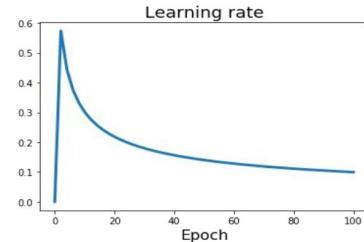
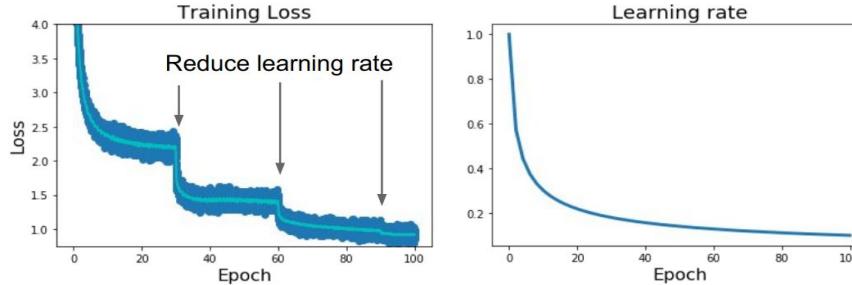
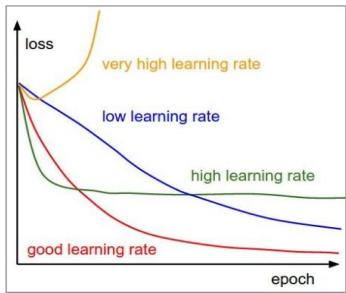
Momentum update:



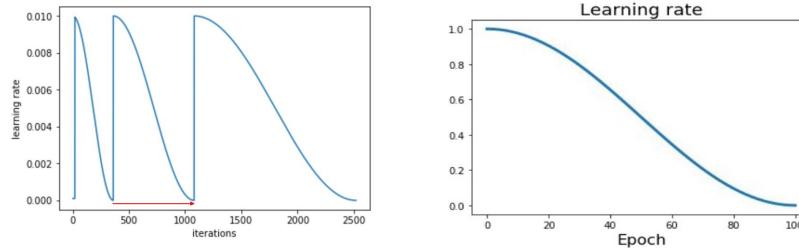
Nesterov Momentum



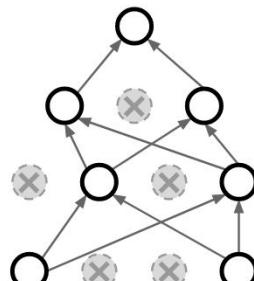
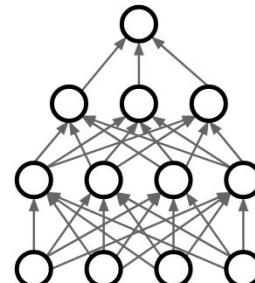
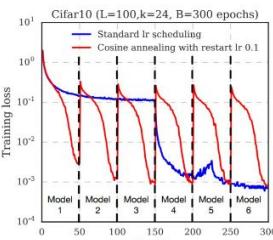
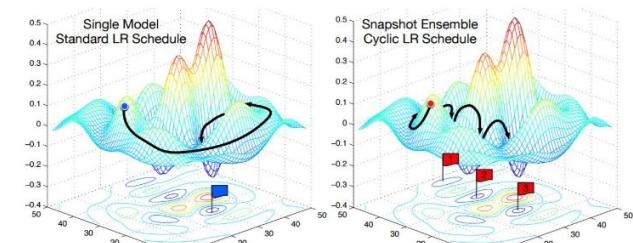
Learning rate Decay



How to find the correct learning rate ?



Model Ensembles



SGDR

Dropout

A there is a lot more

neural networks practitioner
music = loss function



Acknowledgements

- Stanford's CS231n
- The Unreasonable Effectiveness of Recurrent Neural Networks By Andrej Karpathy
- Understanding LSTM Networks By Chris Olah
- A Selective Overview of Deep Learning By Jianqing Fan, Cong Ma, Yiqiao Zhong
- Deep Learning Book by Ian Goodfellow and Yoshua Bengio and Aaron Courville
- Visualizing and Understanding Convolutional Networks By Matthew D. Zeiler and Rob Fergus
- Some Figures from the following papers: GoogleNet, ResNet, AlexNet, Split-brain autoencoders, Context autoencoder, SGDC, DC-GAN, Style-GAN.

The End.