

Neural Networks and Deep Learning

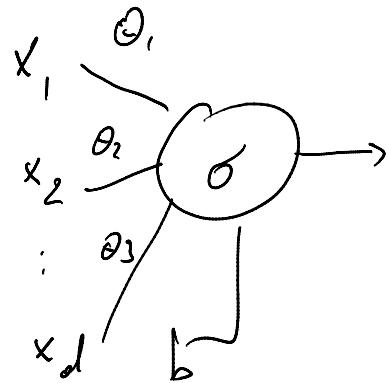
Lecture 3 & 4

Jan Chorowski
Instytut Informatyki
Wydział Matematyki i Informatyki Uniwersytet
Wrocławski
2019-2020

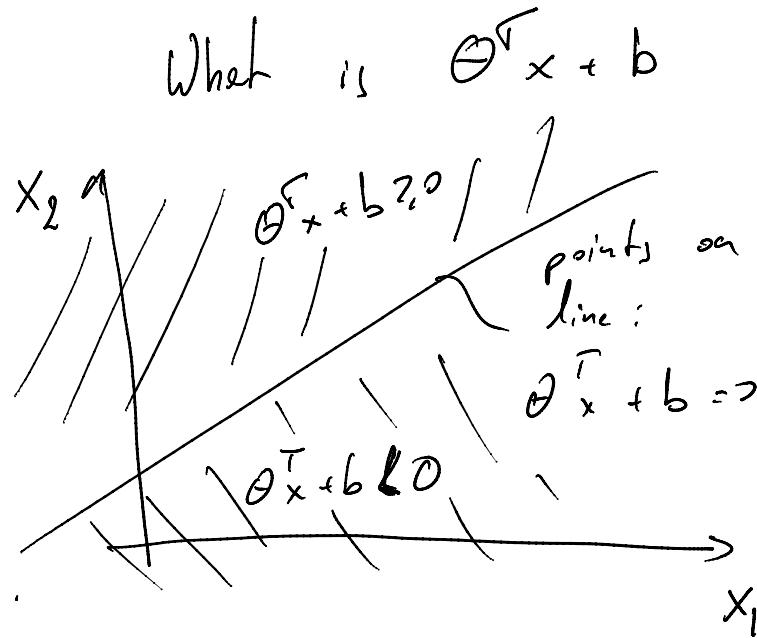
DEPTH MATTERS

Single neuron

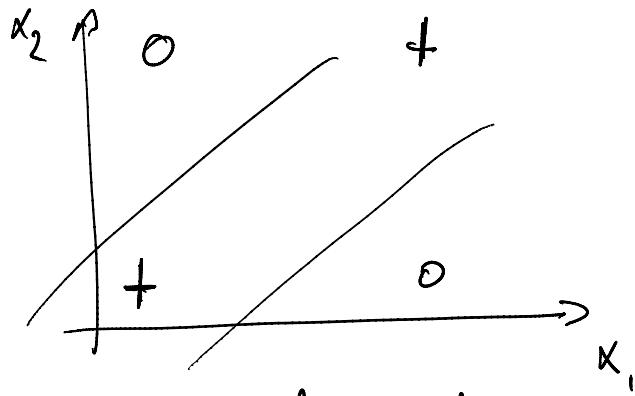
$$f(x) = \sigma(\theta^T x + b)$$



↳ can split space
with 1 hyper-plane.



The XOR Problem



Goal: separate
+ from o

need 2 lines to
do it ?

=> Can't solve with logistic
regression / one hidden neuron

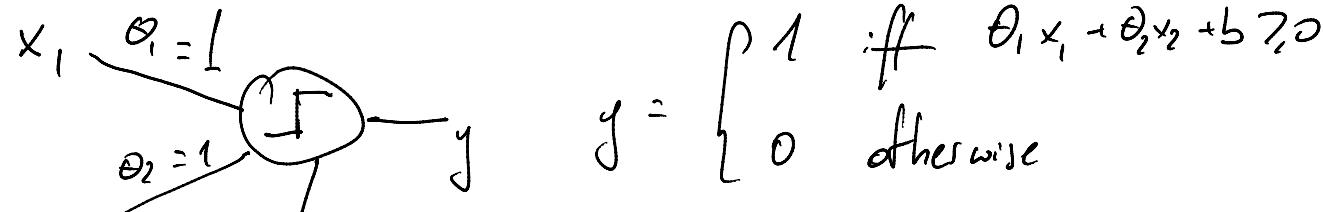
The old approach: hand-design data transformation

We want $x_1 \text{ XOR } x_2$
Intro transformation $\phi\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) \rightarrow \begin{bmatrix} x_1 \text{ OR } x_2 \\ x_1 \text{ AND } x_2 \end{bmatrix}$

$$x_1 \text{ XOR } x_2 \equiv (x_1 \text{ OR } x_2) \text{ AND } \text{Not}(x_1 \text{ AND } x_2)$$

Logic networks

We need : • OR • AND • NOT

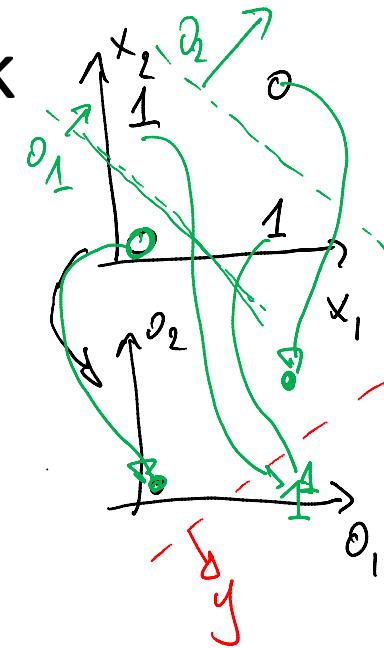
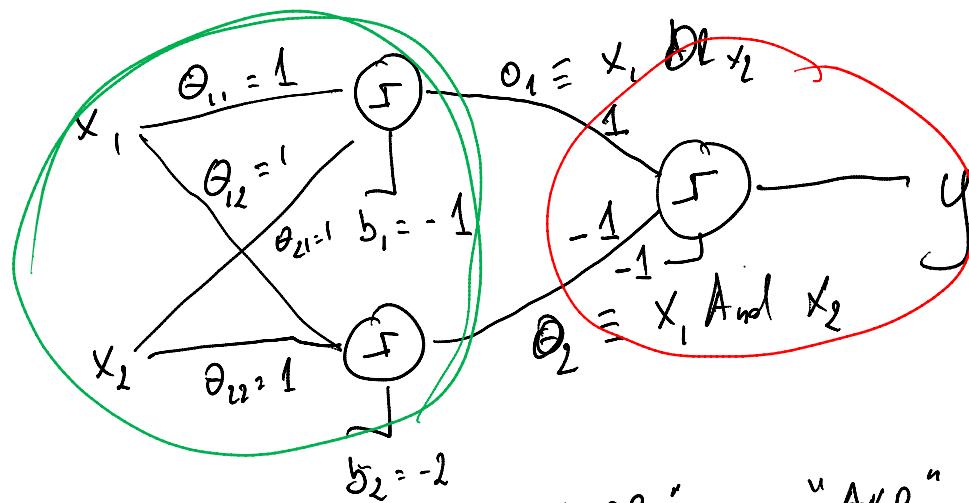


$$y = \begin{cases} 1 & \text{if } \theta_1 x_1 + \theta_2 x_2 + b \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

A truth table for the OR function. The columns are labeled x_1 , x_2 , and y . The rows show all combinations of x_1 and x_2 : (0,0), (0,1), (1,0), (1,1). The output y is 0 for (0,0) and 1 for all other rows. To the right of the table, the text "x₁ OR x₂" is written next to a bracket under the last three rows. To the far right, there is a vertical column of green numbers [0, 1, 0, 1] with the text "x₁ AND x₂" written vertically next to it.

x_1	x_2	y
0	0	0
0	1	1
1	0	1
1	1	1

Full XOR network



$b_2 = -1$		"OR"		"AND"		y
x_1	x_2	o_1	o_2	y		
0	0	0	0	0	0	
0	1	1	0	0	1	
1	0	1	0	0	1	
1	1	1	1	1	0	

$x_1 \oplus x_2$

Hidden layer: a really big deal: Universal Approximation

Thm: Any Boolean Function HAS A DNF FORM

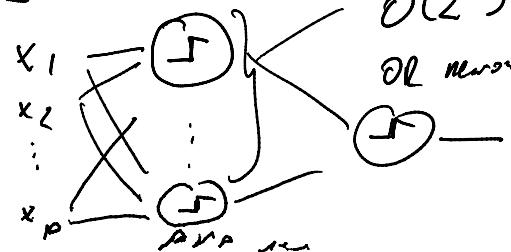
Truth table:

	x_1	x_2	\dots	x_D	y	y^1	y^2	y^D
0	0	0	0	0	0	0	0	0
1	0	0	0	0	1	1	1	1
2	0	0	0	1	1	1	1	1
3	0	0	1	0	1	0	0	0
4	0	0	1	1	0	0	0	0
5	0	1	0	0	0	0	0	0
6	0	1	0	0	1	0	0	0
7	0	1	0	1	1	0	0	0
8	0	1	1	0	0	0	0	0
9	0	1	1	0	1	0	0	0
10	0	1	1	1	1	0	0	0
11	1	0	0	0	0	0	0	0
12	1	0	0	0	1	0	0	0
13	1	0	0	1	1	0	0	0
14	1	0	1	0	0	0	0	0
15	1	0	1	0	1	0	0	0
16	1	0	1	1	1	0	0	0
17	1	1	0	0	0	0	0	0
18	1	1	0	0	1	0	0	0
19	1	1	0	1	1	0	0	0
20	1	1	1	0	0	0	0	0
21	1	1	1	0	1	0	0	0
22	1	1	1	1	1	0	0	0

$$y = \bigvee_{i:y^i=1} (x_1 = x_1^i \wedge x_2 = x_2^i \wedge \dots)$$

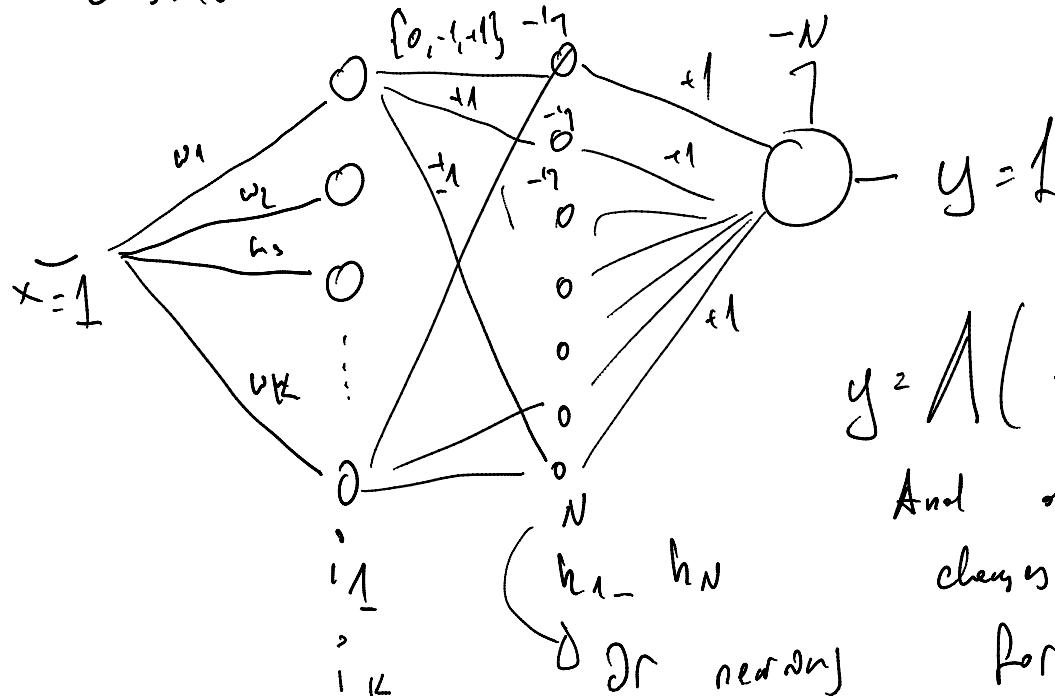
make
a
net

the value of x_1 in
row i : the catch:
 $O(2^D)$ neurons



Hidden layer: a really big deal: Optimization Hardness

Consider this neural net:



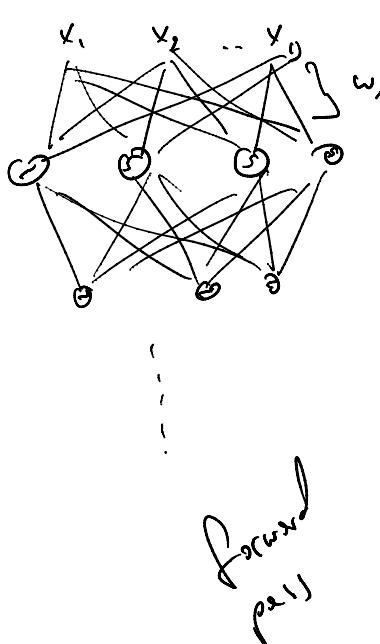
$y^2 \neq 1 (\checkmark)$

And of several changes - CNF
Formulas

$O_p + \text{Loss} = 0 \equiv \text{Solved SAT}$

ERROR BACKPROPAGATION

Backprop intuitions

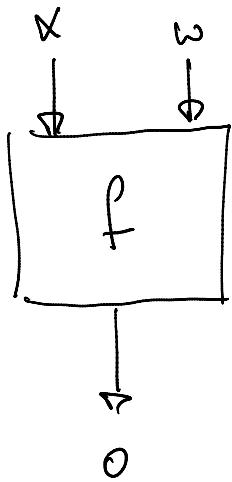


$$\begin{aligned} X &= \\ A^1 &= w^1 x \\ o^1 &= g(A^1) \\ A^2 &= w^2 o^1 \\ o^2 &= g(A^2) \\ \text{Loss} &= L(o^L) \end{aligned}$$

Diagram illustrating the backpropagation process:

- The forward pass consists of the operations: $A^1 = w^1 x$, $o^1 = g(A^1)$, $A^2 = w^2 o^1$, and $o^2 = g(A^2)$.
- The backward pass (backpropagation) involves calculating gradients of the loss function with respect to the inputs. This is shown as:
 $\frac{\delta L}{\delta A^1} = \frac{\delta o^1}{\delta A^1} \frac{\delta L}{\delta o^1}$
 $\frac{\delta L}{\delta A^2} = \frac{\delta o^2}{\delta A^2} \frac{\delta L}{\delta o^2}$
- The Jacobian matrix is indicated by a bracketed term: $\left[\frac{\delta o^2}{\delta A^2} \right]$.
- The backward pass is labeled "backward pass".

Modular implementation



forward: $o = f(x, w)$

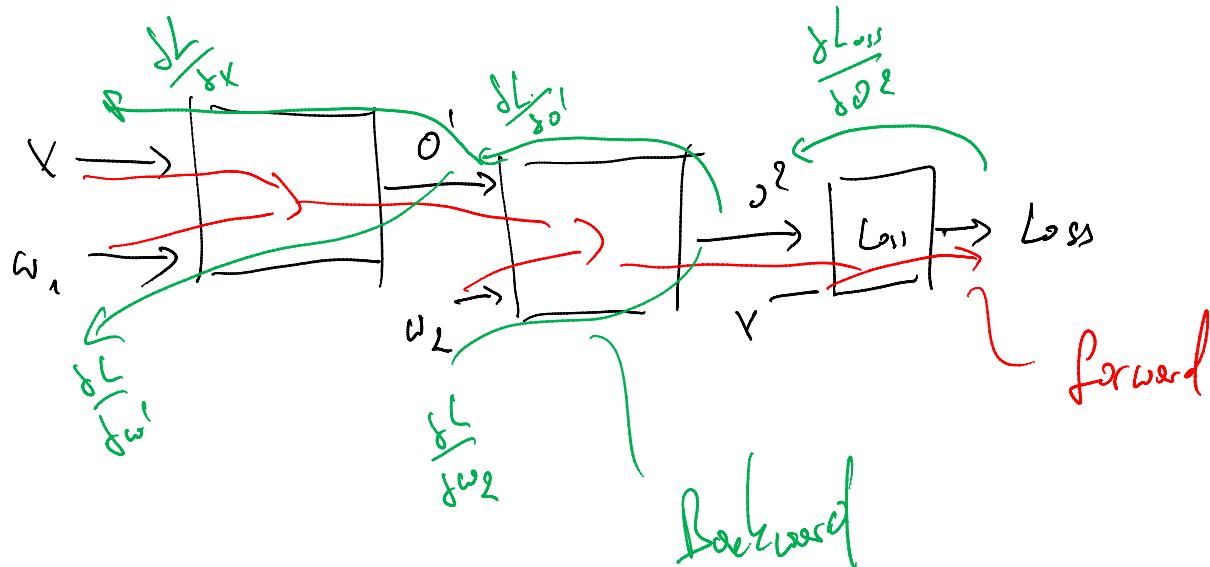
backward :

$$\frac{\delta L}{\delta x} = \frac{\delta o}{\delta x} \frac{\delta L}{\delta o}$$
$$\frac{\delta L}{\delta w} = \frac{\delta o}{\delta w} \frac{\delta L}{\delta o}$$

c.g. a lin. layer

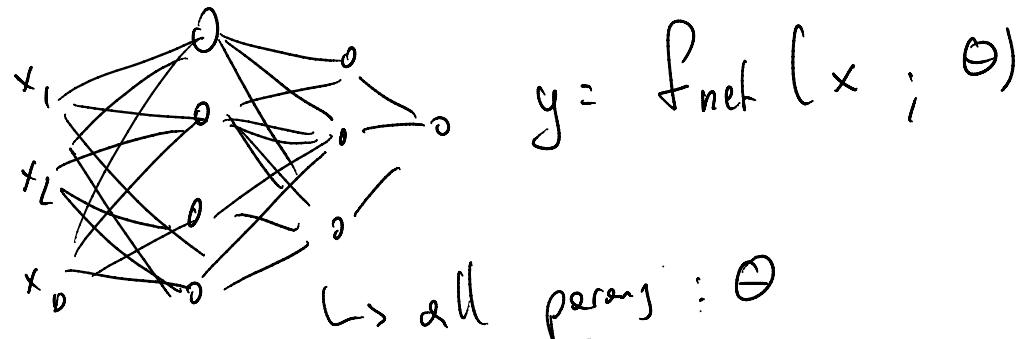
$$o = w \cdot x \quad \frac{\delta L}{\delta x} = W^T \frac{\delta L}{\delta o} \quad \frac{\delta L}{\delta w} = \frac{\delta L}{\delta o} \cdot x$$

Full Net Implementation



TRAINING INTUITIONS

Function fitting



We want: $\min_{\theta} \mathbb{E}_{x,y} [\text{Loss}(f_{\text{net}}(x; \theta), y)]$

Approx # on train set $\{(x^i, y^i) : i=1..N\}$

$$\min_{\theta} L(\theta) = \frac{1}{N} \sum_i \text{Loss}(f_{\text{net}}(x^i; \theta), y^i)$$

$$\text{e.g. } L(\theta) = \frac{1}{N} \sum_i \frac{1}{2} (f_{\text{net}}(x^i; \theta) - y^i)^2$$

Gradient Descent

$$\min_{\theta} L(\theta) = \sum_i^m (f_{\text{net}}(x_i, \theta) - y_i)^2$$

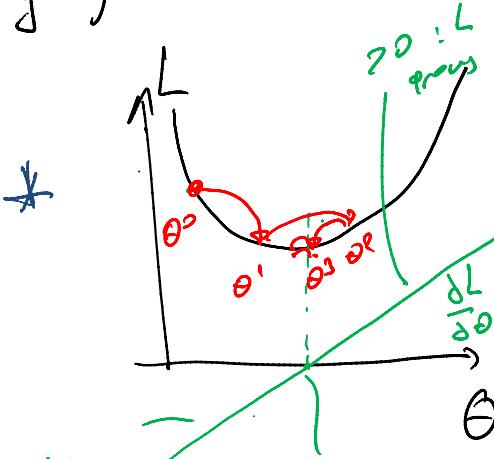
$$\theta \leftarrow \theta^0$$

while not converged

$$\theta \leftarrow \theta - \alpha \frac{\partial L(\theta)}{\partial \theta}$$

$$\text{at opt } \frac{\partial L}{\partial \theta} = 0$$

↳ Grad Desc. stops
by itself



ΔL gets smaller
 θ^0
stationary point

Fine print:
t: Each GD step
one and pass over data

Perceptron

$$y = \Gamma(\theta^T x)$$

$$\theta \leftarrow \theta^0$$

while not converged

$$x^i, y^i \leftarrow \text{sample}(x, y)$$

$$o^i \leftarrow \Gamma(\theta^T x^i)$$

$$\text{if } o^i \neq y^i$$

$$\theta \leftarrow \theta - \alpha \underbrace{(o^i - y^i)}_{\substack{\text{step} \\ \text{size}}} \underbrace{x^i}_{\substack{\text{error in} \\ \text{i-th step}}} \quad \begin{matrix} \nearrow \\ \text{input} \end{matrix}$$

$$\text{case 1: } y^i = 1$$

$$\begin{aligned} o^i \neq y^i &\Rightarrow o^i = 0 \Rightarrow \\ \theta^T x^i &\leq 0 \quad \text{too small} \end{aligned}$$

After update:

$$\begin{aligned} (\theta - \alpha \underbrace{(o^i - y^i)}_{\substack{\text{step} \\ \text{size}}} x^i)^T x^i &= \\ &= \theta^T x^i + \underbrace{\alpha x^i \cdot x^i}_{2,0} \geq \theta^T x^i \end{aligned}$$

Stochastic Gradient Descent

We want: $\min_{\theta} \mathbb{E}_{x,y} [\text{Loss}(f_{\theta}(x, \theta), y)]$

Grad Descent: approx \mathbb{E} by average loss on train data

Stochastic Grad Descent:

approx \mathbb{E} by Loss on a single data sample

• on a "mini batch"

Result: noise ↑ speed ↑ of samples

SGD for linear regression

$$f_{NET}(x, \theta) = \theta^T x$$

$$\text{Loss}(f_{NET}(x, \theta), y) = \frac{1}{2}(\theta^T x - y)^2$$

$$\frac{\delta L(y)}{\delta \theta} = (\theta^T x - y) x$$

Alg:

$$\theta \leftarrow \theta^0$$

while not converged:

x^i, y^i & sample single data sample
very similar to Perceptron

$$\theta^i = \theta^T x$$

$$\theta \leftarrow \theta - \alpha \frac{\delta L(y)}{\delta \theta} = \overbrace{\theta - \alpha}^{\text{step size}} \underbrace{(\theta^i - y^i)}_{\text{error}} \overbrace{x^i}^{\text{input}}$$

SGD: importance of α

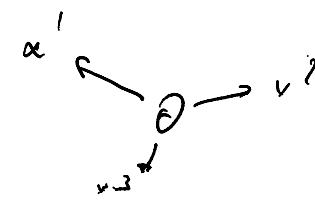
At optimum $\Theta = \arg \min_{\Theta} \mathbb{E}[\text{Loss}(f_{\text{net}}(x, \Theta), y)]$

end $\mathbb{E}\left[\frac{\partial \text{Loss}}{\partial \Theta}\right] = \frac{\partial}{\partial \Theta} \mathbb{E}\left[\frac{\text{Loss}}{\partial \Theta}\right] = 0$

↑ average grad on all
• samples

- samples "pull" Θ in many dir,
average is 0

- SGD looks at a single
sample \Rightarrow it never converges?



SGD: schedules for α

SGD needs a "learning rate schedule"

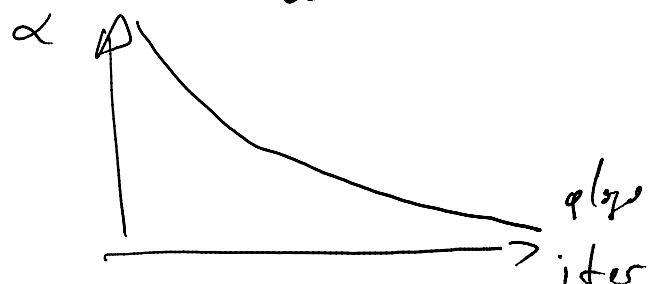
i.e. while True:

$x^i, y^i \leftarrow$ sample

$$\theta \leftarrow \theta - \alpha \frac{\delta L_{\text{tot}}}{\delta \theta}$$

$\alpha \leftarrow$ or smaller val., e.g. 0.999 α

formal reg's (strongly
convex)



$$\begin{cases} \alpha_i = \infty \\ \alpha_i < \infty \quad \text{e.g. } \alpha_i = \frac{1}{i} \end{cases}$$

SGD: practical schedules for α

1 neural net's loss surface is not convex
- many local minima (current hyp: all similar)

- many saddle points / flat regions

- sharp vs flat minima

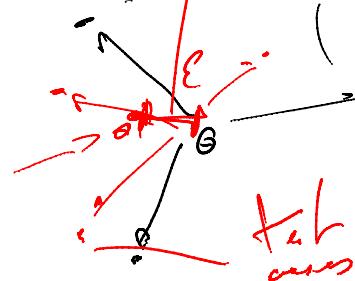
↳ went to a flat minimum, i.e.

$$f(x, \theta) \approx f(x, \theta + \epsilon)$$

train-test
diff.

initialization

opt θ
on test



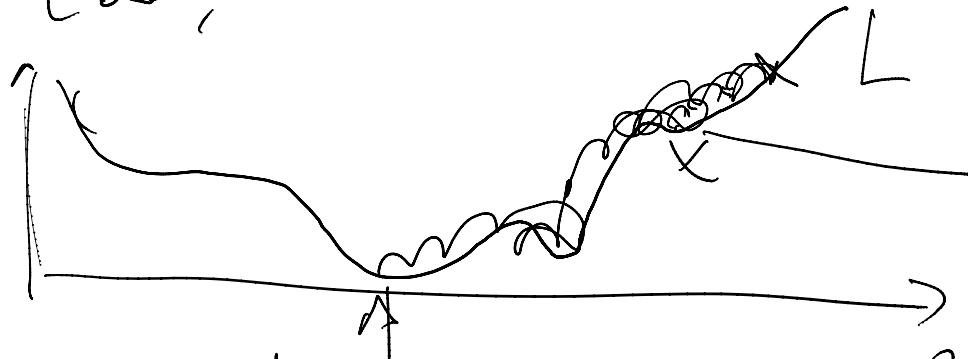
\Rightarrow SGD noise:

- regularizes
- escapes bad/narrow local opt

SGD intuitions:

Intuition:

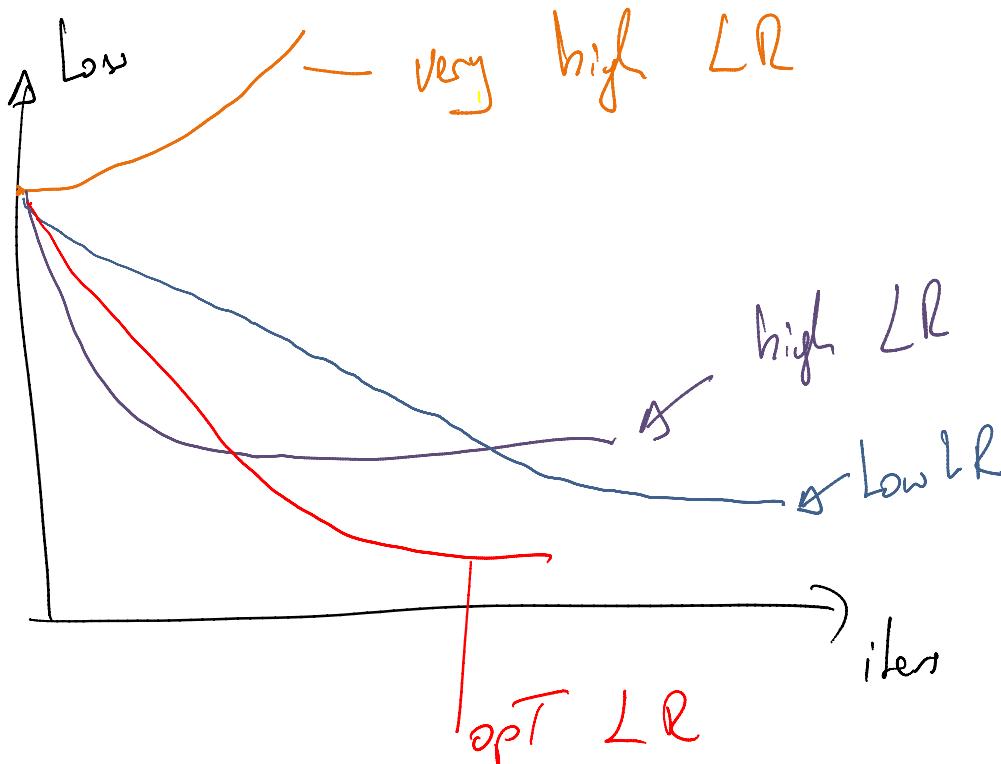
(but, 2D space is nothing like N-D space)



due to noise, SGD
might end up here

GD / 2nd order
method
might
converge
here

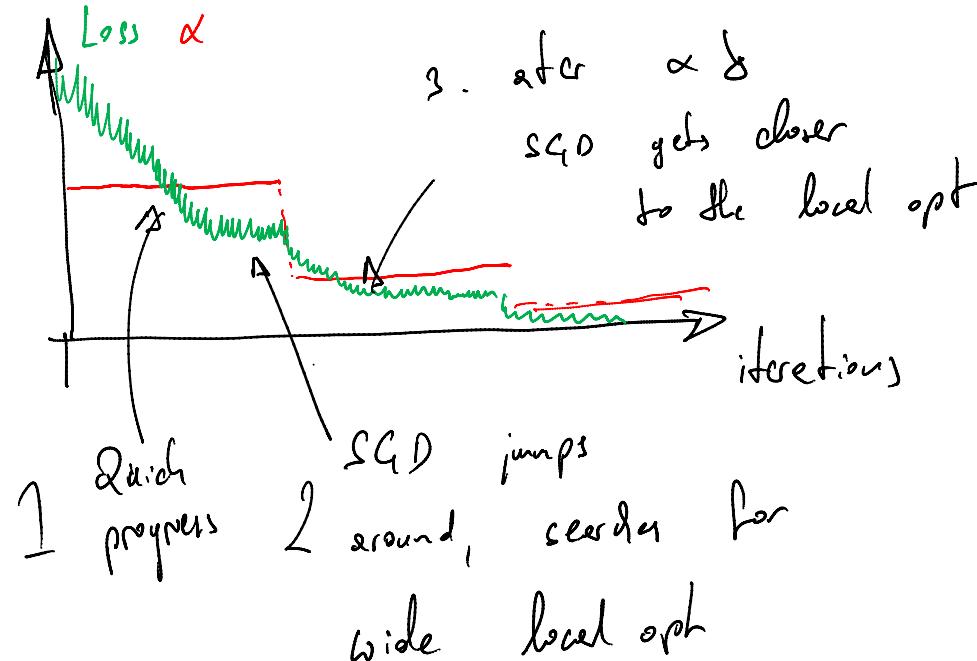
SGD: impact of learning rate



SGD: practical learning rates

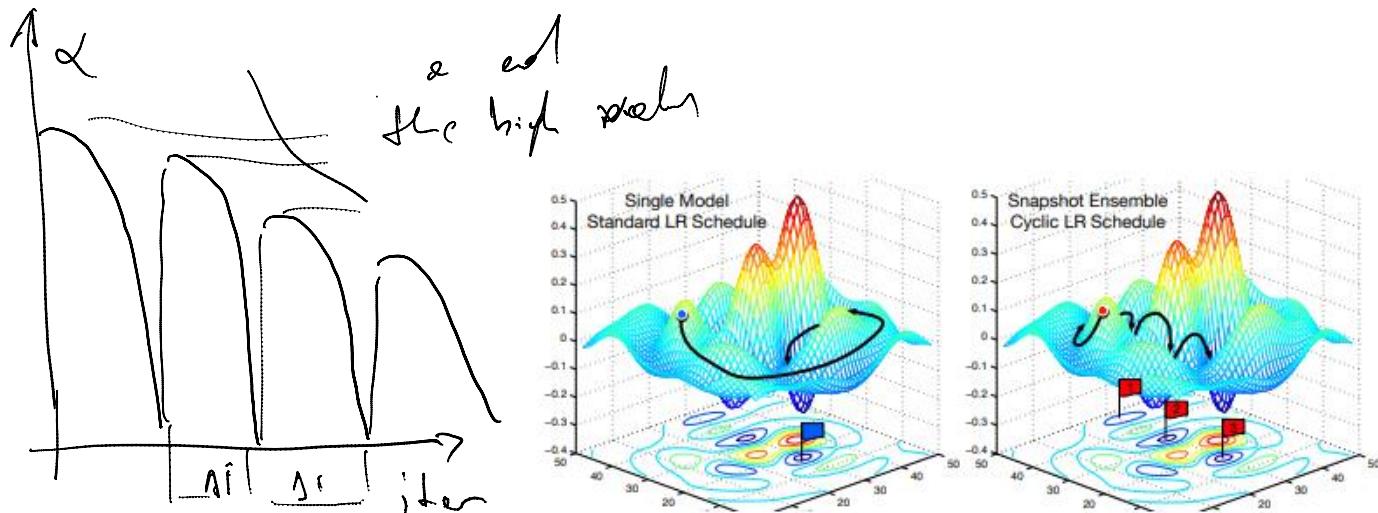


- Stair case



SGD: fancy learning rate schedules

- <https://arxiv.org/pdf/1506.01186.pdf>
- <https://arxiv.org/pdf/1704.00109.pdf>



Speeding up SGD

- Momentum:

<https://distill.pub/2017/momentum/>

- Spread each update over many steps

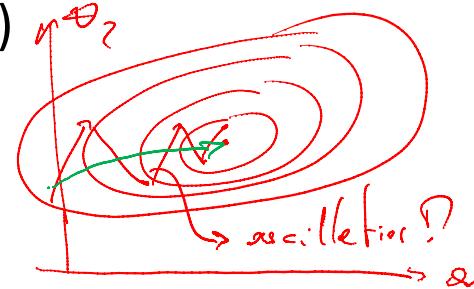
$$\begin{aligned} V_t &\leftarrow \mu_t V_{t-1} - \alpha_t \frac{\partial J}{\partial \Theta} \\ \Theta_t &\leftarrow \Theta_{t-1} + V_t \end{aligned}$$

encrege

- Theoretical justification in the batch setting
(Nesterov accelerated gradient)

- Popular

need large steps →
small step →
but SGD has only 1!



Only use sign of the gradient

- Resilient Propagation (RProp, IRProp):
 - Use only the sign of the gradient
 - Rules for learning rate scaling
 - Only works in BGD
- RMSProp (Rprop for SGD)

$$r_t \leftarrow (1 - \gamma) \left(\frac{\partial J}{\partial \Theta} \right)^2 + \gamma r_{t-1}$$
$$\Theta_t \leftarrow \Theta_t - \frac{\alpha_t}{\sqrt{r_t + \epsilon}} \frac{\partial J}{\partial \Theta}$$

- Keep a running average of the magnitude to get the sign

Adam – the recommended default

- RMSprop with momentum
- <http://arxiv.org/abs/1412.6980>

```
Require:  $\alpha$ : Stepsize
Require:  $\beta_1, \beta_2 \in [0, 1]$ : Exponential decay rates for the moment estimates
Require:  $f(\theta)$ : Stochastic objective function with parameters  $\theta$ 
Require:  $\theta_0$ : Initial parameter vector
 $m_0 \leftarrow 0$  (Initialize 1st moment vector)
 $v_0 \leftarrow 0$  (Initialize 2nd moment vector)
 $t \leftarrow 0$  (Initialize timestep)
while  $\theta_t$  not converged do
     $t \leftarrow t + 1$ 
     $g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$  (Get gradients w.r.t. stochastic objective at timestep  $t$ )
     $m_t \leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$  (Update biased first moment estimate)
     $v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$  (Update biased second raw moment estimate)
     $\hat{m}_t \leftarrow m_t / (1 - \beta_1^t)$  (Compute bias-corrected first moment estimate)
     $\hat{v}_t \leftarrow v_t / (1 - \beta_2^t)$  (Compute bias-corrected second raw moment estimate)
     $\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon)$  (Update parameters)
end while
return  $\theta_t$  (Resulting parameters)
```

Model averaging

- Use many models:
 - Train many models, and average (ensemble)
 - Train 1 net, but once in a while use high LR to jump to another local optimum (“snapshot ensemble”)
- Stabilize the weights by Polyak averaging
 - Keep a running average of parameters:

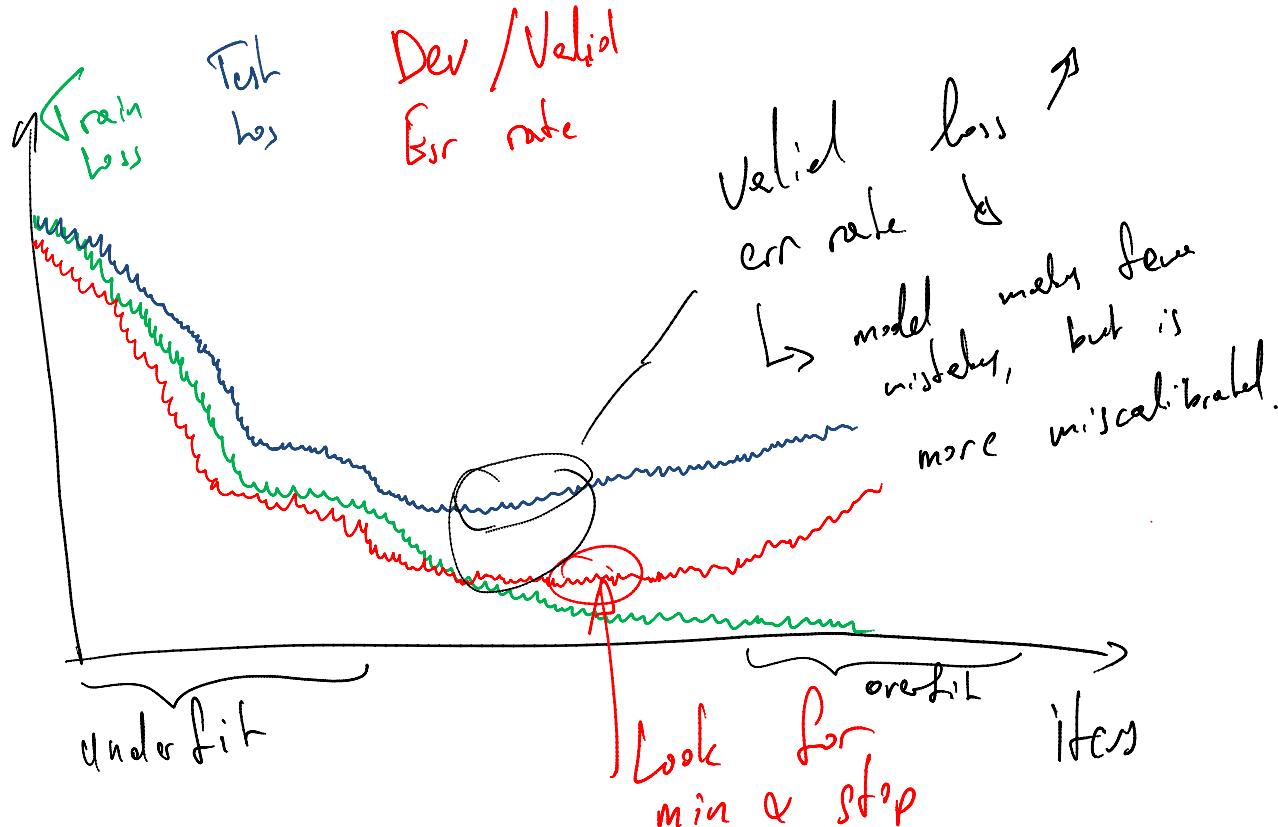
$$\begin{aligned}\Theta &\leftarrow \text{normal train update} \\ \Theta^{ave} &\leftarrow 0.995\Theta^{ave} + 0.005\Theta\end{aligned}$$

Use Θ^{ave} for testing & profit!

What to monitor during training

- Loss
- Magnitude (mean and st. dev) of gradient for (best for each layer)
- Activation values (mean, st. dev.) for each layer (monitors whether units are “stuck”)
- Good practices:
 - Save the model during training.
 - Make the analysis from the saved models -> less clutter in the optimization code.

Early stopping



Have some patience!

With SGD we expect some steps in the wrong direction !?

- Patience algorithm:

- Set patience to some initial value

- While $\text{iters} < \text{patience}$

- if significant improvement of some criterion

- $\text{patience} = \max(\text{patience}, \text{iter}^* 2)$

Can also use patience for learning rate decay!

E.g.:

https://pytorch.org/docs/stable/_modules/torch/optim/lr_scheduler.html#ReduceLROnPlateau

Neural Networks: scale matters

Consider a lin. layer: $O_1 = \omega_1 x$

Consider another one: $O_2 = \frac{1}{10} \omega_2 x$ with $\omega_2 = 10 \omega_1$

Q: How do they differ?

Forward \rightarrow same operations,

$$O_2 = \frac{1}{10} (10 \omega_1) x = \omega_1 x$$

Backward:

$$\frac{\delta \text{Loss}}{\delta \omega_1} = \frac{\delta L}{\delta O_1} x$$

$$\frac{\delta \text{Loss}}{\delta \omega_2} = \frac{\delta L}{\delta O_2} x \cdot \frac{1}{10} ?$$

ω_2 : ωx larger, but ωx smaller grad? ?

$100 \times ??$

→
grad
step
magnitude
decay

Choose initial weights wisely

- Goal: the outputs of the neuron have 0 mean, st. dev 1.
- Assume inputs are normalized.

$$O = \sum_{i=1}^n w_i x_i$$

$\underbrace{\qquad\qquad\qquad}_{\substack{\text{mean } 0, \text{ std dev } 1}} \quad \underbrace{\qquad\qquad\qquad}_{\substack{\text{mean } 0, \text{ std dev } = \delta}}$

$$\text{std dev}(O) = \sqrt{n \cdot \text{std dev}^2(w_i x_i)} = \sqrt{n} \delta \Rightarrow \delta \approx \frac{1}{\sqrt{n}}$$

"fan-in" of a neuron

n independent random vars, each with variance δ^2

Then compensate for nonlinearity, e.g. ReLU: half neurons = 0 \Rightarrow sum goes over $\frac{n}{2}$ neurons, ReLU

Choose initial weights wisely

- Goal: the outputs of the neuron has 0 mean, st. dev 1.
- Normalize inputs!
- Typical rules of thumb (for tanh() activ. function):

- $b_{init} = \bar{0}$

- $W_{k_{init}} \sim Uniform\left[-\frac{\sqrt{3}}{\sqrt{n_k}}, \frac{\sqrt{3}}{\sqrt{n_k}}\right]$

why?

$$\text{std (tanh)}(-\alpha, +\alpha) = \frac{2\alpha}{\sqrt{1-\alpha^2}} \approx \frac{1}{\sqrt{n}}$$

$$\Rightarrow \alpha = \frac{\sqrt{n}}{\sqrt{4 - \sqrt{n}}} \approx \frac{\sqrt{2}}{\sqrt{n}}$$

- Or also account for fan-out:

- $W_{k_{init}} \sim Uniform\left[-\frac{\sqrt{6}}{\sqrt{n_k+n_{k+1}}}, \frac{\sqrt{6}}{\sqrt{n_k+n_{k+1}}}\right]$



- For ReLU use $\sqrt{2}$ x larger ranges (since $\frac{1}{2}$ units will be 0), so called Kaiming He initialization

https://pytorch.org/docs/stable/nn.init.html#torch.nn.init.kaiming_uniform

Batch normalization

<http://arxiv.org/abs/1502.03167>

- Normalize all activations in each minibatch!

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{scale and shift}$$

- BN helps because it normalizes the scale of values that flow through the network.

NB: it was originally said to prevent covariate shift (during training the upper layer's inputs are constantly changing), this is currently contested

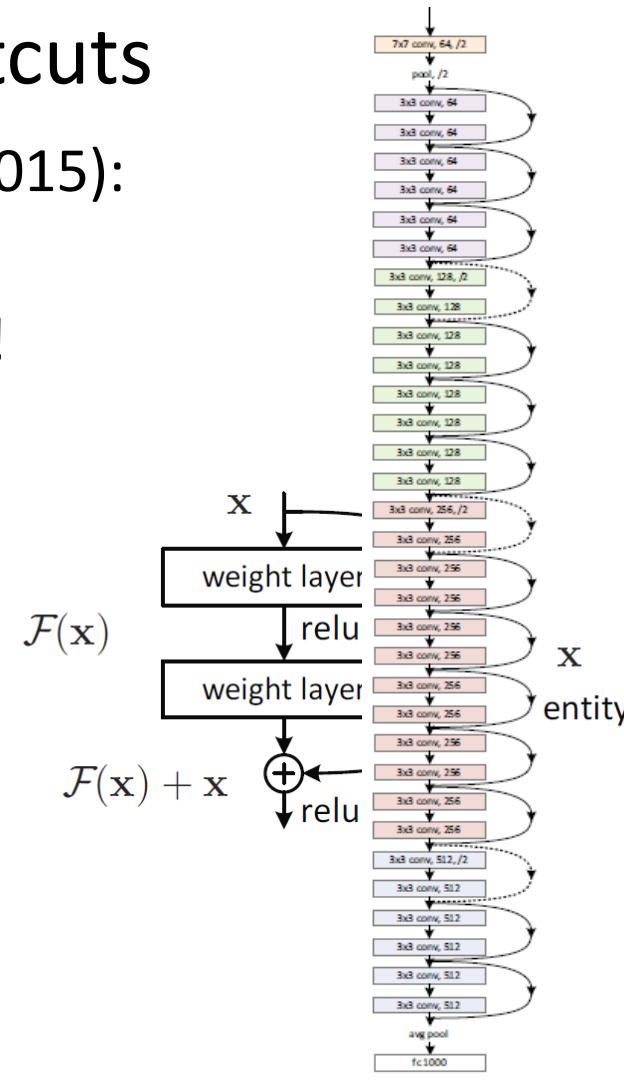
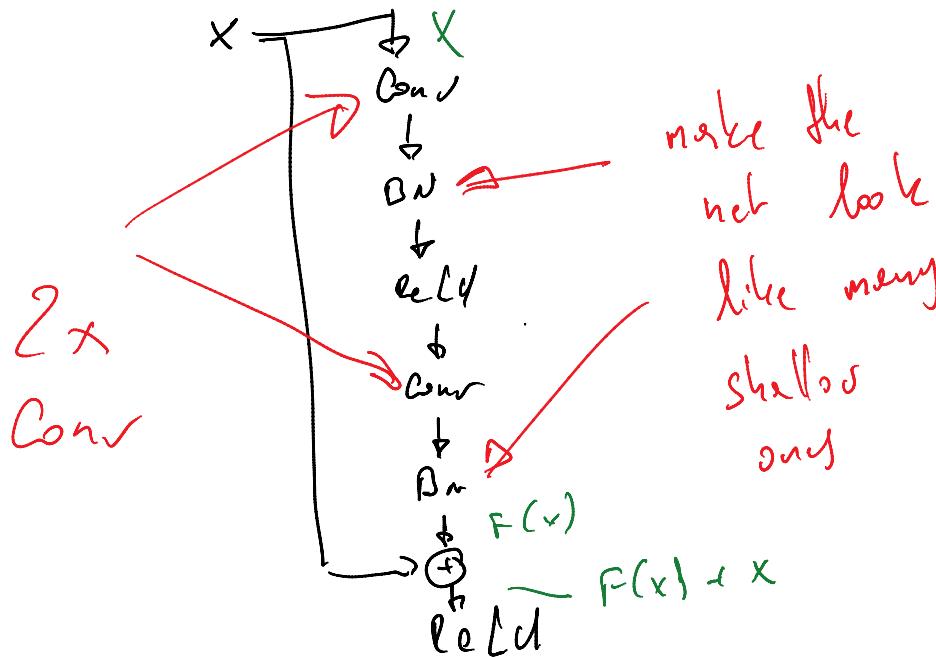
Avoid flat regions of the sigmoids

- For output layer and classification: use SoftMax and cross-entropy loss
- For inner layers:
 - Prefer tanh over log. sigmoid (tanh gives mean 0)
 - The contemporary default: ReLU
http://machinelearning.wustl.edu/mlpapers/papers/icml2010_NairH10
 - Or Maxout <http://arxiv.org/abs/1302.4389>



Use gradient shortcuts

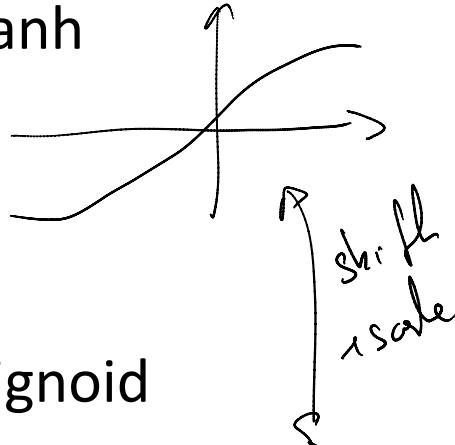
- Residual networks (SOTA Imagenet 2015):
<https://arxiv.org/abs/1512.03385>
 - Networks with more than 150 layers!



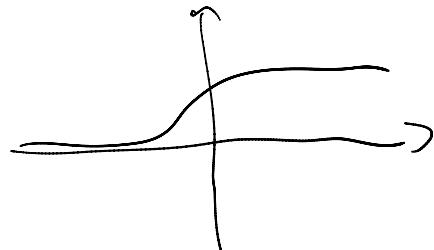
A set of activation functions

Classical:

Tanh



Sigmoid

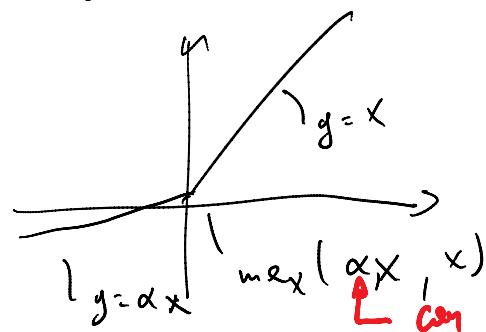


Contemporary

ReLU



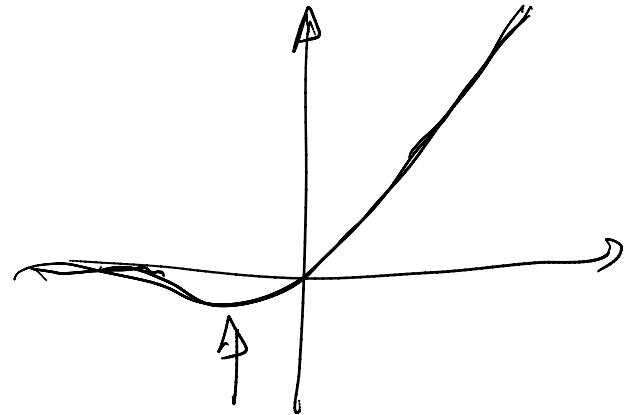
Leaky ReLU



Newcomers

Swish $x\sigma(\beta x)$

(arxiv.org/abs/1710.05941)



can train as well

Babysitting the training process

http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture6.pdf
from slide 61

Where to get more information

- “**Efficient backprop**” by Y. Le Cun
- „Practical Recommendation for Gradient-Based Training of Deep Architectures” Y. Bengio
- Slides by Geoff Hinton in his Coursera lectures:
<https://www.coursera.org/course/neuralnets>
and especially
http://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf
- „Stochastic Gradient Descent Tricks” L. Bottou
- Slides from Stanford:
http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture6.pdf
http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture7.pdf