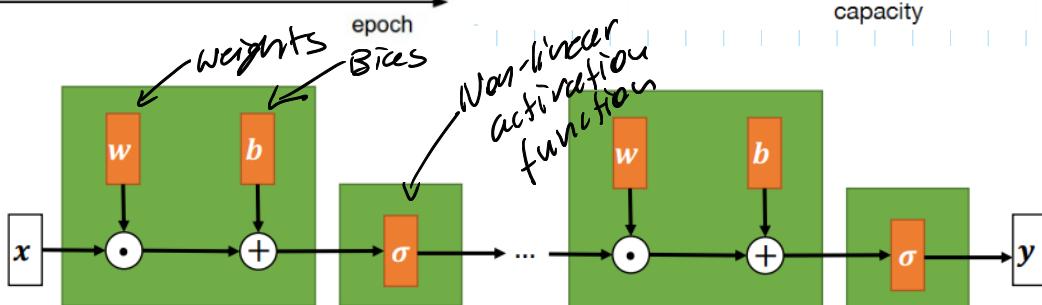
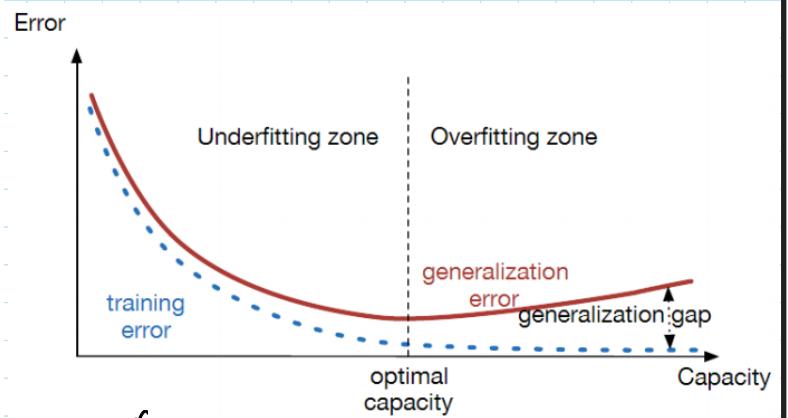
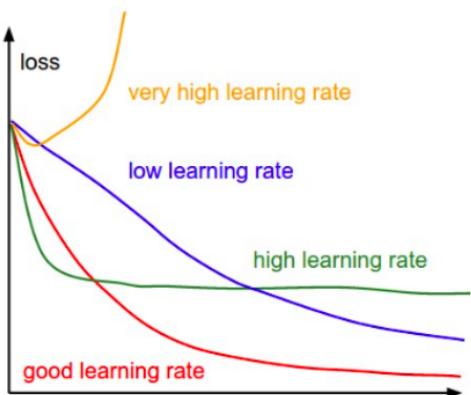


Machine Perception

- Maximum Likelihood:
 - $L(\theta) = \prod p(x_i | \theta)$
 - $\log(L(\theta)) = \sum \log p(x_i | \theta)$
 - Derive θ and set to 0

Activation function

- Must be non-linear, or the resulting function is just a linear mapping of the input
- If network is linear, we can summarize network to only one block
- Sigmoid function:
 - Logistic activation functions
 - Binary classification
 - Output between 0,1]
- Softmax function:
 - Multiclass
 - Summing all outputs = 1
- ReLU:
 - Remove negative values



$$\hookrightarrow x^{(l)} = \sigma(w^{(l)^T} \cdot x^{(l-1)} + b^{(l)})$$

- XOR example

- Use neural network
- Given two binary inputs, return XOR
- Linear activation function can only linearly separate input space

Assuming a linear function:

$$J(\Theta) = \frac{1}{(N=4)} \sum_{x \in X} (f^*(x^{(i)}) - f(x^{(i)}; \Theta))^2 \text{ with } f(x; w, b) = x^T w + b$$

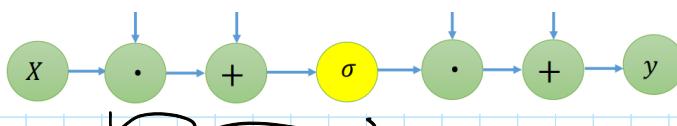
1	0
0	1

Solving via normal equation:

$$w = 0 \text{ and } b = \frac{1}{2}$$

- With non-linear activation function like ReLU we transform input space

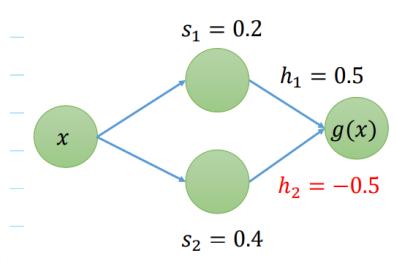
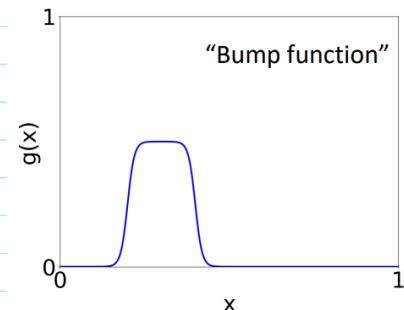
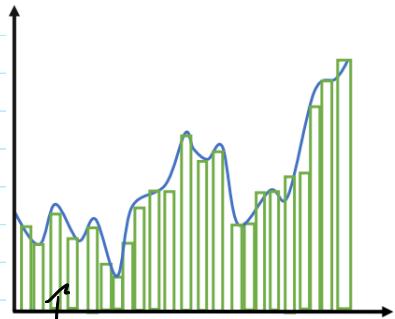
A feed-forward neural network with a single hidden layer and continuous non-linear activation function can approximate any continuous function with arbitrary precision.



Transformation of input space to get separability

- Universal Approximation theorem

- Feed forward neural network
 - One single hidden layer
 - Non-linear activation function
- } Able to approximate any continuous function with arbitrary precision



We are learning these green bumps

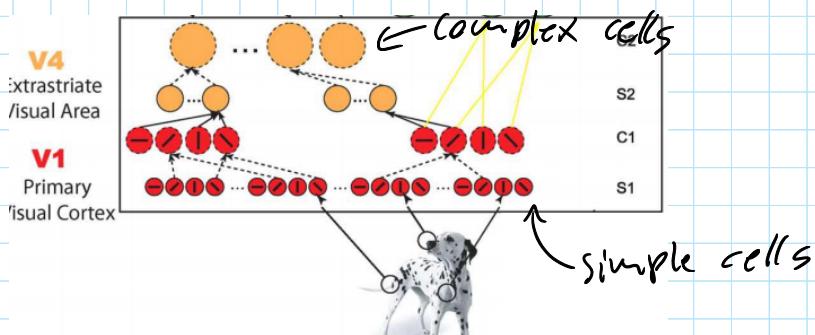
- Deeper networks work better than wider networks because we apply the activation function multiple times

Backpropagation

- Chain Rule: $\frac{\partial f(g(x))}{\partial x} = \frac{\partial f(g(x))}{\partial g(x)} \cdot \frac{\partial g(x)}{\partial x}$
- new weight(i) = old weight(i) - Learning rate $\times \left(\frac{\text{error}}{\partial \text{weight}_i} \right)$
- Weighted input: $z = x \cdot w \Leftrightarrow z'(x) = w / z'(w) = x$
- ReLu: $R = \max(0, z) \Leftrightarrow R'(z) = \begin{cases} 0 & z \leq 0 \\ 1 & z > 1 \end{cases}$
- Loss: $C = \frac{1}{2} (\hat{y} - y)^2 \Leftrightarrow C'(\hat{y}) = (\hat{y} - y)$
- If there are multiple paths to the weight, we must add all derived paths

Convolutional Neural Network

- Simple cells:
- Large cells
- Come first in network
- Complex cells:
- Combine simple cells output



CNN: learn filters

Given a transform T , a function f , defining a translation or shift operation and input vectors u, v scalars α, β

Definition: Transform T is *linear* if,

$$T(\alpha u + \beta v) = \alpha T(u) + \beta T(v)$$

Definition: A transform T is *invariant to f* if, *shift operation, e.g. shift all by one to right*

$$T(f(u)) = f(T(u))$$

Definition: A transform T is *equivariant to f* ,

$$T(f(u)) = f(T(u))$$

Any linear, shift-equivariant transform T can be written as a convolution

Correlation

$$I'(i, j) = \sum_{m=-k}^k \sum_{n=-k}^k K(m, n) I(i+m, j+n)$$

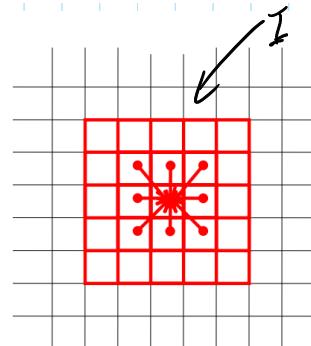
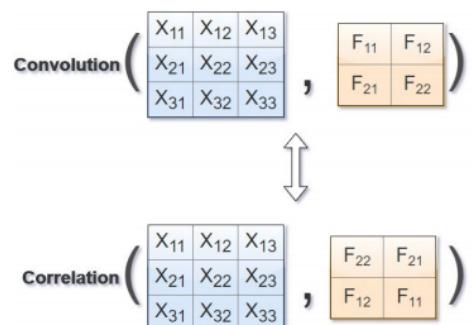
Don't really matter, because we learn the filters

Convolution

$$\begin{aligned} I'(i, j) &= \sum_{m=-k}^k \sum_{n=-k}^k K(m, n) I(i-m, j-n) \\ &= \sum_{m=-k}^k \sum_{n=-k}^k K(-m, -n) I(i+m, j+n) \end{aligned}$$

So if $K(i, j) = K(-i, -j)$, then Correlation == Convolution

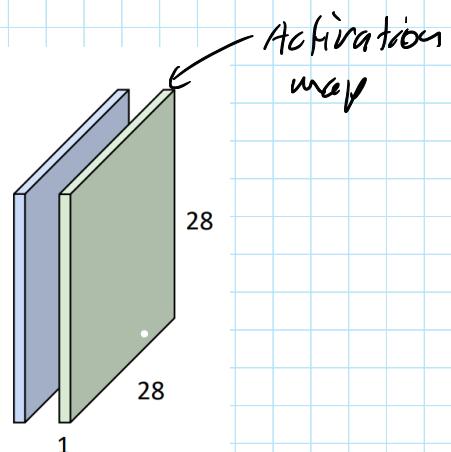
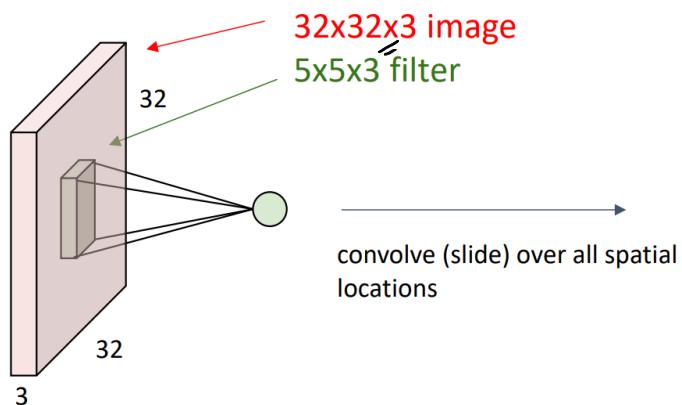
$$O(i, j) = \sum_{m=-k}^k \sum_{n=-k}^k F(m, n) X(i-m, j-n)$$



Kernel

K(-1,-1)	K(0,-1)	K(1,-1)
K(-1,0)	K(0,0)	K(1,0)
K(-1,1)	K(0,1)	K(1,1)

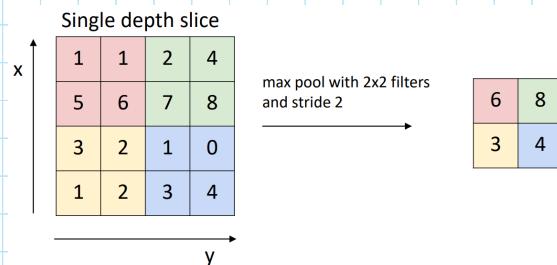
$$\begin{aligned} I'(x, y) = & K(1,1)I(x-1, y-1) + K(0,1)I(x, y-1) + K(-1,1)I(x+1, y-1) \\ & + K(1,0)I(x-1, y) + K(0,0)I(x, y) + K(-1,0)I(x+1, y) \\ & + K(-1,1)I(x-1, y+1) + K(0,1)I(x-1, y) + K(-1,-1)I(x+1, y+1) \end{aligned}$$



Pooling layer

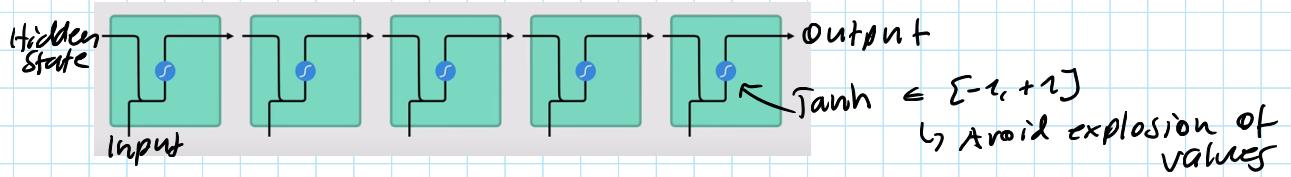
- Downscales image

- Operates on each activation map independently



- Dense layer : - last step before classification
- Aggregate information

- Recurrent neural network

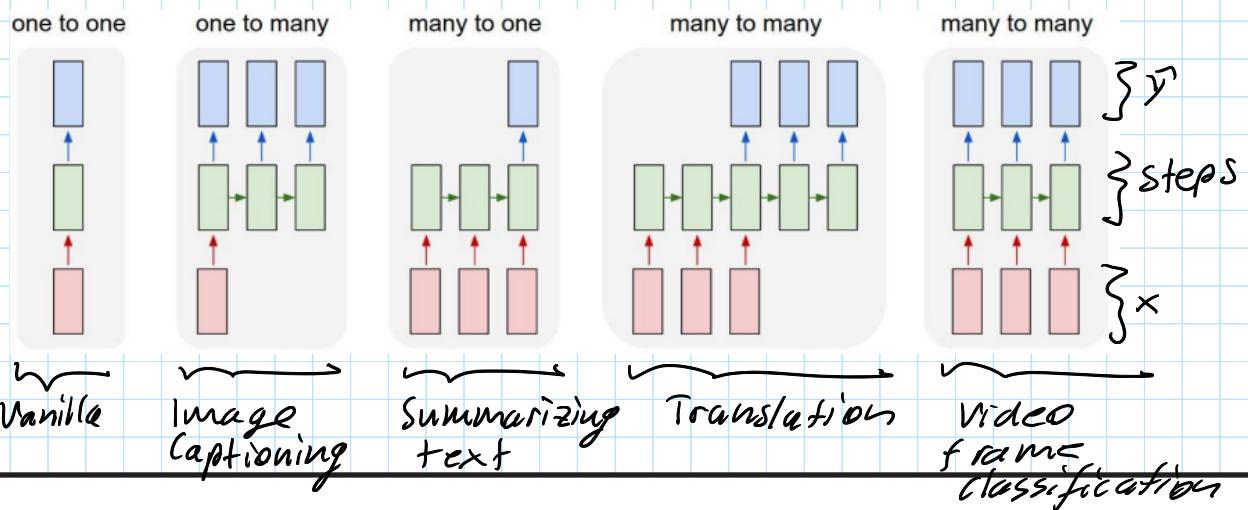


- Suffers from vanishing gradient problem
 - ↳ Gradient is getting smaller during back propagation through RNN
 - ↳ If gradient is small, layer doesn't learn
 - ↳ If earlier layer can't learn, RNNs forget what it has seen in the long term
 - ↳ "short term memory effect"
- Hidden state flows from one step to the next
- RNN work sequentially, one info at a time
- Hidden state is initialized as zero vector
- At each run, the hidden state is modified with the new input

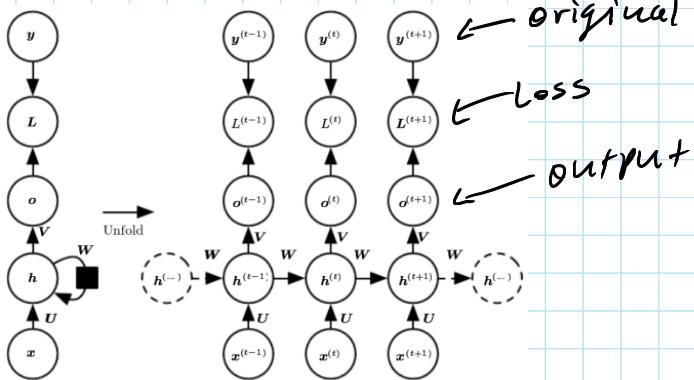
$$\cdot h^+ = \tanh \left(\underbrace{W_{hh} \cdot h^{+1}}_{\text{Hidden state}} + \underbrace{W_{xh} \cdot x^+}_{\text{Current input}} + b \right)$$

$$\cdot \hat{y}^+ = \text{softmax}(W_{hy} \cdot h^+ + c \leftarrow \text{Bias})$$

$$\cdot \text{Loss: } L = \sum_{t=1}^T L(y^+, \hat{y}^+)$$



- Parameters are shared between time-steps



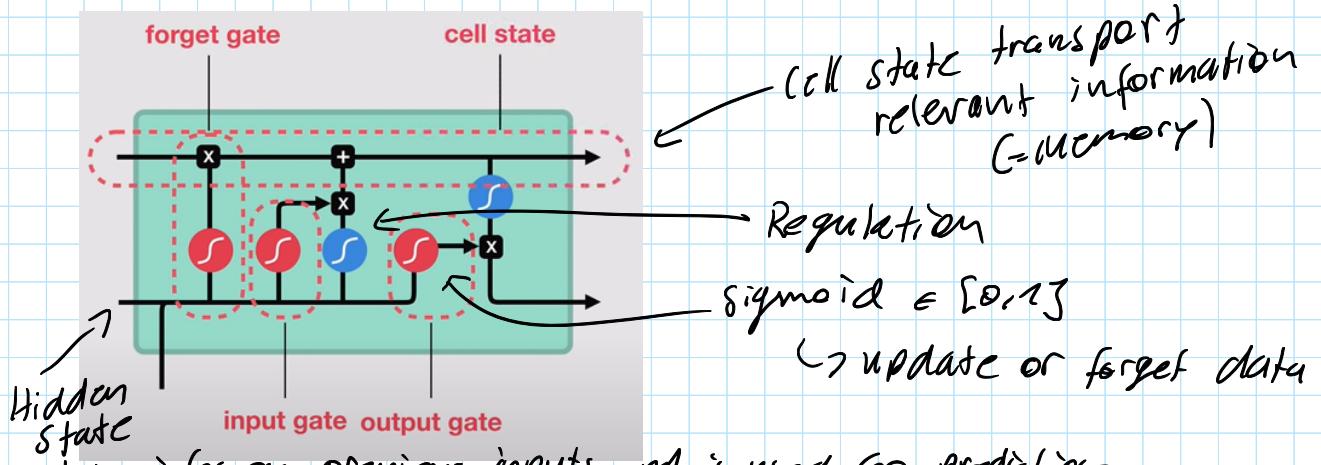
$$\frac{\partial L^t}{\partial W} = \sum_{k=1}^t \frac{\partial L^t}{\partial y^t} \frac{\partial y^t}{\partial h^t} \frac{\partial h^t}{\partial h^k} \frac{\partial h^k}{\partial W}$$

$$\frac{\partial h^t}{\partial h^k} = \prod_{i=k+1}^t \frac{\partial h^i}{\partial h^{i-1}}$$

\hookrightarrow Exploding/vanishing gradient
gradient clipping LSTM or GRU

- Long Short Term Memory (=LSTM)

- More evolved RNNs that doesn't suffer from vanishing gradient



- Gates are also neural networks that learn, what info to keep or forget

- forget gate: what info should be forgotten (from the previous cell state ($0 \Rightarrow$ forget))

- input gate: Decides how much new input we should take ($0 \Rightarrow$ not important)

- output gate: How much memory affects the output

$$c^t = \underbrace{\text{forget-layer}(h^{t-1} + x^t)}_{\text{forget gate output}} \cdot c^{t-1} + \underbrace{\text{input-layer}(h^{t-1} + x^t)}_{\text{input gate output}} \cdot (h^{t-1} + x^t)$$

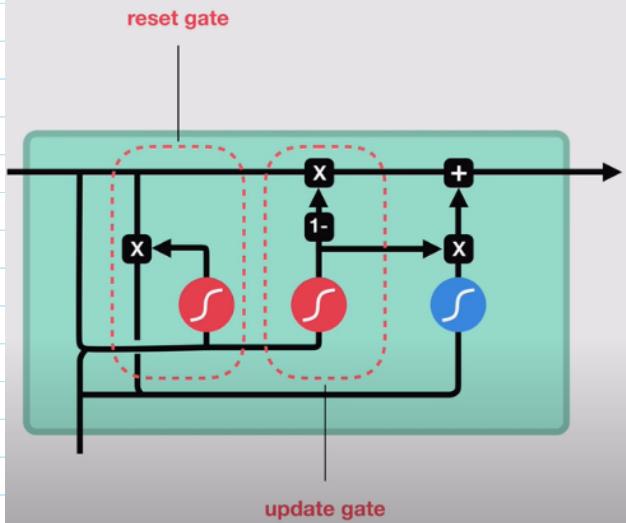
$$h^t = \underbrace{\text{output-layer}(h^{t-1} + x^t)}_{\text{output gate output}} \cdot \tanh(c^t)$$

- Gated recurrent unit (=GRU)

- Cell state is stored in hidden state
- Fewer parameters than LSTM
- It lacks output gate, because the memory is already in the hidden state

- The internal state of a standard RNN is entirely deterministic

↳ They are often augmented with random latent variables to increase modelling capacity and to better capture uncertainty



Object detection

- As regression problem: $\text{Dog}(x, y, w, h)$
- Need to test many positions and scales
- As classification: 1. Find "blobby" image regions
2. classify regions of interest

Semantic segmentation: Label each pixel

1. Extract patch

2. Run through CNN

- Get smaller output

picture due to Pooling

Convolutions at original image resolution is too expensive

Downsampling → Upsampling → Output

Strided convolutions Unpooling

Max-pooling

Unpooling

Nearest Neighbor	
1	2
3	4
1	1
3	3
2	2
4	4
3	3
3	3

Input: 2 x 2 Output: 4 x 4

"Bed of Nails"	
1	2
3	4
1	0
0	0
0	0
3	0
0	0
0	0

Input: 2 x 2 Output: 4 x 4

Max Pooling	
Remember which element was max!	
1	2
3	5
1	2
7	3
6	3
2	1
4	8

Input: 4 x 4 Output: 2 x 2

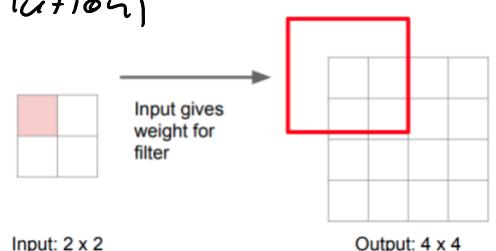
Max Unpooling	
Use positions from pooling layer	
1	2
3	4
0	0
0	1
0	0
0	0
3	0
0	0
0	4

Input: 2 x 2 Output: 4 x 4

Remembers, which element was max while max-pooling

Learnable upsampling (=Deconvolution)

- Use transpose convolution
- Learn filter



- Regularization

- Any changes, which improve generalization error
- Avoid overfitting
- Constraining regularizing coefficients towards zero to get a more flexible model
- Manipulate data to avoid overfitting
 - Data standardization: better convergence

$$X_s = \frac{X - \mu}{\sigma}$$

Where $\mu = \frac{1}{n} \sum_{i=1}^n x_i$ and $\sigma = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \mu)^2}$

- Data augmentation: generate synthetic data by transforming existing samples
 - add noise
- Manipulate model to avoid overfitting
 - Penalize big weights:
 - $L(\theta) = L(\theta) + \alpha R(\theta)$
 - L1: Lasso: sparse solution
 - L2: ridge: shrinkage
 - Dropout: get rid of some connections

Dropout	Bagging
Models share weight	Models are independent
Trains partly only small percentage of its models	Trains all models until convergence

- pretraining: Train first on related task B with large dataset, finetune on task-specific dataset
- Use pretrained network-architecture

- Generative Modelling

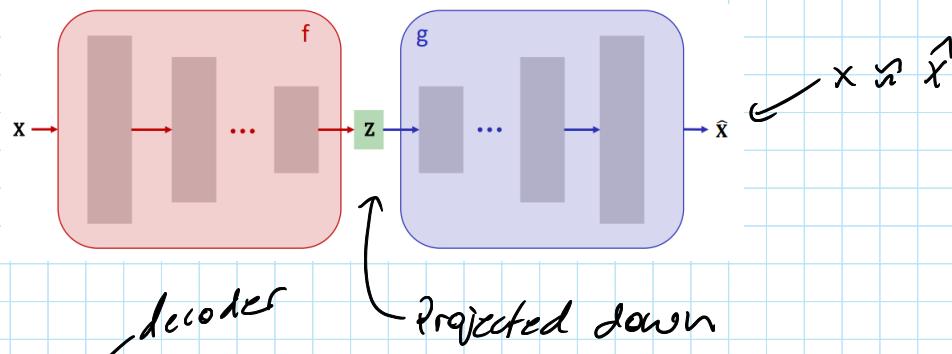
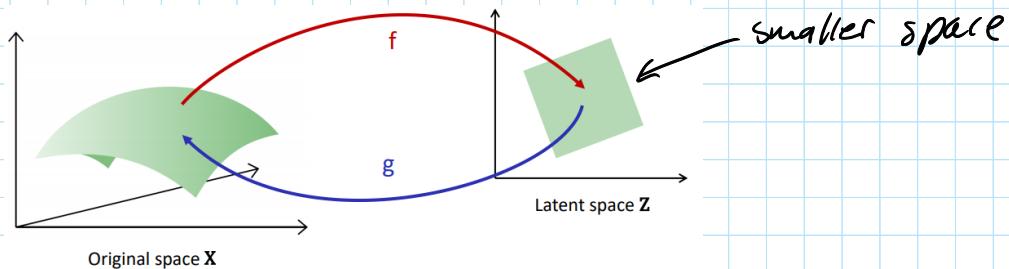
- Given training data, generate new samples, drawn from training distribution
 $\hookrightarrow p_{\text{model}}(x) \approx p_{\text{data}}(x)$

- Autoencoder

- Transform high dimensional input to low-dim.
so that each dimension encodes

Ausschensmerkmal

- Consists of **encoder** / **decoder**



$$\hat{\theta}_f, \hat{\theta}_g = \underset{\theta_f, \theta_g}{\operatorname{argmin}} \sum_n^N \|x_n - \hat{x}_n\|^2$$

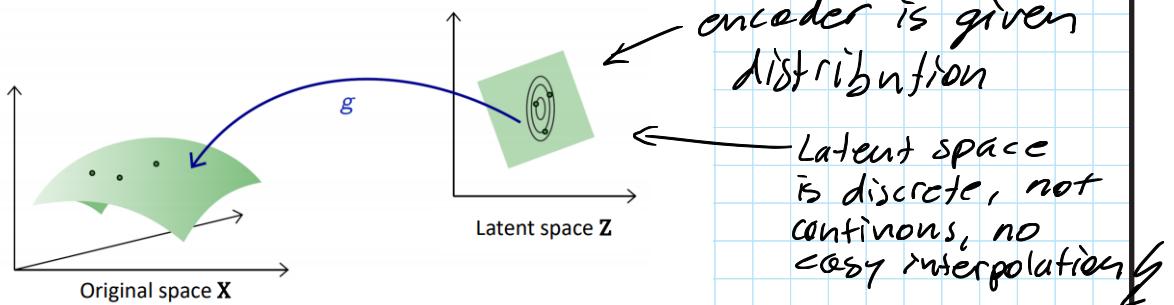
$$\hat{\theta}_f, \hat{\theta}_g = \underset{\theta_f, \theta_g}{\operatorname{argmin}} \sum_n^N \|x_n - g(f(x_n))\|^2$$

encoder

- Used for human motion prediction
- By introducing a simple density function (e.g. Gaussian) over the latent space z , the decoder g can generate new data

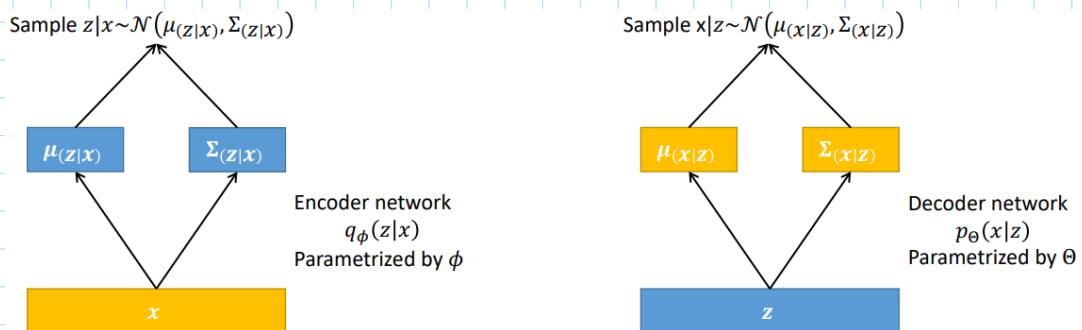
$$f(x) \sim N(\hat{\mu}, \hat{\sigma} I)$$

where $\hat{\mu}$ and $\hat{\sigma}$ are mean and standard deviation estimated from training data.

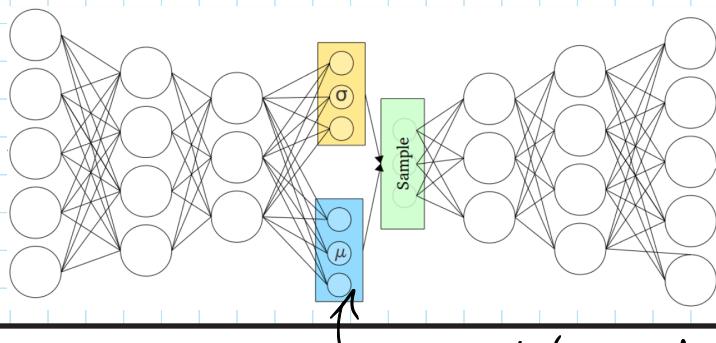


Variational Autoencoder (=VAEs)

- Probabilistic version of autoencoder \Rightarrow continuous latent space
- generate data with natural variation
- We want to estimate parameter Θ^*
- $\hookrightarrow p_\theta(x) = \int_z p_\theta(x|z) \cdot p_\theta(z) \cdot dz$
- We use additional encoder network $q_\phi(z|x)$ to approximate $p_\theta(z|x)$
- $\hookrightarrow p_\theta(z|x) = \frac{p_\theta(x|z) \cdot p_\theta(z)}{p_\theta(x)}$ $\Rightarrow p_\theta(x) = \frac{p_\theta(x|z) \cdot p_\theta(z)}{p_\theta(z|x)}$ prior

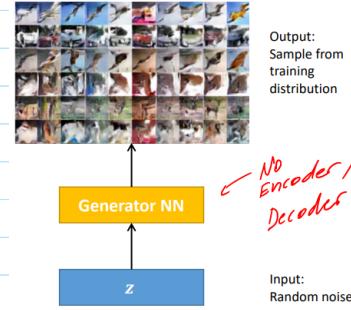
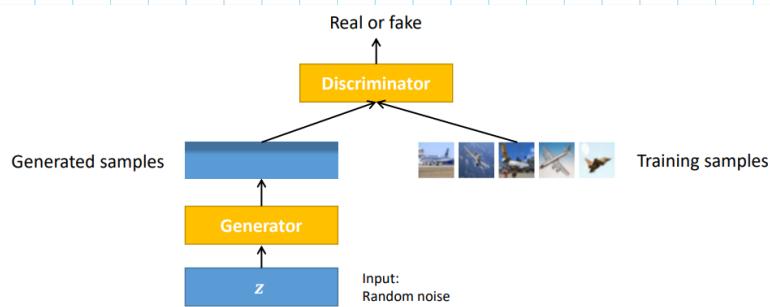


- Generates blurry image
- Used for generating Data, z is sampled from prior



Learned by encoder, sampled from by decoder

- Generative Adversarial Networks (GANs)



- Two networks play a two player game
- generator: - tries to fool discriminator
 - $G: \mathbb{R}^D \rightarrow X$
- Discriminator: - tries to distinguish images
 - $D: X \rightarrow [0, 1] \quad (0 = \text{fake})$
 - If $D = 0.5$ every time, generator has won

Assuming G is fixed, to train D given a set of real samples:

$$x^n \sim p_d, \quad n = 1, \dots, N$$

We generate

$$z^n \sim \mathcal{N}(0, 1), \quad n = 1, \dots, N$$

and form a training dataset

$$\mathbb{D} = \{(x^{(1)}, 1), \dots, (x^{(n)}, 1), (G(z^{(1)}), 0), \dots, (G(z^{(n)}), 0)\}$$

- Difficult to train: Nash-Equilibrium in two-player-game
 1. Make multiple steps in generator
 2. Only make one step in discriminator

Objective:

$$\min_{\Theta_g} \max_{\Theta_d} \mathbb{E}_{x \sim p_d(x)} [\log D_{\Theta_d}(x)] + \mathbb{E}_{z \sim p_z(z)} [\log (1 - D_{\Theta_d}(G_{\Theta_g}(z)))]$$

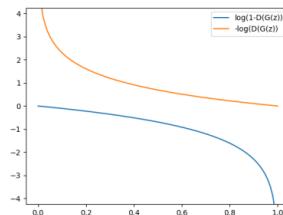
Alternate between:

1. Gradient ascent on D

$$\max_{\Theta_d} \mathbb{E}_{x \sim p_d(x)} [\log D_{\Theta_d}(x)] + \mathbb{E}_{z \sim p_z(z)} [\log (1 - D_{\Theta_d}(G_{\Theta_g}(z)))]$$

2. Instead: gradient ascent on G:

$$\max_{\Theta_g} \mathbb{E}_{z \sim p_z(z)} [\log (D_{\Theta_d}(G_{\Theta_g}(z)))]$$



- Cycle GAN

- Input \rightarrow Output $\xrightarrow{\text{Input}}$
- Like in language, translate DE \rightarrow EN and then EN \rightarrow DE to check, if the meaning is preserved

- Mode collapse

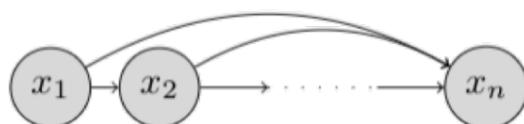
- Generator collapses and produces only limited varieties of samples
- Generator focuses only on one mode



- Autoregressive Models

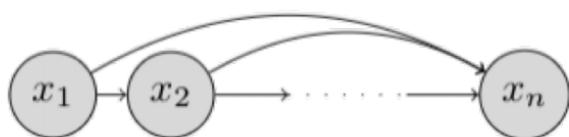
- Regression model uses data from same input variable at previous time steps, it is referred to as autoregression
 - Generative model models joint distribution $P(x, y)$
 - Use x_1, x_2, x_3 to generate x_4
 - Use x_2, x_3, x_4 to generate x_5
- } like a repeat

$$p(x) = \prod_1^n p(x_i | x_1, \dots, x_{i-1}) = \prod_1^n p(x_i | x_{<i})$$



} Attempt 1:
Tabular approach

$$p_{\theta_i}(x_i | x_{<i}) = \text{Ber}(f(x_1, \dots, x_{i-1}))$$



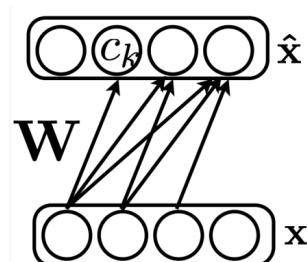
} Attempt 2

Fully Visible Belief Networks

$$\hat{x}_i = p(x_i = 1 | x_{<i})$$

Modelled via logistic regression:

$$f_i(x_1, x_2, \dots, x_{i-1}) = \sigma(a_0^i + a_1^i x_1 + \dots + a_{i-1}^i x_{i-1})$$



Neural Autoregressive Density Estimator (NADE)

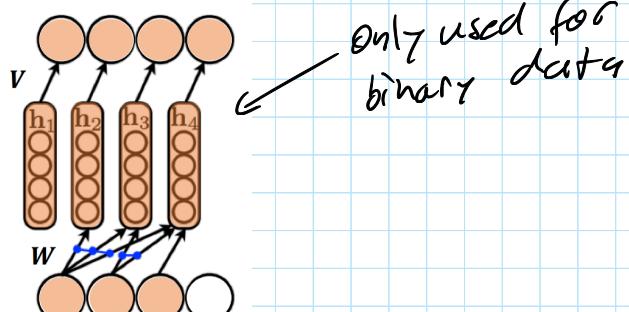
autoencoder-like neural network to learn $p(x_i = 1 | x_{<i})$:

$$\begin{aligned} h^i &= \sigma(b + W_{<i} x_{<i}) \\ \hat{x}_i &= \sigma(c_i + V_{i,*} h^{(i)}) \end{aligned}$$

each conditional is modeled by the same neural network

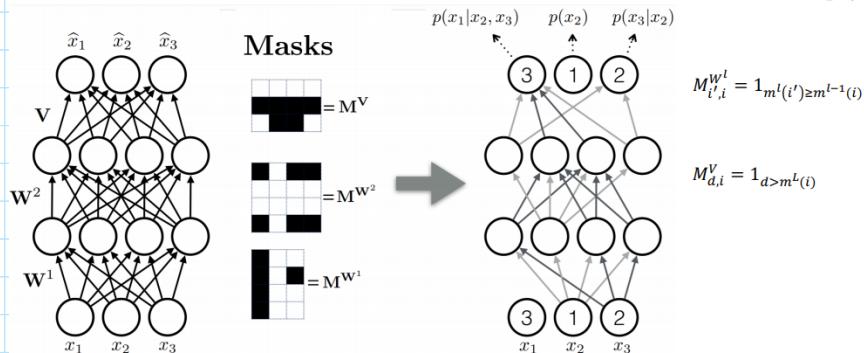
We can leverage the fact that:

$$(b + W_{<i+1} x_{<i+1}) - (b + W_{<i} x_{<i}) = W_{i,i+1} x_{i+1}$$



- Made (=Masked Autoencoder Distr. Estimator)
 - No computational path between output x_d and any of the input units must exist
 - Use masks to remove some connections

Constrain Autoencoder such that output can be used as conditionals $p(x_i|x_{<i})$

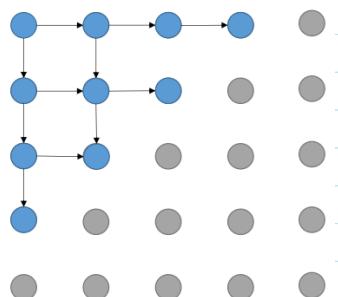


PixelRNN [van der Oord et al. 2016]

Generate image pixels starting from corner

Dependency on previous pixels modeled using an RNN (LSTM)

Issue: Sequential generation is slow – due to explicit pixel dependencies



} faster, but still sequentially!

PixelCNN [van der Oord et al. 2016b]

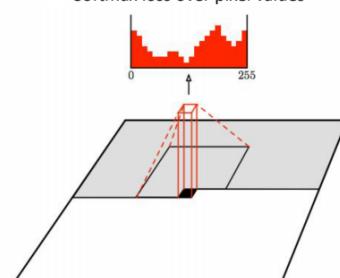
Still generate image pixels starting from corner

Dependency on previous pixels now modeled using a CNN over context region

Training maximizes likelihood of training images

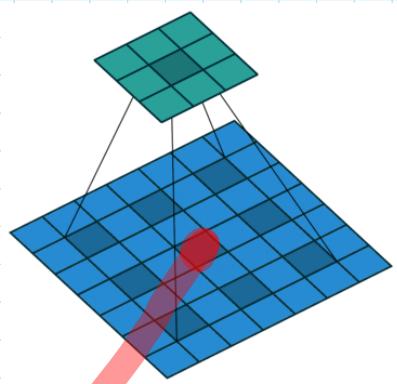
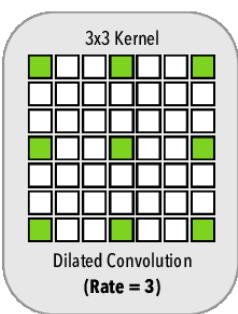
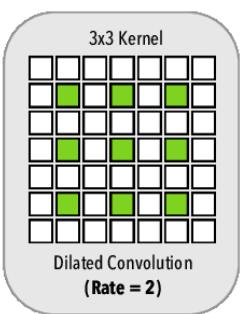
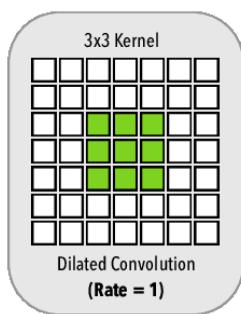
$$p(x) = \prod_1^n p(x_i|x_1, \dots, x_{i-1})$$

Softmax loss over pixel values

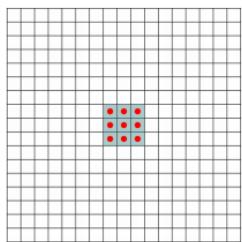


[image source: van der Oord et al.]

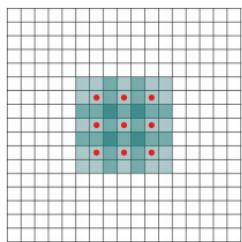
- Dilated Convolution



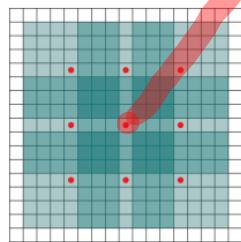
Receptive field size can increase exponentially with dilated layers



Layer 1



Layer 2

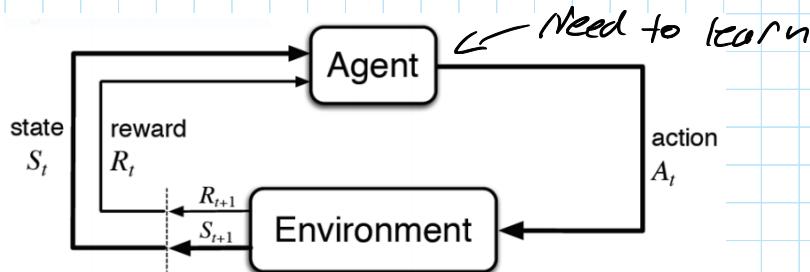


Layer 3

- Red: Foot print of one dilated convolution
Shaded: Receptive field of convolution
Darker: Receptive field overlap between foot prints

- Reinforced Learning

- Learn how to act
- Maximize numerical award given



- We model environment as Markov Decision process

An MDP consists of $(\mathcal{S}, \mathcal{A}, \mathcal{R}, \mathbb{P}, \gamma)$

- a set of states \mathcal{S}
- a set of actions \mathcal{A}
- a reward function $r: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ (deterministic, can also be a distribution)
- a transition function $p: \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$ (deterministic, can also be a distribution)
- initial state $s_0 \in \mathcal{S}$
- Discount factor: γ

- $V_\pi: \mathcal{S} \rightarrow \mathbb{R}$ of a policy π returns the expected cumulative return for each state s when following policy π

Dynamic Programming (=Tabular Learning => Tab[le])

- We need whole transition matrix

(\hookrightarrow Model of the world ($=MDP$))

- Value iteration:
 1. compute optimal values for each state V_π
 2. obtain policy π^* based on V_π

- Policy Iteration:
 1. Initialize arbitrary policy π^*
 2. Repeat until policy doesn't improve
 - 2.1. evaluate value function V_π
 - 2.2. Improve policy

- 2.1. evaluate value function V_π
- 2.2. Improve policy

Temporal Difference Learning

For each step with action a from s to s' that we take with our policy, we compute the difference to our current estimate and update our value function:

$$\Delta V(s) = r(s, a) + \gamma V(s') - V(s)$$

$$V(s) \leftarrow V(s) + \alpha \Delta V(s)$$

Anpassen des
expected returns

- Random policy: - choose action randomly
 - Exploration, find unknown states
- Greedy policy: - choose best action
 - Can get stuck in local min.
 - Exploration
- To get balance between exploration/exploitation
we use epsilon-greedy-strategy
 - ↳ $\epsilon=1 \Rightarrow$ agent explores at start
 - ↳ Through time, agent will get greedy

- Q-Learning ($Q \hat{=} \text{Quality value}$)

For each step with action a from s to s' that we take with our policy we compute the difference to our current estimate and update our value function like this:

$$\Delta Q(S, A) = R_{t+1} + \gamma \max_a \{Q(S', a)\} - Q(S, A)$$

$$Q(S, A) \leftarrow Q(S, A) + \alpha \Delta Q(S, A)$$

Get Δ

Learning rate

use this to
update the
fields

- Random policy for initialization
- Use initial Q-function and update it

by exploration/immediate reward

$$Q(s, a) = r(s, a) + \gamma \max_a Q(s', a)$$

The Bellman Equation

Reward from new
state s'

The above equation states that the Q Value yielded from being at state s and selecting action a , is the immediate reward received, $r(s, a)$, plus the highest Q Value possible from state s' (which is the state we ended up in after taking action a from state s). We'll receive the highest Q Value from s' by choosing the action that maximizes the Q Value. We also introduce γ , usually called the discount factor, which controls the importance of longterm rewards versus the immediate one.

Not done yet — our greedy algorithm has a serious problem: if you keep selecting the same best-actions, you'll never try anything new, and you might miss a more rewarding approach just because you never tried it.

To solve this, we use an ε -greedy approach: for some $0 < \varepsilon < 1$, we choose the *greedy action* (using our table) with a probability $p = 1-\varepsilon$, or a *random action* with probability $p = \varepsilon$. We thus give the Agent a chance to *explore* new opportunities.

This algorithm is known as *Q Learning* (or *Q-Table*). Congratulations! you just learned your very first Reinforcement Learning algorithm!

- Deep Reinforced Learning (-Deep Q-Networks)
- We need to learn • Policy: $\pi: S_t \rightarrow A_t$
• Value: $V: S_t \rightarrow V(S_t)$
- ~ Used if state space is large

Enter Deep Learning! We combine Q Learning and Deep Learning, which yields *Deep Q Networks*. The idea is simple: we'll replace the the Q Learning's table with a neural network that tries to approximate Q Values. It is usually referred to as the *approximator* or the *approximating function*, and denoted as $Q(s, a; \theta)$, where θ represents the trainable weights of the network.

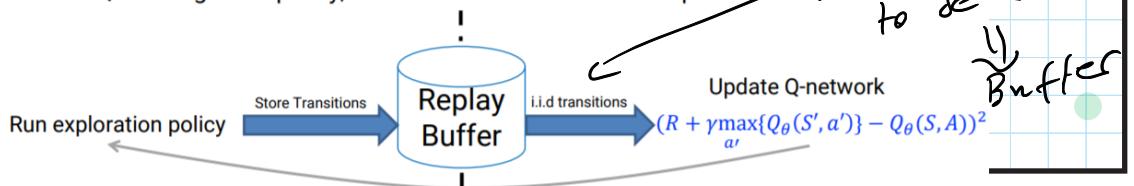
Moving on: Training. In Reinforcement Learning, the training set is created as we go; we ask the Agent to try and select the best action using the current network — and we record the *state, action, reward* and the *next state* it ended up at. We decide on a batch size b , and after every time b new records were recorded, we select b records *at random* (!!) from the memory, and train the network. The memory buffers used are usually referred to as *Experience Replay*. Several types of such memories exist — one very common is a cyclic memory buffer. This makes sure the Agent keeps training over its new behavior rather than things that might no longer be relevant.

The Q-learning updates reduce to **stochastic gradient descent** on

$$\text{Loss}(\theta) = (R + \gamma \max_{a'} \{Q_\theta(S', a')\} - Q_\theta(S, A))^2$$

How can we address this?

- Use a **replay buffer** to store generated samples
- During updates, randomly sample transitions from the buffer
- Since Q-learning is off-policy, we are allowed to use old samples



Sigmoid (Logistic)

- + Has a nice interpretation as the firing rate of a neuron: from not firing at all (0) to fully-saturated firing at an assumed maximum frequency (1).
- + Finite range (0,1): stable training.
- Saturates towards either end and large activations will get very small gradients.
- Not zero-centered. Always positive.



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

$$f'(x) = f(x)(1 - f(x))$$

Hyperbolic Tangent (tanh)

Scaled sigmoid function

$$\tanh(x) = 2\sigma(2x) - 1$$

Like the sigmoid function

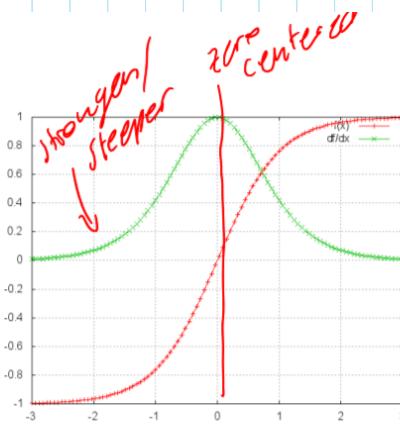
- + has finite range (-1,1),

- its activations saturate,

but unlike the sigmoid function

- + its output is zero-centered,

- + has stronger and steeper gradients.

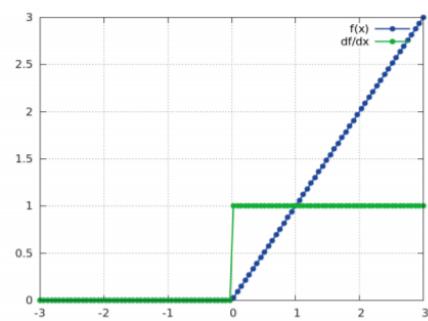


$$f(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$f'(x) = 1 - f(x)^2$$

Rectified Linear Unit (ReLU)

- range is [0, inf): it can blow up the activation and destabilize training.
- + range is [0, inf): It was presented to greatly accelerate the convergence.
- + inexpensive operation compared to sigmoid and tanh.
- + sparsity of the activations.
- “dying” ReLU problem: units with negative activations get no update.



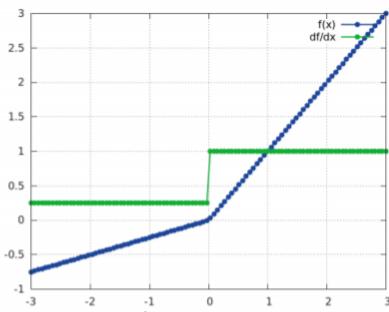
$$f(x) = \max(0, x)$$

$$f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$$

Leaky ReLU

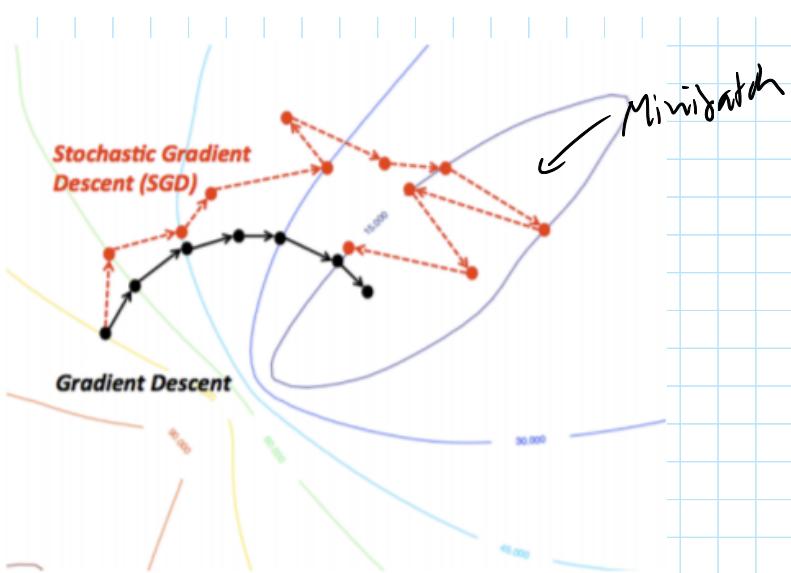
+ Attempts to solve "dying" ReLU by simply introducing a small negative slope for negative activations.

α is a small positive constant.



$$f(x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$$

$$f'(x) = \begin{cases} \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$$



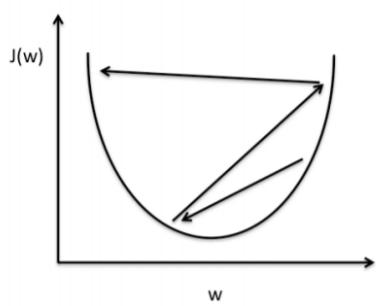
Optimization: AdaGrad (= Adaptive gradient Descent)

Adagrad modifies the general learning rate η at each time step t for every parameter θ_i based on the past gradients that have been computed for θ_i

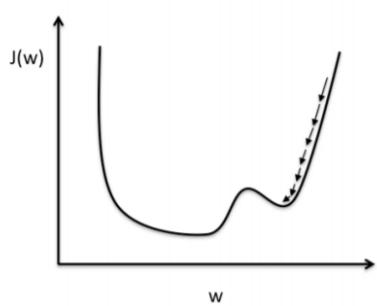
$$g_{t,i} = \nabla_{\theta} J(\theta_i)$$

$$\theta_{t+1,i} = \theta_{t,i} - \frac{\eta}{\sqrt{G_{t,ii} + \epsilon}} \cdot g_{t,i}$$

\curvearrowleft Avoids overshooting



Large learning rate: Overshooting.



Small learning rate: Many iterations until convergence and trapping in local minima.