

3. Basic Feedforward Network

CS 519: Deep Learning, Winter 2017

Fuxin Li

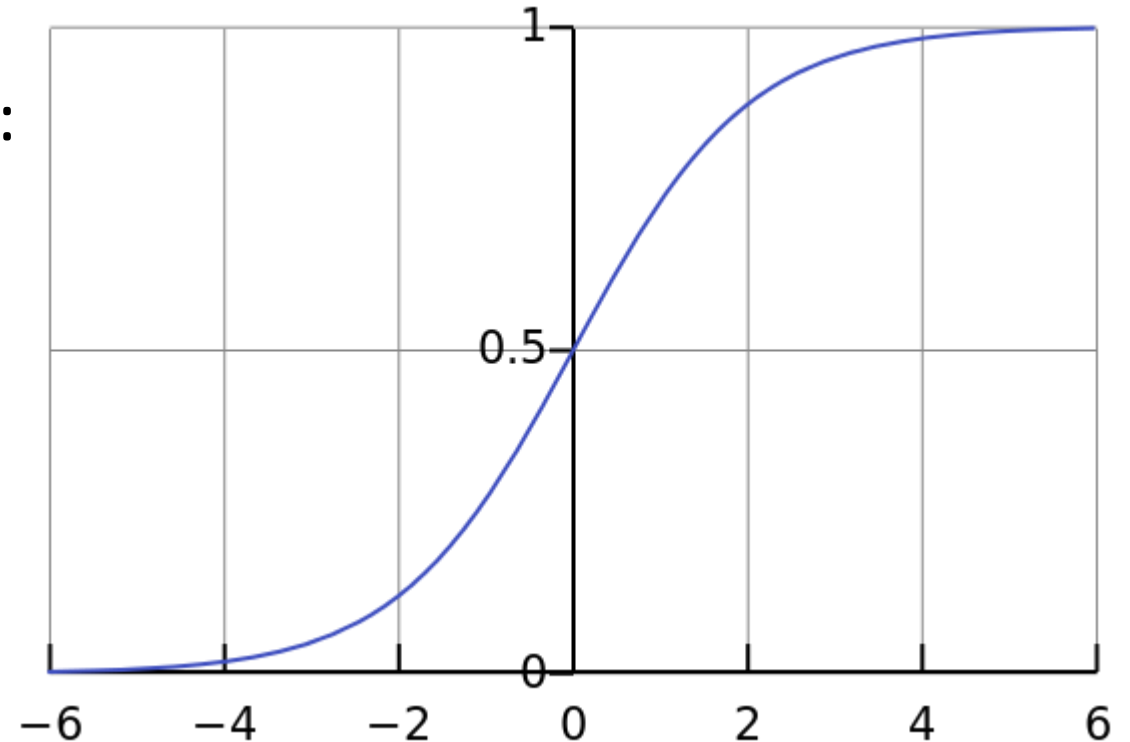
With materials from Zolt Kira, Roger Grosse, Nitish Srivastava, Michael Nielsen

Linear Classifier and the Perceptron Algorithm

- $f(\mathbf{x}) = \sigma(\mathbf{w}^\top \mathbf{x} + b)$
- σ : Sigmoid function: $\sigma(x) = \frac{1}{1+e^{-x}}$
- The connection to logistic regression:
 - Assume binomial distribution with parameter \hat{p}
 - Assume the logit transform is linear:

$$\log \frac{\hat{p}}{1 - \hat{p}} = \mathbf{w}^\top \mathbf{x} + b$$

$$\Rightarrow \hat{p} = \sigma(f(\mathbf{x}))$$



Maximum Log-Likelihood

- MLE of the binomial likelihood:

$$\sum_{i=1}^n y_i^* \log \hat{p}_i + (1 - y_i^*) \log(1 - \hat{p}_i)$$

- where $y_i^* \in \{0,1\}$, if $y_i \in \{-1,1\}$, then $y_i^* = \frac{1+y_i}{2}$

$$\log \hat{p}_i = -\log(1 + e^{-f(x_i)})$$

$$\log(1 - \hat{p}_i) = -\log(1 + e^{f(x_i)})$$

$$y_i^* \log \hat{p}_i + (1 - y_i^*) \log(1 - \hat{p}_i) = -\log(1 + e^{-y_i f(x_i)})$$

Gradient descent optimization

- Optimize \mathbf{w} , b with gradient descent

$$\min_{\mathbf{w}, b} \sum_i \log(1 + e^{-y_i(\mathbf{w}^\top \mathbf{x}_i + b)})$$

$$\begin{aligned} \nabla \mathbf{w} &= \sum_i \frac{-y_i e^{-y_i(\mathbf{w}^\top \mathbf{x}_i + b)}}{1 + e^{-y_i(\mathbf{w}^\top \mathbf{x}_i + b)}} \mathbf{x}_i \\ &= \sum_i -y_i^* (1 - \hat{p}(\mathbf{x}_i)) - (1 - y_i^*) \hat{p}(\mathbf{x}_i) \mathbf{x}_i \end{aligned}$$

$$\nabla b = \sum_i \frac{-y_i e^{-y_i(\mathbf{w}^\top \mathbf{x}_i + b)}}{1 + e^{-y_i(\mathbf{w}^\top \mathbf{x}_i + b)}}$$

XOR problem and linear classifier

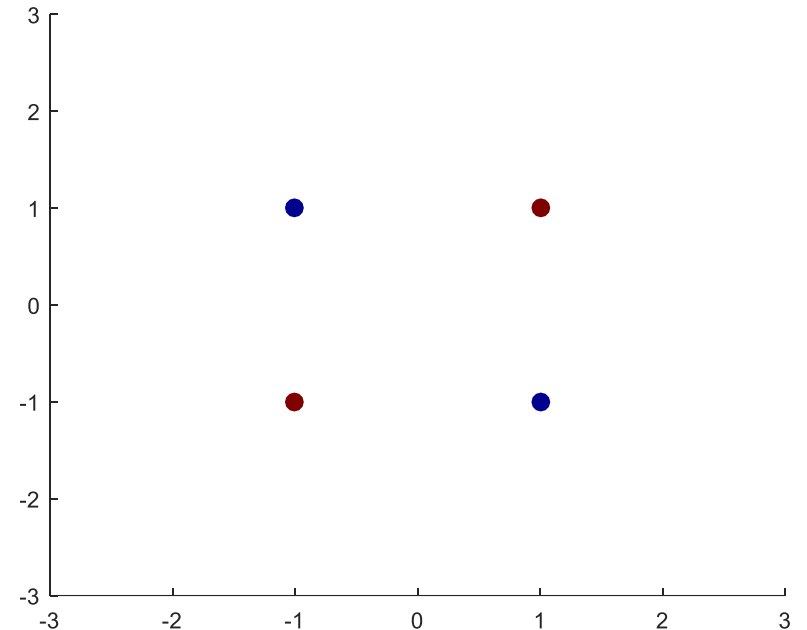
- 4 points: $X = [(-1,-1), (-1,1), (1,-1), (1,1)]$
- $Y = [-1 \ 1 \ 1 \ -1]$
- Try using binomial log-likelihood loss:

$$\min_{\mathbf{w}} \sum_i \log(1 + e^{-y_i(\mathbf{w}^\top \mathbf{x}_i + b)})$$

- Gradient:

$$\nabla \mathbf{w} = \sum_i \frac{-y_i e^{-y_i(\mathbf{w}^\top \mathbf{x}_i + b)}}{1 + e^{-y_i(\mathbf{w}^\top \mathbf{x}_i + b)}} \mathbf{x}_i$$

$$\nabla b = \sum_i \frac{-y_i e^{-y_i(\mathbf{w}^\top \mathbf{x}_i + b)}}{1 + e^{-y_i(\mathbf{w}^\top \mathbf{x}_i + b)}}$$



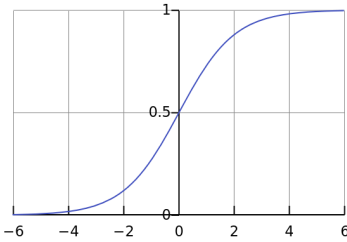
Try $\mathbf{w} = 0, b = 0$,
what do you see?

With 1 hidden layer

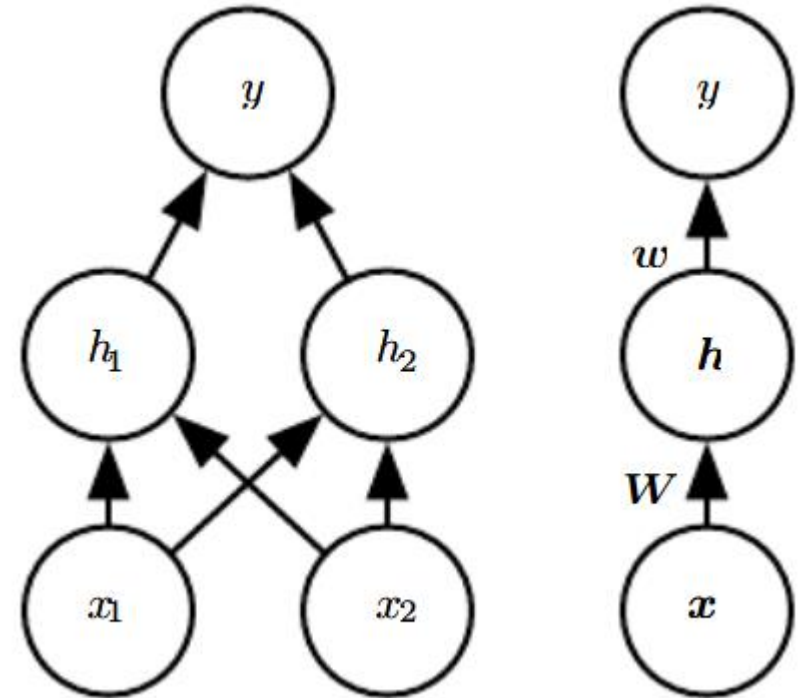
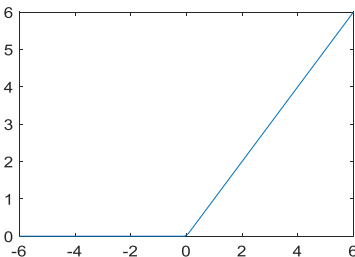
- A hidden layer makes a nonlinear classifier

$$f(x) = \mathbf{w}^T g(\mathbf{W}^T \mathbf{x} + \mathbf{c}) + b$$

- g needs to be nonlinear
- Sigmoid: $\text{Sigm}(x) = 1/(1 + e^{-x})$

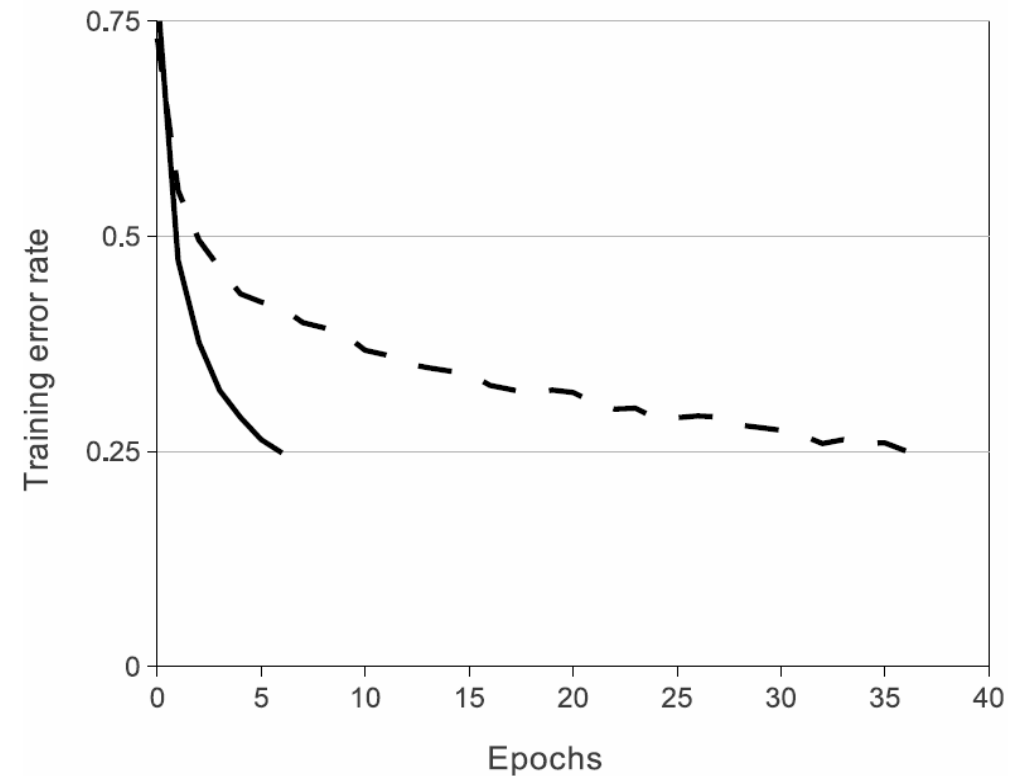


- RELU: $g(x) = \max(0, x)$



ReLU the (only) major improvement in DL

- ReLU vs. Sigmoid convergence rate
- ReLU is 7 times faster to converge than Sigmoid!



Taking gradient

$$\min_{\mathbf{W}, \mathbf{w}} E(f) = \sum_i L(f(\mathbf{x}_i), y_i)$$

$$f(\mathbf{x}) = \mathbf{w}^\top g(\mathbf{W}^\top \mathbf{x} + \mathbf{c}) + b$$

- What is $\frac{\partial E}{\partial \mathbf{W}}$?

- Consider chain rule: $\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$

Note: Vectorized Computations

On the left are the computations performed by a network. Write them in terms of matrix and vector operations. Let $\sigma(\mathbf{v})$ denote the logistic sigmoid function applied elementwise to a vector \mathbf{v} . Let \mathbf{W} be a matrix where the (i, j) entry is the weight from visible unit j to hidden unit i .

$$z_i = \sum_j w_{ij} x_j$$

$$h_i = \sigma(z_i)$$

$$y = \sum_i v_i h_i$$

$$\mathbf{z} = \mathbf{W}\mathbf{x}$$

$$\mathbf{h} = \sigma(\mathbf{z})$$

$$y = \mathbf{v}^T \mathbf{h}$$

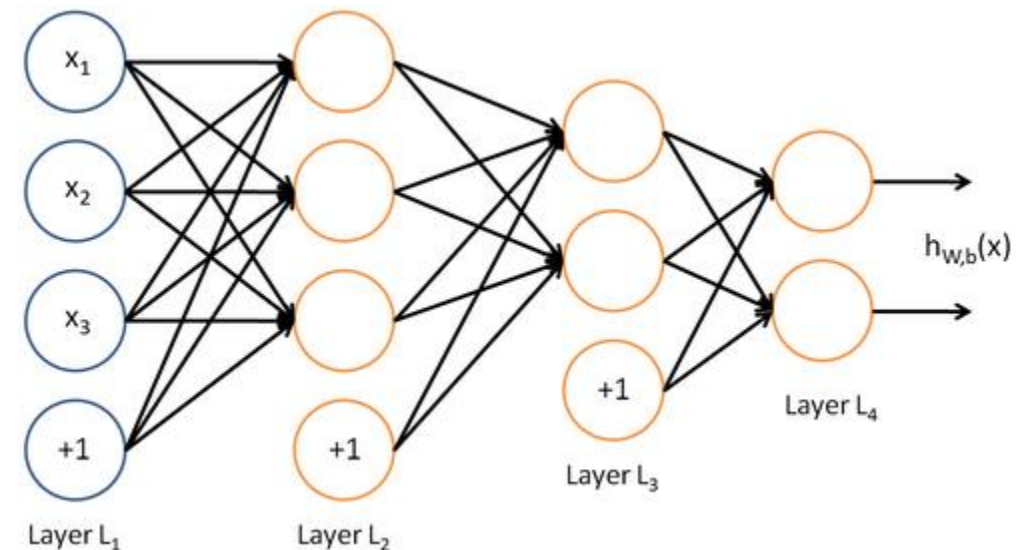
Backpropagation

- Save the gradients and the gradient products that have already been computed to avoid computing multiple times
- In a multiple layer network:
 - (Ignore constant terms)

$$f(x) = \mathbf{w}_n^\top g \left(\mathbf{W}_{n-1}^\top g \left(\mathbf{W}_{n-2}^\top g \left(\dots \left(\mathbf{W}_1^\top g(x) \right) \right) \right) \right)$$

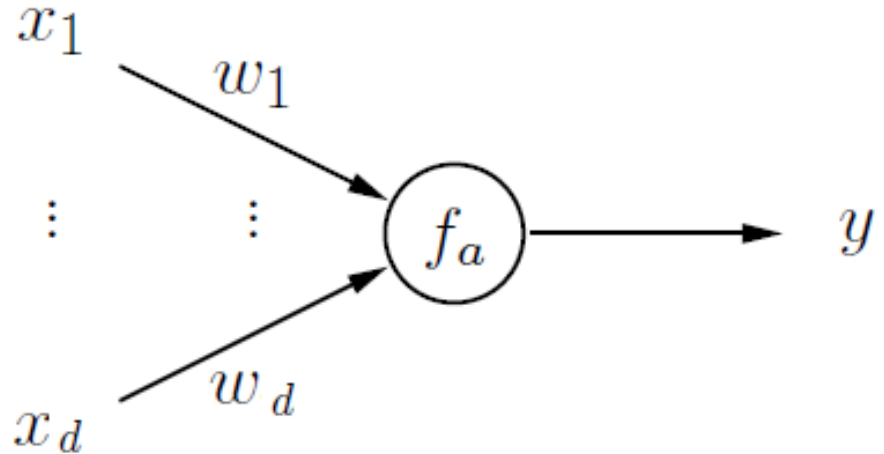
$$\begin{aligned} \frac{\partial E}{\partial \mathbf{W}_k} &= \frac{\partial E}{\partial f_k} g(f_{k-1}(\mathbf{x})) \\ &= \frac{\partial E}{\partial f_{k+1}} \frac{\partial f_{k+1}}{\partial f_k} g(f_{k-1}(\mathbf{x})) \end{aligned}$$

Define: $f_k(\mathbf{x}) = \mathbf{w}_k^\top g(f_{k-1}(\mathbf{x})), f_0(\mathbf{x}) = \mathbf{x}$



Modules

- Each layer can be seen as a module
- Given input, return
 - Output $f_a(\mathbf{x})$
 - Network gradient $\frac{\partial f_a}{\partial \mathbf{x}}$
 - Gradient of module parameters $\frac{\partial f_a}{\partial \mathbf{w}_a}$
- During backprop, propagate/update
 - Backpropagated gradient $\frac{\partial E}{\partial f_a}$

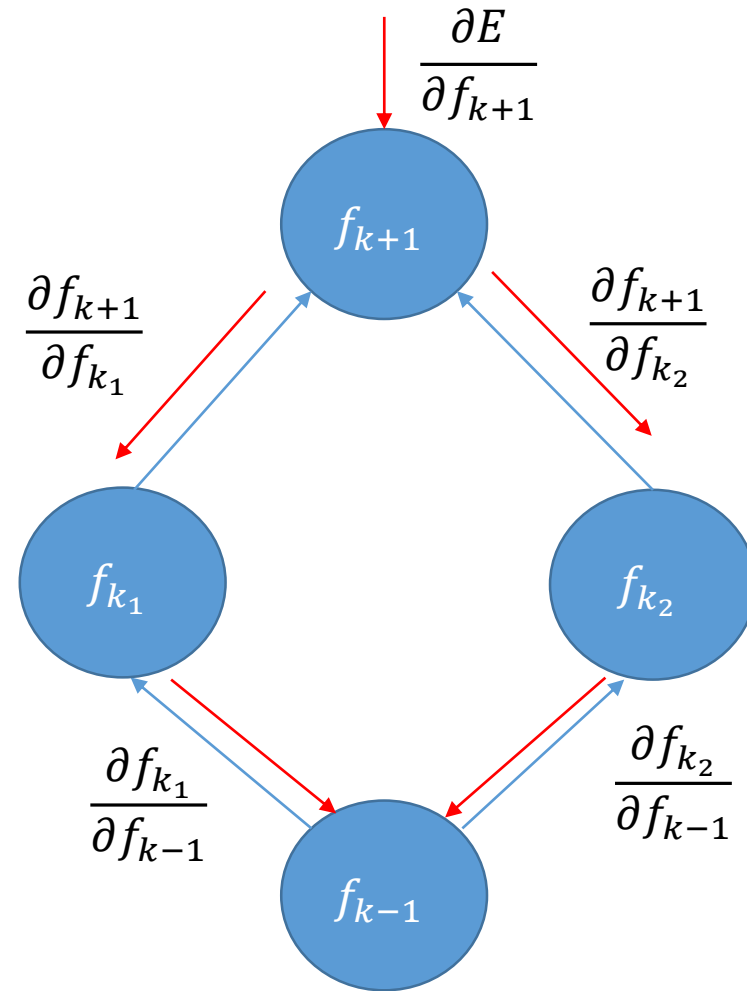


$$\frac{\partial E}{\partial \mathbf{W}_k} = \frac{\partial E}{\partial f_k} g(f_{k-1}(x)) = \underbrace{\frac{\partial E}{\partial f_{k+1}}}_{\text{Backprop signal}} \underbrace{\frac{\partial f_{k+1}}{\partial f_k}}_{\text{Network gradient}} \underbrace{g(f_{k-1}(x))}_{\text{Gradient of parameters}}$$

Note: $\frac{\partial E}{\partial f_k} = \frac{\partial E}{\partial f_{k+1}} \frac{\partial f_{k+1}}{\partial f_k}$

Multiple Inputs and Multiple Outputs

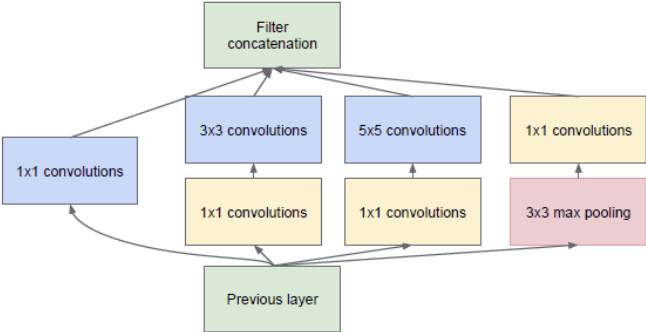
- $$\frac{\partial E}{\partial f_{k-1}} = \frac{\partial E}{\partial f_{k+1}} \frac{\partial f_{k+1}}{\partial f_{k_1}} \frac{\partial f_{k_1}}{\partial f_{k-1}} + \frac{\partial E}{\partial f_{k+1}} \frac{\partial f_{k+1}}{\partial f_{k_2}} \frac{\partial f_{k_2}}{\partial f_{k-1}}$$



Different DAG structures

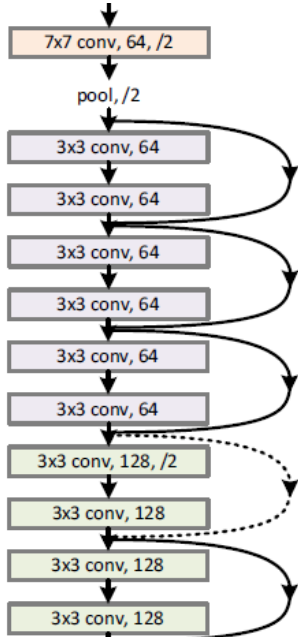
- The backpropagation algorithm would work for any DAGs
- So one can imagine different architectures than the plain layerwise one

Inception

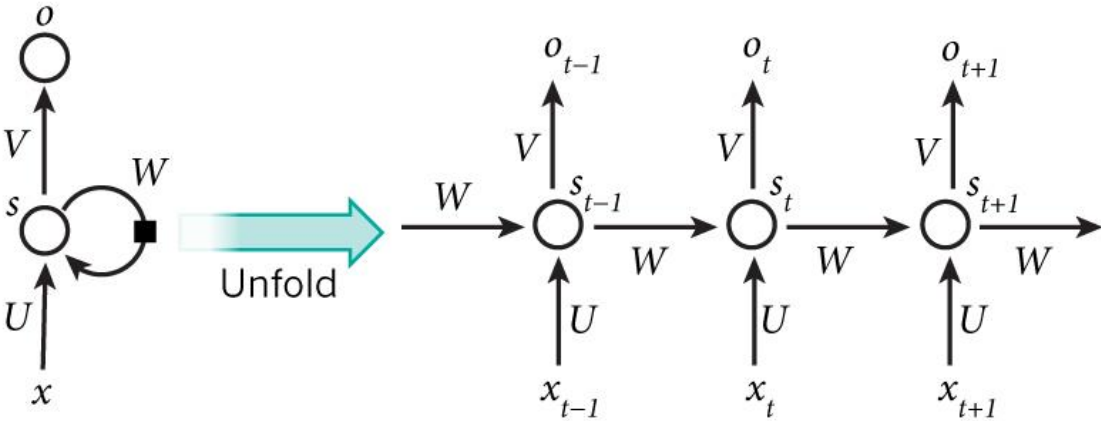


(b) Inception module with dimension reductions

Residual



RNN



Loss functions

- Regression:

- Least squares $L(f) = (f(x) - y)^2$

- L1 loss $L(f) = |f(x) - y|$

- Huber loss $L(f) = \begin{cases} \frac{1}{2} (f(x) - y)^2, & |f(x) - y| < \delta \\ \delta \left(|f(x) - y| - \frac{1}{2} \delta \right), & \text{otherwise} \end{cases}$

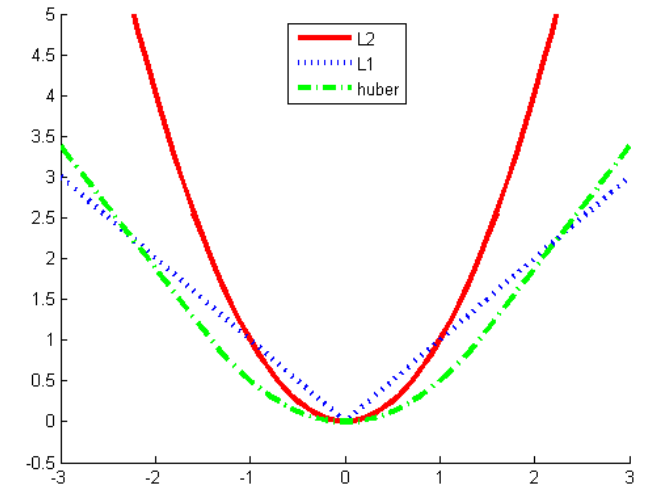
- Binary Classification

- Hinge loss $L(f) = \max(1 - yf(x), 0)$

- Binomial log-likelihood $L(f) = \ln(1 + \exp(-2yf(x)))$

- Cross-entropy $L(f) = -y^* \ln \text{sigm}(f) - (1 - y^*) \ln(1 - \text{sigm}(f))$,

- $y^* = (y + 1)/2$



Multi-class: Softmax layer

- Multi-class logistic loss function

$$P(y = j|\mathbf{x}) = \frac{e^{\mathbf{x}^T \mathbf{w}_j}}{\sum_{k=1}^K e^{\mathbf{x}^T \mathbf{w}_k}}$$

- Log-likelihood:

- Loss function is minus log-likelihood

$$-\log P(y = j|x) = -\mathbf{x}^T \mathbf{w}_j + \log \sum_k e^{\mathbf{x}^T \mathbf{w}_k}$$

Softmax Intuition

- “Soft” argmax function
- e.g. $\mathbf{x}^\top \mathbf{w}_1 = 5, \mathbf{x}^\top \mathbf{w}_2 = 3$
 - $P(y = 1|\mathbf{x}) = 88\%$
 - $P(y = 2|\mathbf{x}) = 12\%$
- e.g. $\mathbf{x}^\top \mathbf{w}_1 = 15, \mathbf{x}^\top \mathbf{w}_2 = 10$
 - $P(y = 1|\mathbf{x}) = 99.33\%$
 - $P(y = 2|\mathbf{x}) = 0.67\%$
- e.g. $\mathbf{x}^\top \mathbf{w}_1 = 15, \mathbf{x}^\top \mathbf{w}_2 = 12, \mathbf{x}^\top \mathbf{w}_3 = 13$
 - $P(y = 1|\mathbf{x}) = 84.4\%$
 - $P(y = 3|\mathbf{x}) = 11.4\%$
 - $P(y = 2|\mathbf{x}) = 4.2\%$

$$P(y = j|\mathbf{x}) = \frac{e^{\mathbf{x}^\top \mathbf{w}_j}}{\sum_{k=1}^K e^{\mathbf{x}^\top \mathbf{w}_k}}$$

Subgradients

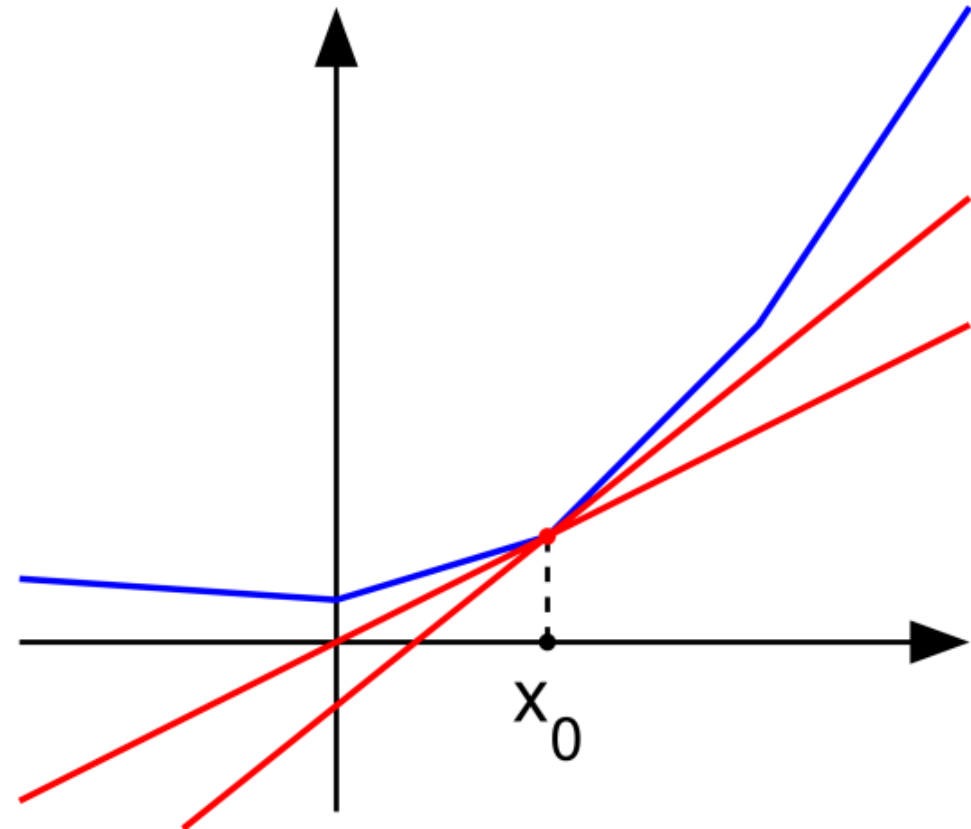
- What if the function is non-differentiable?

- Subgradients:

- For **convex** $f(x)$ at x_0 :
- If for any y

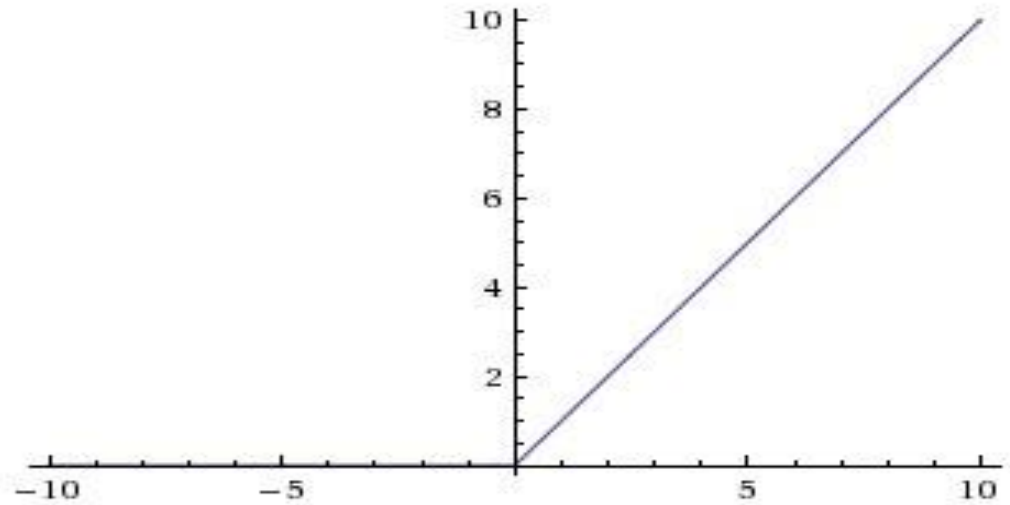
$$f(y) \geq f(x) + g^T(y - x)$$

- g is called a subgradient
- Subdifferential: ∂f : set of all subgradients
- Optimality condition: $0 \in \partial f$



The RELU unit

- $f(x) = \max(x, 0)$
- Convex
- Non-differentiable
- Subgradient: $\frac{df}{dx} = \begin{cases} 1, & x > 0 \\ [0, 1], & x = 0 \\ 0, & x < 0 \end{cases}$



Subgradient descent

- Similar to gradient descent

$$x^{(k+1)} = x^{(k)} - \alpha_k g^{(k)}$$

- Step size rules:

- Constant step size: $\alpha_k = \alpha$.

- Square summable: $\alpha_k \geq 0$, $\sum_{k=1}^{\infty} \alpha_k^2 < \infty$, $\sum_{k=1}^{\infty} \alpha_k = \infty$.

- Usually, a large constant that drops slowly after a long while

- e.g. $\frac{100}{100+k}$

Universal Approximation Theorems

- Many universal approximation theorems proved in the 90s
- Simple statement: every continuous function can be approximated by a 1-hidden layer neural network with arbitrarily high precision

Formal statement [\[edit\]](#)

The theorem^{[2][3][4][5]} in mathematical terms:

Let $\varphi(\cdot)$ be a nonconstant, **bounded**, and **monotonically-increasing continuous** function. Let I_m denote the m -dimensional **unit hypercube** $[0, 1]^m$. The space of continuous functions on I_m is denoted by $C(I_m)$. Then, given any function $f \in C(I_m)$ and $\varepsilon > 0$, there exists an integer N and real constants $v_i, b_i \in \mathbb{R}$, where $i = 1, \dots, N$ such that we may define:

$$F(x) = \sum_{i=1}^N v_i \varphi(w_i^T x + b_i)$$

as an approximate realization of the function f where f is independent of φ ; that is,

$$|F(x) - f(x)| < \varepsilon$$

for all $x \in I_m$. In other words, functions of the form $F(x)$ are **dense** in $C(I_m)$.

It obviously holds replacing I_m with any compact subset of \mathbb{R}^m .

Universal Approximation Theorems

- The approximation does not need many units if the function is kinda nice. Let

$$C_f = \int_{\mathbf{R}^d} \|\omega\| |\tilde{f}(\omega)| d\omega$$

- Then for a 1-hidden layer neural network with n hidden nodes, we have for a finite ball with radius r ,

$$\int_{B_r} (f(x) - f_n(x))^2 d\mu(x) \leq \frac{4r^2 C_f^2}{n}$$