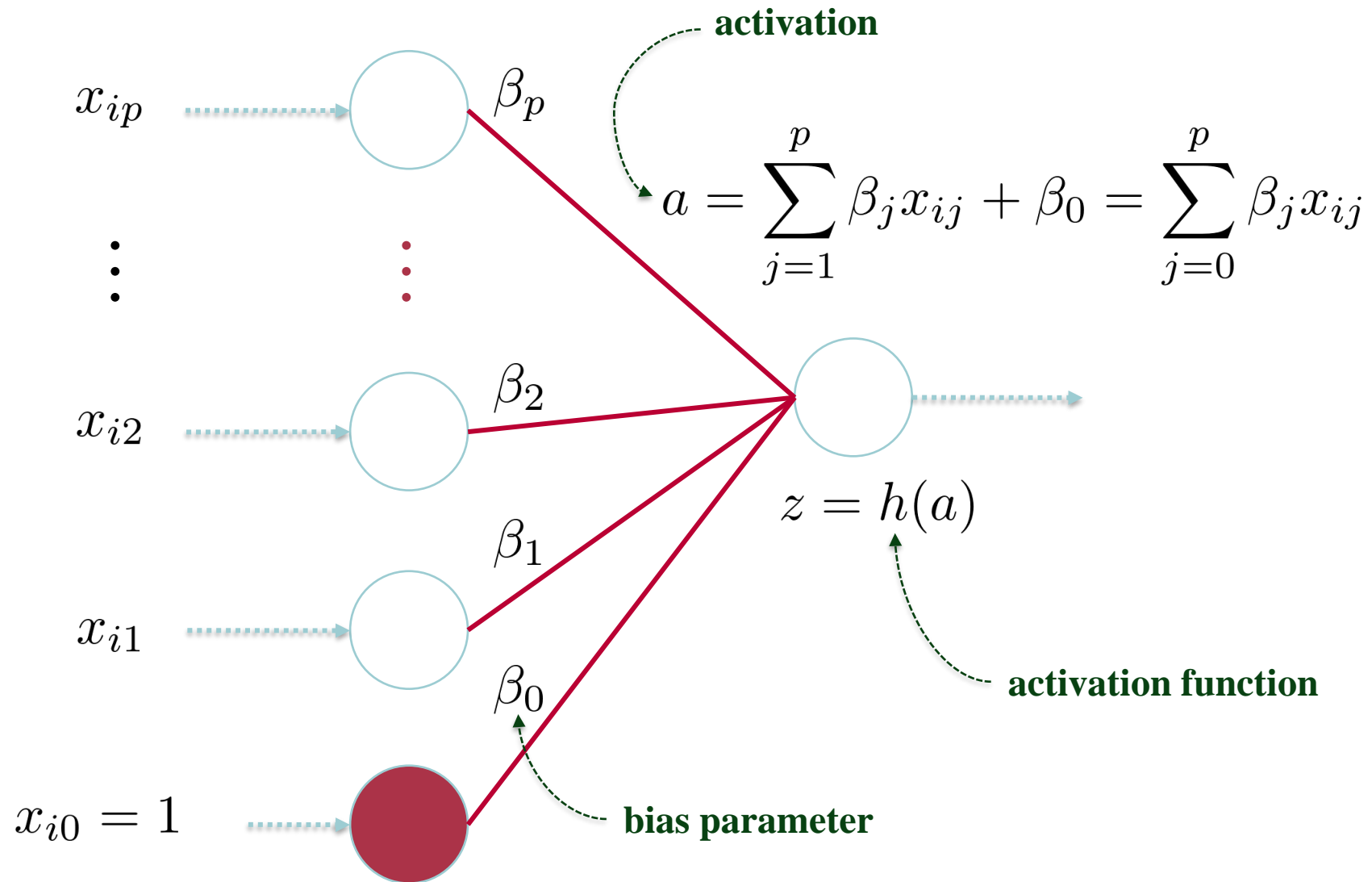


The background of the slide features a complex, abstract graphic. On the left, several teal lines enter from the edge, some of which are stepped or 'bent' like circuit traces. These lines lead towards a central cluster of nodes. The nodes are small circles, mostly teal, but with a horizontal band of orange nodes in the middle. These nodes are interconnected by a web of thin, light-blue lines. From the right side of this central cluster, several teal lines emerge and flow outwards, some appearing as smooth, wavy ribbons. A single, prominent orange line also flows out from the top right of the central node cluster. The overall composition suggests a flow of information or data through a network structure.

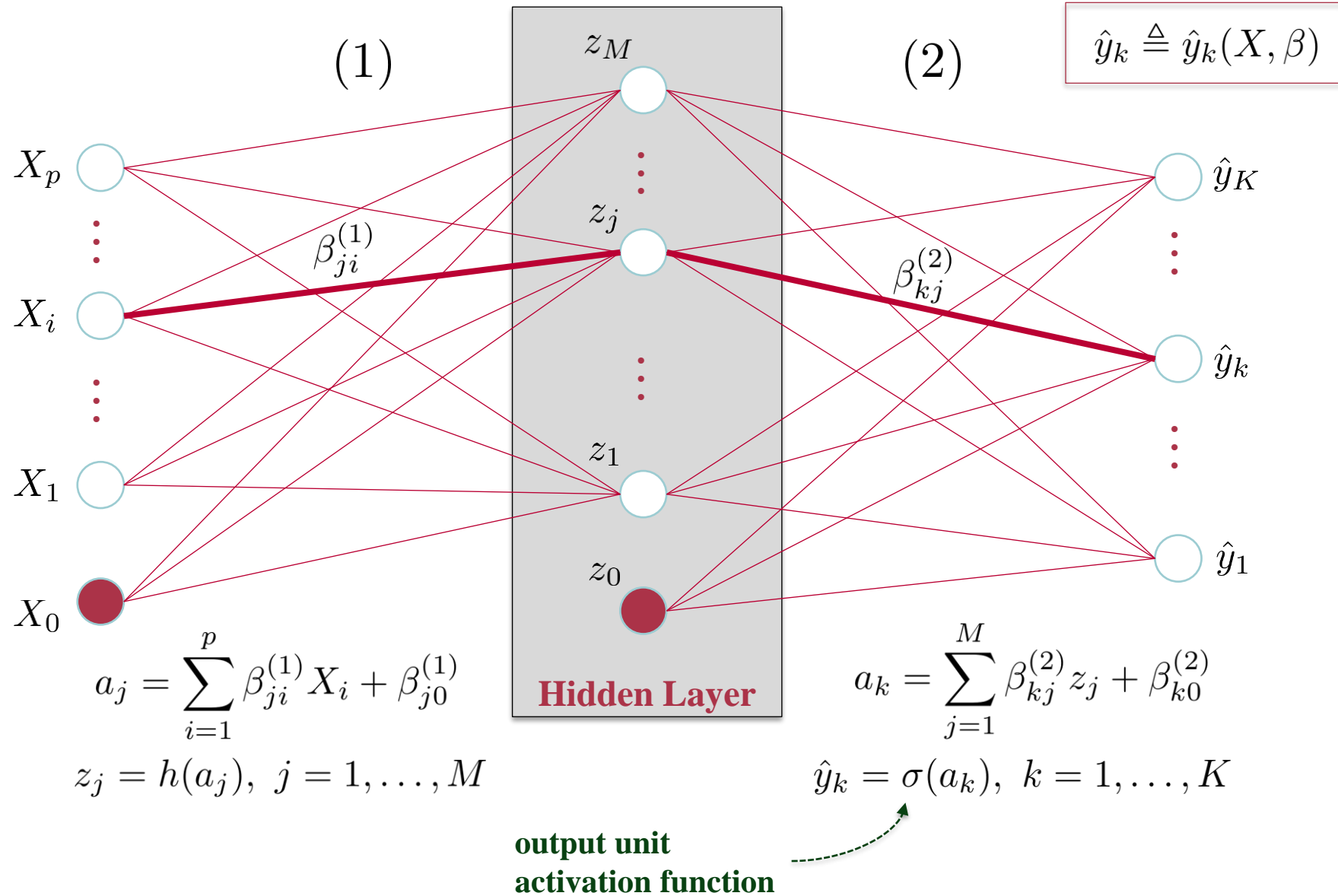
Neural Networks

Özgür Martin

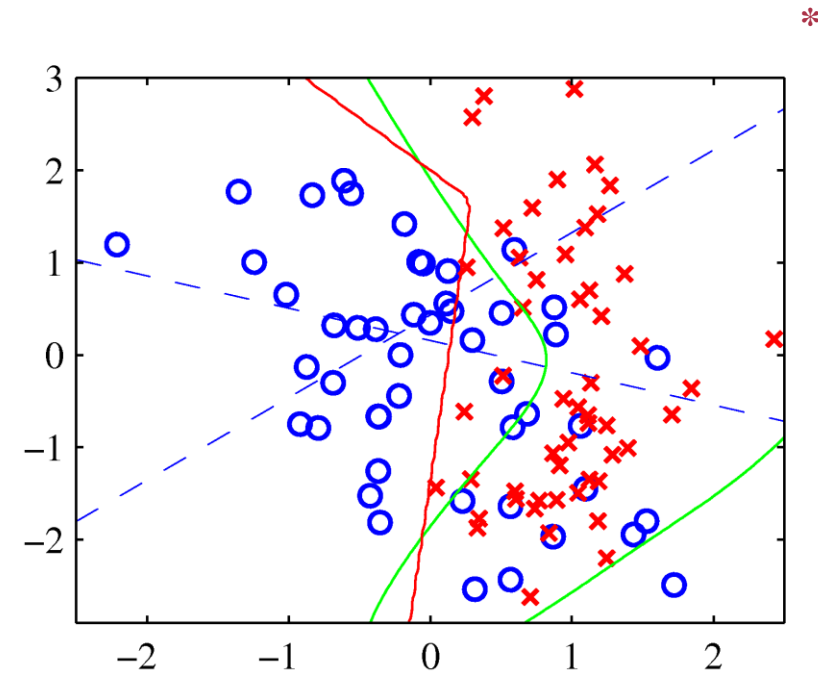
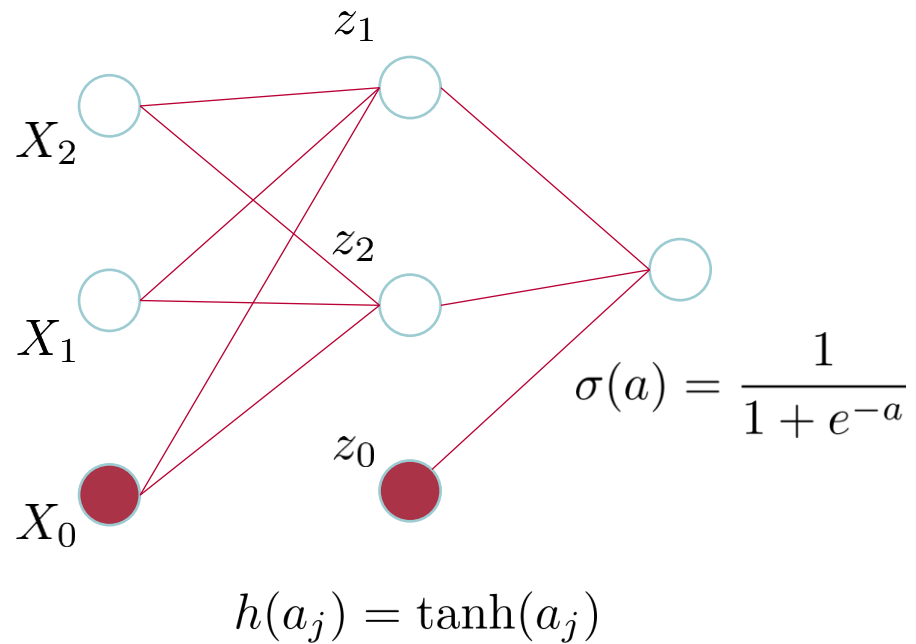
Processing Units



Three-Layer (Feed-Forward) Neural Network



Binary Classification Example



*Pattern Recognition and Machine Learning, C. M. Bishop, Springer, 2006, pg. 232.

Activation Functions

$$h(a_j) = \frac{1}{1 + e^{-a_j}}$$

Sigmoid Function

$$h(a_j) = \tanh(a_j)$$

“tanh” Function

$$h(a_j) = \max\{0, a_j\}$$

Rectified Linear Unit (ReLU)

$$\sigma(a_k) = a_k$$

Identity Function

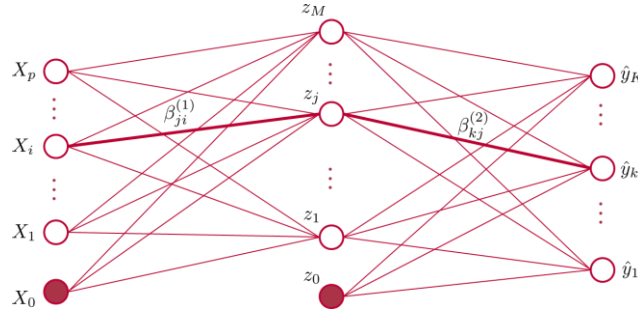
$$\sigma(a_k) = \frac{1}{1 + e^{-a_k}}$$

Sigmoid Function

$$\sigma(a_k) = \frac{e^{a_k}}{\sum_{j=1}^K e^{a_j}}$$

Softmax Function

Network Function



$$a_j = \sum_{i=1}^p \beta_{ji}^{(1)} X_i + \beta_{j0}^{(1)} \quad a_k = \sum_{j=1}^M \beta_{kj}^{(2)} z_j + \beta_{k0}^{(2)}$$

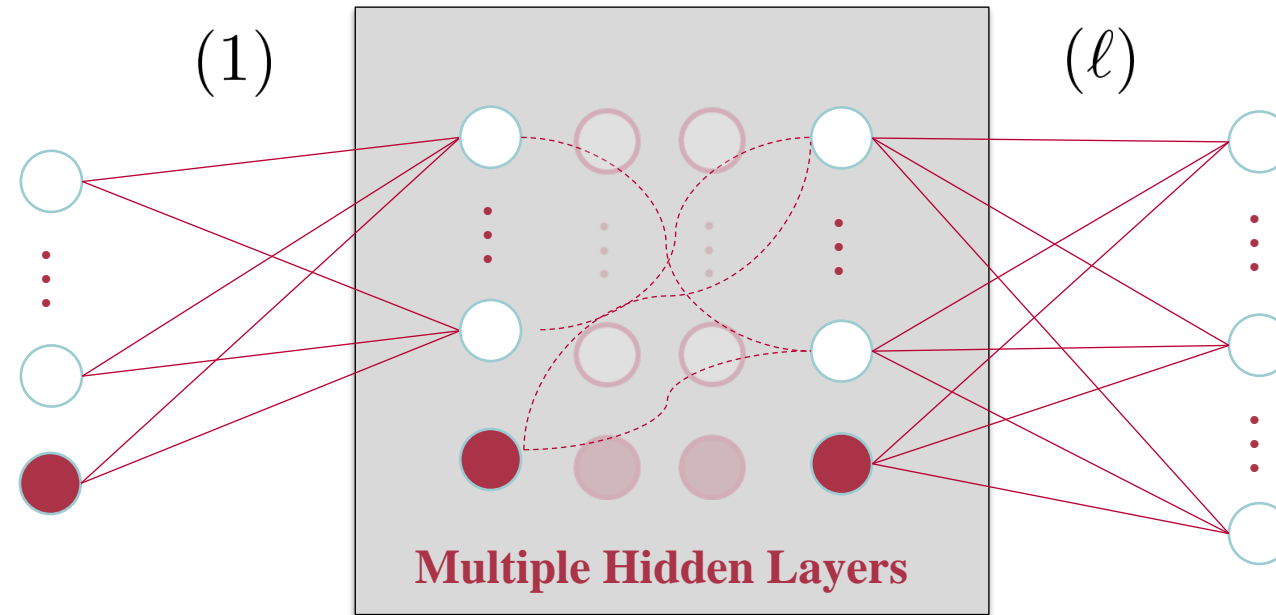
$$z_j = h(a_j), \quad j = 1, \dots, M \quad \hat{y}_k = \sigma(a_k), \quad k = 1, \dots, K$$

$$\hat{y}_k(X, \beta) = \sigma \left(\sum_{j=1}^M \beta_{kj}^{(2)} h \left(\sum_{i=1}^p \beta_{ji}^{(1)} X_i + \beta_{j0}^{(1)} \right) + \beta_{k0}^{(2)} \right)$$

$$X_0 = 1, \quad z_0 = 1$$

$$\hat{y}_k(X, \beta) = \sigma \left(\sum_{j=0}^M \beta_{kj}^{(2)} h \left(\sum_{i=0}^p \beta_{ji}^{(1)} X_i \right) \right)$$

Network Function



$$\hat{y}_k(X, \beta) = \sigma \left(\sum_j \beta_{kj}^{(\ell)} h \left(\sum_s \beta_{js}^{(\ell-1)} h \left(\dots h \left(\sum_i \beta_{ji}^{(1)} X_i \right) \dots \right) \right) \right)$$

Network Training

Error Functions

Regression

$$(x_i, y_i), i = 1, 2, \dots, n$$

$$x_i \in \mathbb{R}^p \quad y_i \in \mathbb{R}^K$$

$$E(\beta) = \sum_{i=1}^n E_i(\beta)$$

Sum of Squares

$$K = 1$$

$$E(\beta) = \frac{1}{2} \sum_{i=1}^n (\hat{y}(x_i, \beta) - y_i)^2$$

$$K > 1$$

$$E(\beta) = \frac{1}{2} \sum_{i=1}^n \|\hat{y}(x_i, \beta) - y_i\|^2 = \frac{1}{2} \sum_{i=1}^n \sum_{k=1}^K (\hat{y}_k(x_i, \beta) - y_{ik})^2$$

$$\hat{y}_k = \sigma(a_k) = a_k$$

$$\frac{\partial E_i}{\partial a_k} = \hat{y}_k(x_i, \beta) - y_{ik}$$

Network Training

Error Functions

K -Separate Binary Classification

$$(x_i, y_i), i = 1, 2, \dots, n$$
$$x_i \in \mathbb{R}^p \quad y_i \in \{0, 1\}^K$$

$$E(\beta) = \sum_{i=1}^n E_i(\beta)$$

Cross Entropy

$$K = 1$$

$$E(\beta) = - \sum_{i=1}^n (y_i \ln(\hat{y}(x_i, \beta)) + (1 - y_i) \ln(1 - \hat{y}(x_i, \beta)))$$

$$K > 1$$

$$E(\beta) = - \sum_{i=1}^n \sum_{k=1}^K (y_{ik} \ln(\hat{y}_k(x_i, \beta)) + (1 - y_{ik}) \ln(1 - \hat{y}_k(x_i, \beta)))$$

$$\hat{y}_k = \sigma(a_k) = \frac{1}{1 + e^{-a_k}}$$

$$\frac{\partial E_i}{\partial a_k} = \hat{y}_k(x_i, \beta) - y_{ik}$$

Network Training

Error Functions

Multiclass Classification

$(x_i, y_i), i = 1, 2, \dots, n \quad x_i \in \mathbb{R}^p$

y_i is a unit vector in \mathbb{R}^K

$$E(\beta) = \sum_{i=1}^n E_i(\beta)$$

Cross Entropy

$$E(\beta) = - \sum_{i=1}^n \sum_{k=1}^K y_{ik} \ln(\hat{y}_k(x_i, \beta))$$

$$\hat{y}_k = \sigma(a_k) = \frac{e^{a_k}}{\sum_{j=1}^K e^{a_j}}$$

$$\frac{\partial E_i}{\partial a_k} = \hat{y}_k(x_i, \beta) - y_{ik}$$

Network Training

Optimization

$$\min_{\beta} \frac{1}{n} \sum_{i=1}^n E_i(\beta)$$

$$\frac{1}{n} \sum_{i=1}^n \nabla E_i(\beta) = 0$$

(Batch) Gradient Descent

$$\beta_{k+1} = \beta_k - \frac{\alpha_k}{n} \sum_{i=1}^n \nabla E_i(\beta_k)$$

k : iteration
 α_k : learning rate

Stochastic Gradient Descent (SGD)

$$\beta_{k+1} = \beta_k - \alpha_k \nabla E_j(\beta_k), \quad j \in \{1, \dots, n\}$$

Mini-batch SGD

$$\beta_{k+1} = \beta_k - \frac{\alpha_k}{|\mathcal{J}|} \sum_{j \in \mathcal{J}} \nabla E_j(\beta_k), \quad \mathcal{J} \subseteq \{1, \dots, n\}$$

Network Training

Backpropagation

$$\nabla E(\beta) = \sum_{i=1}^n \nabla E_i(\beta) = 0 \qquad \frac{\partial E}{\partial \beta_{ts}^{(l)}} = \sum_{i=1}^n \frac{\partial E_i}{\partial \beta_{ts}^{(l)}} \quad ?$$

Recall

$$E(\beta) = - \sum_{i=1}^n \sum_{k=1}^K y_{ik} \ln(\hat{y}_k(x_i, \beta)) \qquad E(\beta) = \frac{1}{2} \sum_{i=1}^n \sum_{k=1}^K (\hat{y}_k(x_i, \beta) - y_{ik})^2 \qquad \frac{\partial E_i}{\partial a_k} = \hat{y}_k(x_i, \beta) - y_{ik}$$

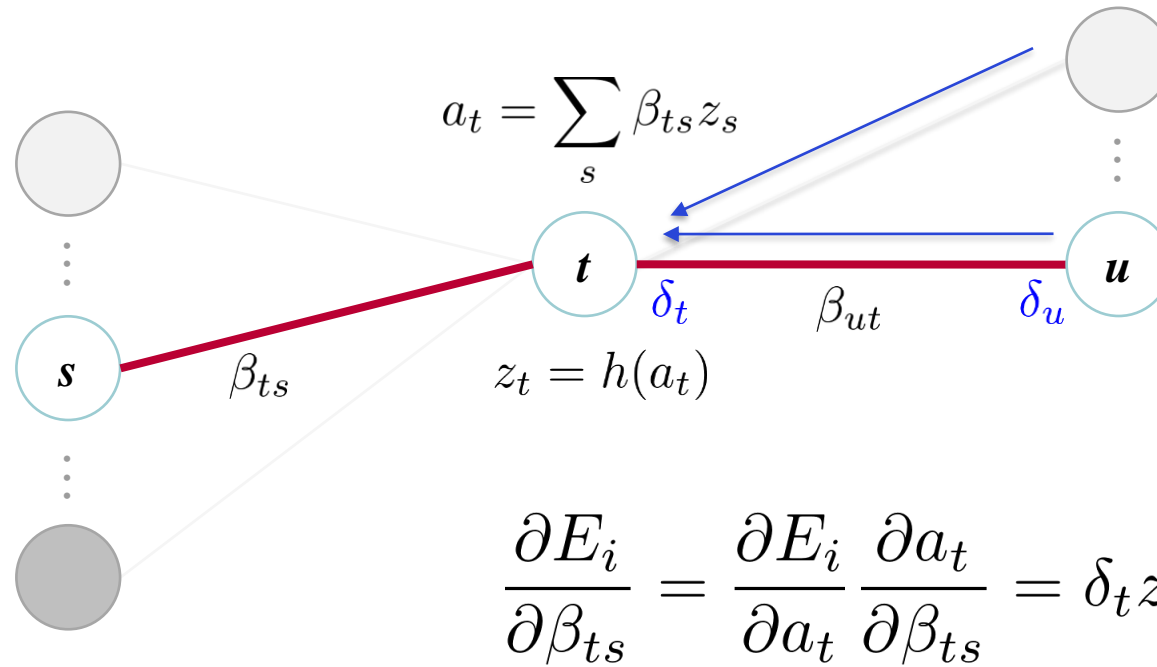
$$\hat{y}_k(X, \beta) = \sigma \left(\sum_j \beta_{kj}^{(\ell)} h \left(\sum_s \beta_{js}^{(\ell-1)} h \left(\dots h \left(\sum_i \beta_{ji}^{(1)} X_i \right) \dots \right) \right) \right)$$

Apply chain rule over and over again?

Yes but try it on the network!

Network Training

Backpropagation



$$\hat{y}_{ik} \triangleq \hat{y}_k(x_i, \beta)$$
$$\frac{\partial E_i}{\partial a_k} = \hat{y}_{ik} - y_{ik}$$

$$\frac{\partial E_i}{\partial \beta_{ts}} = \frac{\partial E_i}{\partial a_t} \frac{\partial a_t}{\partial \beta_{ts}} = \delta_t z_s$$

$$\delta_t \triangleq \frac{\partial E_i}{\partial a_t}$$

$$\delta_t = \frac{\partial E_i}{\partial a_t} = \sum_u \frac{\partial E_i}{\partial a_u} \frac{\partial a_u}{\partial a_t} = h'(a_t) \sum_u \delta_u \beta_{ut}$$

$$\delta_k = \frac{\partial E_i}{\partial a_k} = \hat{y}_{ik} - y_{ik}, \quad k = 1, \dots, K$$

Network Training

Backpropagation

Deep Learning: An Introduction for Applied Mathematicians

Catherine F. Higham* Desmond J. Higham[†]

January 19, 2018

Abstract

Multilayered artificial neural networks are becoming a pervasive tool in a host of application fields. At the heart of this deep learning revolution are familiar concepts from applied and computational mathematics; notably, in calculus, approximation theory, optimization and linear algebra. This article provides a very brief introduction to the basic ideas that underlie deep learning from an applied mathematics perspective. Our target audience includes postgraduate and final year undergraduate students in mathematics who are keen to learn about the area. The article may also be useful for instructors in mathematics who wish to enliven their classes with references to the application of deep learning techniques. We focus on three fundamental questions: what is a deep neural network? how is a network trained? what is the stochastic gradient method? We illustrate the ideas with a short MATLAB code that sets up and trains a network. We also show the use of state-of-the art software on a large scale image classification problem. We finish with references to the current literature.

Network Training

Backpropagation

Apply an input vector x_i and propagate it through the network using

$$a_t = \sum_s \beta_{ts} z_s \quad z_t = h(a_t) \quad \hat{y}_k = \sigma(a_k)$$

Evaluate the following for the output units:

$$\delta_k = \frac{\partial E_i}{\partial a_k} = \hat{y}_{ik} - y_{ik}, \quad k = 1, \dots, K$$

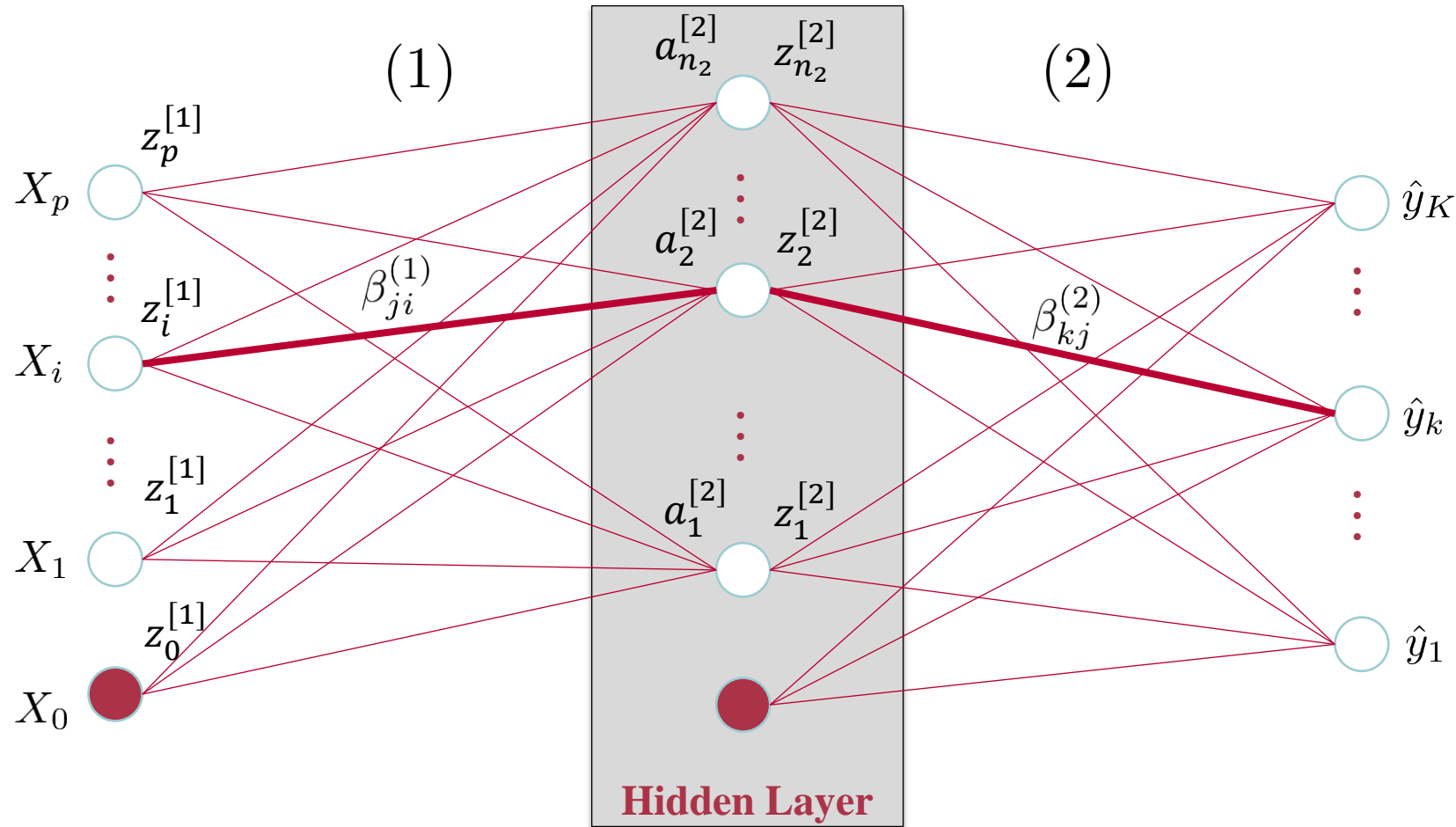
Backpropagate the following for the hidden units:

$$\delta_t = \frac{\partial E_i}{\partial a_t} = \sum_u \frac{\partial E_i}{\partial a_u} \frac{\partial a_u}{\partial a_t} = h'(a_t) \sum_u \delta_u \beta_{ut}$$

Evaluate the components of the gradient:

$$\frac{\partial E_i}{\partial \beta_{ts}} = \frac{\partial E_i}{\partial a_t} \frac{\partial a_t}{\partial \beta_{ts}} = \delta_t z_s$$

Network Training



$$l = 1, 2, \dots, L$$

$$a^{[l]} = (a_0^{[l]}, a_1^{[l]}, \dots, a_{n_l}^{[l]})$$

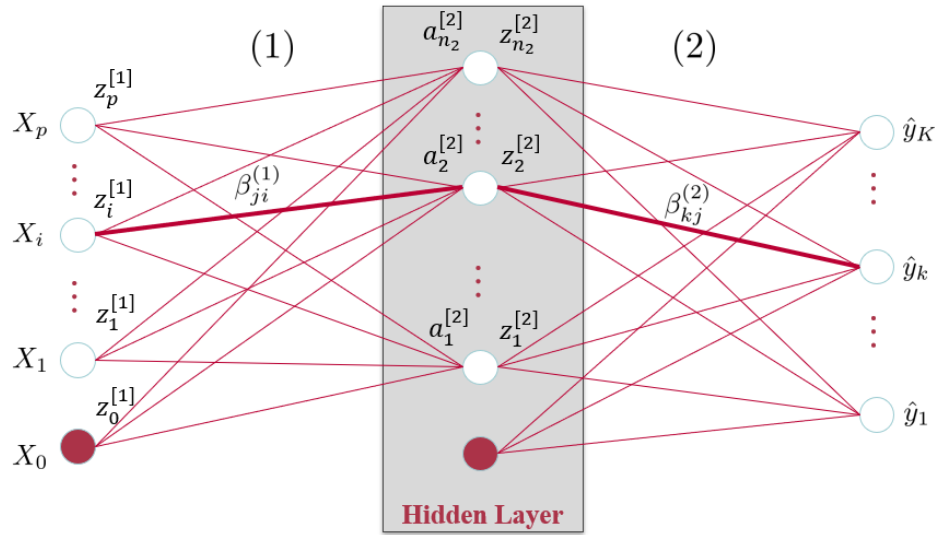
$$a_i^{[l]} = \sum_{j=0}^{n_{l-1}} \beta_{ji}^{(l-1)} z_j^{[l-1]}$$

$$a^{[l]} = W^{[l-1]} z^{[l-1]}$$

$$z^{[l]} = (z_0^{[l]}, z_1^{[l]}, \dots, z_{n_k}^{[l]})$$

$$z^{[l]} = h(a^{[l]})$$

Network Training



$$a^{[l]} = (a_0^{[l]}, a_1^{[l]}, \dots, a_{n_l}^{[l]})$$

$$a_i^{[l]} = \sum_{j=0}^{n_{l-1}} \beta_{ji}^{(l-1)} z_j^{[l-1]}$$

$$a^{[l]} = W^{[l-1]} z^{[l]}$$

$$z^{[l]} = (z_0^{[l]}, z_1^{[l]}, \dots, z_{n_k}^{[l]})$$

$$z^{[l]} = h(a^{[l]})$$

For counter 1 upto Niter

Choose an integer k uniformly at random from $\{1, 2, \dots, N\}$

x_k is current training data point

$$z^{[1]} = x_k$$

For $l = 2$ up to L

$$a^{[l]} = W^{[l-1]} z^{[l-1]}$$

$$z^{[l]} = h(a^{[l]})$$

$$D^{[l]} = \text{diag}(h'(a^{[l]}))$$

end

$$\delta^{[L]} = D^{[L]}(z^{[L]} - y_k)$$

For $l = L - 1$ down to 2

$$\delta^{[l]} = D^{[l]}(W^{[l]})^T \delta^{[l+1]}$$

end

For $l = L - 1$ down to 1

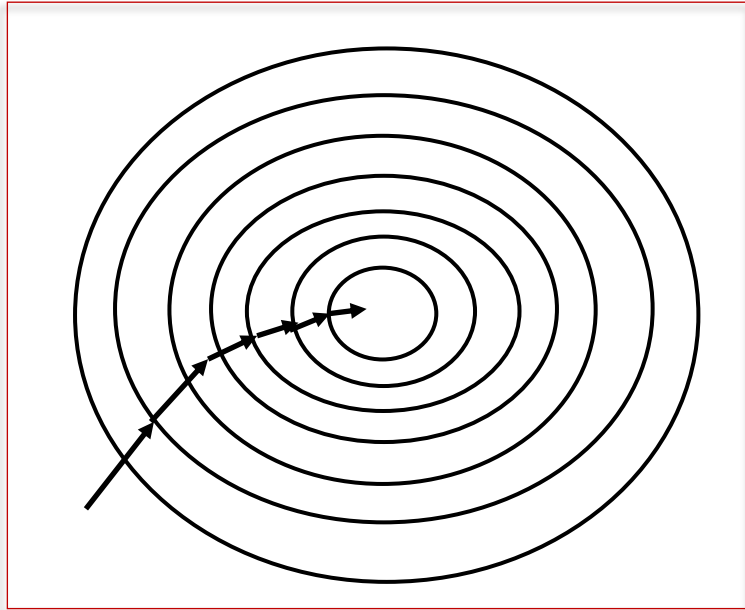
$$W^{[l]} = W^{[l]} - \eta \delta^{[l]} a^{[l-1]T}$$

end

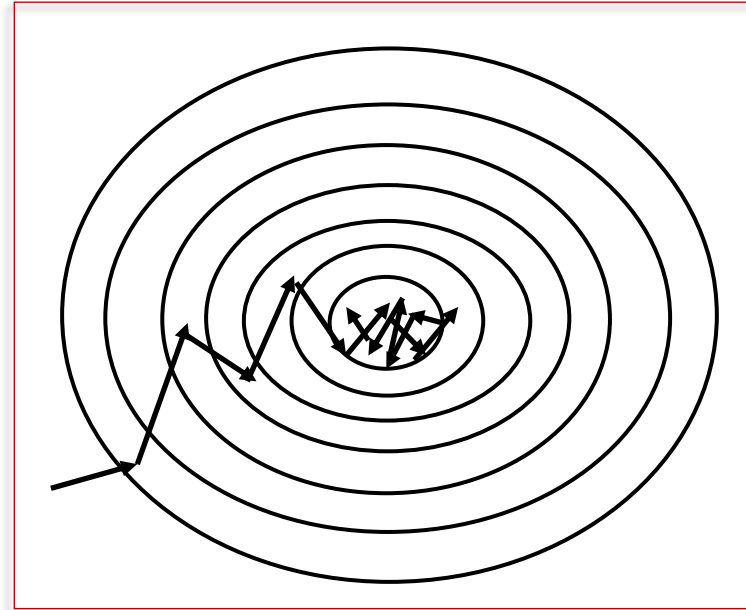
end

Network Training

GD vs. SGD



Gradient Descent

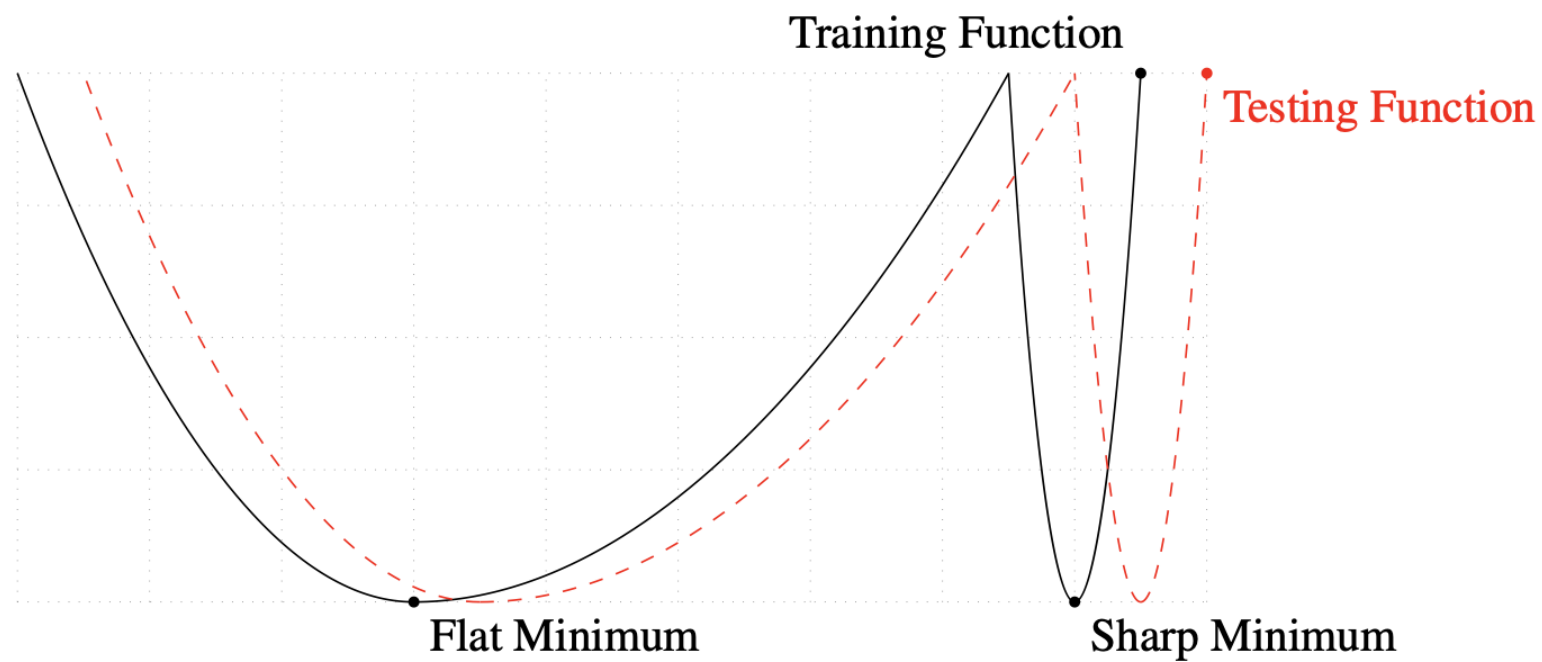


Stochastic Gradient Descent

$$\beta_{k+1} = \beta_k - \frac{\alpha_k}{|\mathcal{J}|} \sum_{j \in \mathcal{J}} \nabla E_j(\beta_k), \quad \mathcal{J} \subseteq \{1, \dots, n\}$$

Network Training

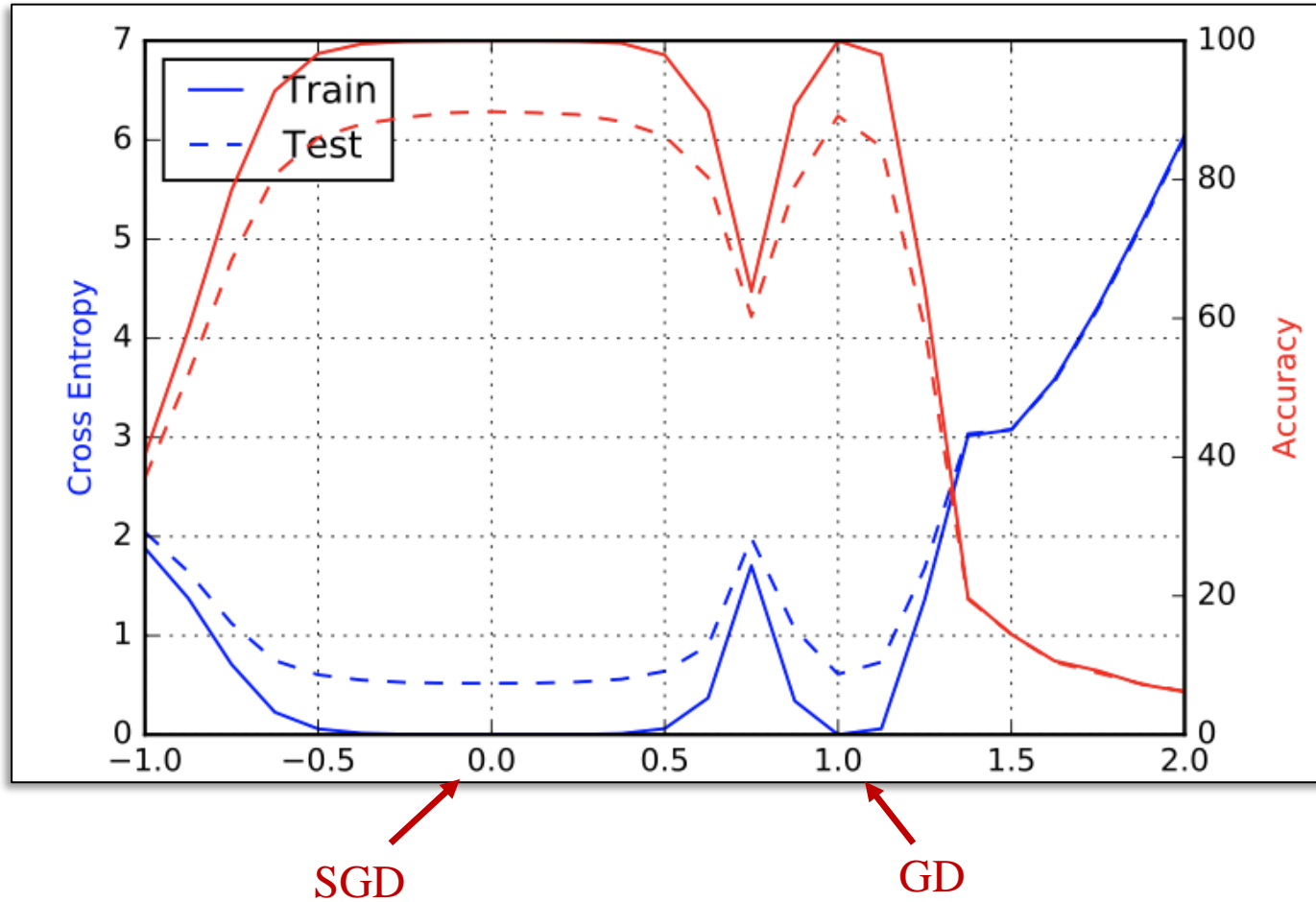
Flat vs. Sharp Minima



*Keskar et al. 2017

Network Training

Flat vs. Sharp Minima



Deep Convolutional
Neural Net on CIFAR-10

*Keskar et al. 2017

$$\min_{\beta} \underbrace{\frac{1}{n} \sum_{i=1}^n E_i(\beta)}_{F(\beta)} \equiv_{\beta \leftrightarrow w} \min_w F(w)$$

$$g(w_k, \xi_k) = \frac{1}{|\mathcal{J}|} \sum_{j \in \mathcal{J}_{\xi_k}} \nabla E_j(w_k), \quad \mathcal{J}_{\xi_k} \subseteq \{1, \dots, n\}$$

$$w_{k+1} = w_k - \alpha_k g(w_k, \xi_k)$$

Assumption 1. Let $F : \mathbb{R}^p \rightarrow \mathbb{R}$ be continuously differentiable and the gradient $\nabla F : \mathbb{R}^p \rightarrow \mathbb{R}^p$ be Lipschitz continuous with Lipschitz constant L . That is,

$$\|\nabla F(w) - \nabla F(v)\|_2 \leq L\|w - v\|_2 \quad \text{for all } w, v \in \mathbb{R}^p.$$

Proposition 1. Under Assumption 1, $F : \mathbb{R}^p \rightarrow \mathbb{R}$ satisfies

$$F(w) \leq F(v) + \nabla F(v)^\top (w - v) + \frac{L}{2} \|w - v\|_2^2 \quad \text{for any } w, v \in \mathbb{R}^p.$$

Proof. Exercise...

Definition. $\mathbb{E}_{\xi_k} [\cdot] = \mathbb{E} [\cdot \mid w_k]$

Lemma 1. Under Assumption 1, the iterates of SGD satisfies

$$\mathbb{E}_{\xi_k} [F(w_{k+1})] - F(w_k) \leq -\alpha_k \nabla F(w_k)^\top \mathbb{E}_{\xi_k} [g(w_k, \xi_k)] + \frac{1}{2} \alpha_k^2 L \mathbb{E}_{\xi_k} [\|g(w_k, \xi_k)\|_2^2]$$

for all $k = 1, \dots, T$.

Proof. By Assumption 1 and Proposition 1, we have

$$\begin{aligned} F(w_{k+1}) - F(w_k) &\leq \nabla F(w_k)^\top (w_{k+1} - w_k) + \frac{L}{2} \|w_{k+1} - w_k\|_2^2 \\ &\leq -\alpha_k \nabla F(w_k)^\top g(w_k, \xi_k) + \frac{1}{2} \alpha_k^2 L \|g(w_k, \xi_k)\|_2^2. \end{aligned}$$

Taking expectation with respect to ξ_k and noting that w_k does not depend on ξ_k , we obtain

$$\begin{aligned} \mathbb{E}_{\xi_k} [F(w_{k+1})] - F(w_k) &\leq -\alpha_k \mathbb{E}_{\xi_k} [F(w_k)^\top g(w_k, \xi_k)] + \frac{1}{2} \alpha_k^2 L \|g(w_k, \xi_k)\|_2^2 \\ &= -\alpha_k \nabla F(w_k)^\top \mathbb{E}_{\xi_k} [g(w_k, \xi_k)] + \frac{1}{2} \alpha_k^2 L \mathbb{E}_{\xi_k} [\|g(w_k, \xi_k)\|_2^2] \quad \blacksquare \end{aligned}$$

Assumption 2. There exists $M \geq 0$ and $M_G \geq 1$ such that

$$\mathbb{E}_{\xi_k} [\|g(w_k, \xi_k)\|_2^2] \leq M + M_G \|\nabla F(w_k)\|_2^2$$

Lemma 2. Under Assumption 1 and Assumption 2, the iterates of SGD satisfies

$$\mathbb{E}_{\xi_k} [F(w_{k+1})] - F(w_k) \leq -\left(1 - \frac{1}{2}\alpha_k L M_G\right)\alpha_k \|\nabla F(w_k)\|_2^2 + \frac{1}{2}\alpha_k^2 L M$$

for all $k = 1, \dots, T$.

Proof. Using Lemma 1 and the assumptions leads to the desired result:

$$\begin{aligned} \mathbb{E}_{\xi_k} [F(w_{k+1})] - F(w_k) &\leq -\alpha_k \nabla F(w_k)^\top \mathbb{E}_{\xi_k} [g(w_k, \xi_k)] + \frac{1}{2}\alpha_k^2 L \mathbb{E}_{\xi_k} [\|g(w_k, \xi_k)\|_2^2] \\ &\leq -\alpha_k \|\nabla F(w_k)\|_2^2 + \frac{1}{2}\alpha_k^2 L (M + M_G \|\nabla F(w_k)\|_2^2) \\ &= -\left(1 - \frac{1}{2}\alpha_k L M_G\right)\alpha_k \|\nabla F(w_k)\|_2^2 + \frac{1}{2}\alpha_k^2 L M \quad \blacksquare \end{aligned}$$

Definition (Total Expectation). $\mathbb{E}[F(w_k)] = \mathbb{E}_{\xi_1} [\mathbb{E}_{\xi_2} [\dots \mathbb{E}_{\xi_{k-1}} [F(w_k)]]]$

Theorem 1 (nonconvex, fixed stepsize). Under Assumption 1 and Assumption 2, suppose SGD is run with fixed stepsize, $\alpha_k = \alpha$, $k = 1, \dots, T$ satisfying $0 < \alpha \leq \frac{1}{LM_G}$. Then, we have

$$\mathbb{E} \left[\frac{1}{T} \sum_{k=1}^T \|\nabla F(w_k)\|_2^2 \right] \leq \alpha LM + \frac{2(F(w_1) - F_*)}{T\alpha},$$

where F_* is the minimum value of F .

Theorem 1 (nonconvex, fixed stepsize). Under Assumption 1 and Assumption 2, suppose SGD is run with fixed stepsize, $\alpha_k = \alpha$, $k = 1, \dots, T$ satisfying $0 < \alpha \leq \frac{1}{LM_G}$. Then, we have

$$\mathbb{E} \left[\frac{1}{T} \sum_{k=1}^T \|\nabla F(w_k)\|_2^2 \right] \leq \alpha LM + \frac{2(F(w_1) - F_*)}{T\alpha},$$

where F_* is the minimum value of F .

Proof. Using Lemma 2 and taking the total expectation yields

$$\begin{aligned} \mathbb{E} [F(w_{k+1})] - \mathbb{E} [F(w_k)] &\leq -\left(1 - \frac{1}{2}\alpha LM_G\right)\alpha \mathbb{E} [\|\nabla F(w_k)\|_2^2] + \frac{1}{2}\alpha^2 LM \\ &\leq -\frac{1}{2}\alpha \mathbb{E} [\|\nabla F(w_k)\|_2^2] + \frac{1}{2}\alpha^2 LM. \end{aligned}$$

Then, summing over $k = 1, \dots, T$ leads to

$$F_* - F(w_1) \leq \mathbb{E} [F(w_{T+1})] - F(w_1) \leq -\frac{1}{2}\alpha \sum_{k=1}^T \mathbb{E} [\|\nabla F(w_k)\|_2^2] + \frac{1}{2}T\alpha^2 M.$$

This implies

$$\sum_{k=1}^T \mathbb{E} [\|\nabla F(w_k)\|_2^2] \leq T\alpha LM + \frac{2(F(w_1) - F_*)}{\alpha}. \quad \blacksquare$$

Theorem 2 (nonconvex, diminishing stepsize). Under Assumption 1 and Assumption 2, suppose SGD is run with stepsizes satisfying

$$\sum_{k=1}^{\infty} \alpha_k = \infty \quad \text{and} \quad \sum_{k=1}^{\infty} \alpha_k^2 < \infty.$$

Then, we have

$$\lim_{T \rightarrow \infty} \mathbb{E} \left[\sum_{k=1}^T \alpha_k \|\nabla F(w_k)\|_2^2 \right] < \infty.$$

Therefore,

$$\mathbb{E} \left[\frac{1}{A_T} \sum_{k=1}^T \alpha_k \|\nabla F(w_k)\|_2^2 \right] \rightarrow 0 \quad \text{as} \quad T \rightarrow \infty,$$

where $A_T = \sum_{k=1}^T \alpha_k$.

Proof. Since $\alpha_k \rightarrow 0$ as $k \rightarrow \infty$, we can assume that there exists $k \in \mathbb{N}$ satisfying $\alpha_k LM_G \leq 1$. Using now Lemma 2 and taking total expectation gives

$$\begin{aligned}\mathbb{E}[F(w_{k+1})] - \mathbb{E}[F(w_k)] &\leq -\left(1 - \frac{1}{2}\alpha_k LM_G\right)\alpha_k \mathbb{E}[\|\nabla F(w_k)\|_2^2] + \frac{1}{2}\alpha_k^2 LM \\ &\leq -\frac{1}{2}\alpha_k \mathbb{E}[\|\nabla F(w_k)\|_2^2] + \frac{1}{2}\alpha_k^2 LM\end{aligned}$$

Then, summing over $k = 1, \dots, T$ leads to

$$F_* - F(w_1) \leq \mathbb{E}[F(w_{T+1})] - F(w_1) \leq -\frac{1}{2} \sum_{k=1}^T \alpha_k \mathbb{E}[\|\nabla F(w_k)\|_2^2] + \frac{1}{2} LM \sum_{k=1}^T \alpha_k^2.$$

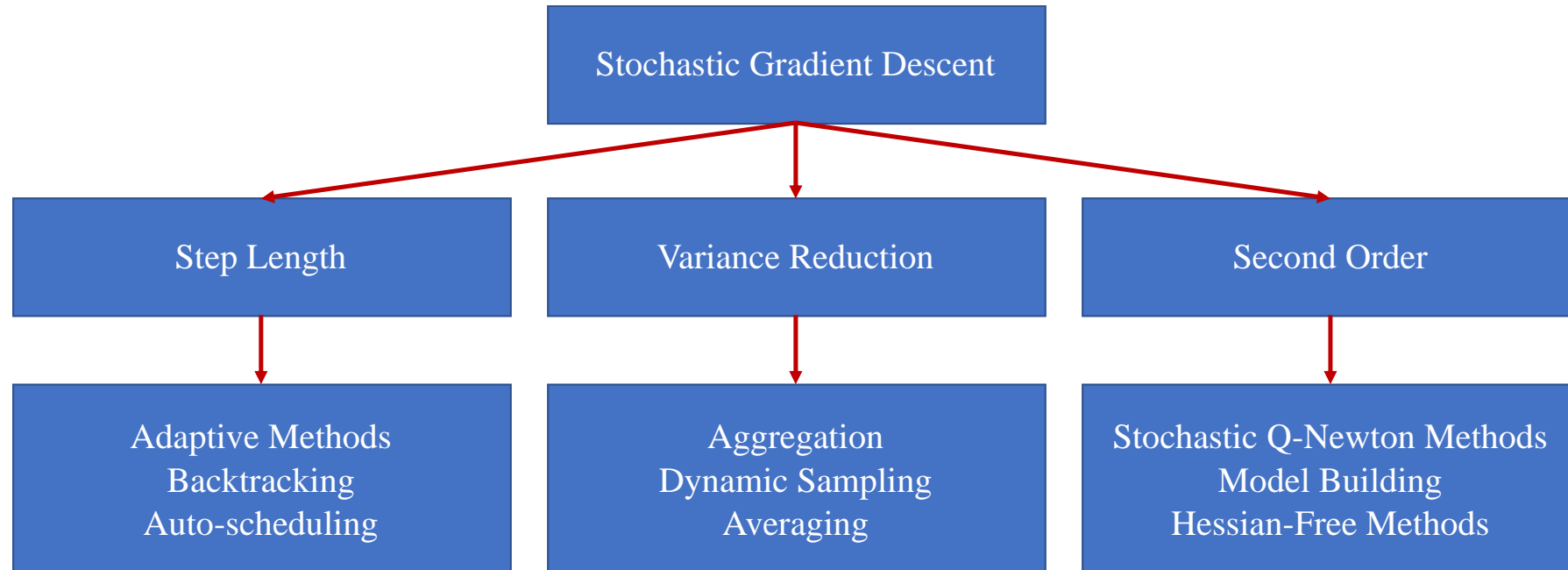
This implies

$$\sum_{k=1}^T \alpha_k \mathbb{E}[\|\nabla F(w_k)\|_2^2] \leq \underbrace{2(F(w_1) - F_*) + LM \sum_{k=1}^T \alpha_k^2}_{< \infty \text{ as } T \rightarrow \infty}.$$

The second result also follows since $A_T \rightarrow \infty$ as $T \rightarrow \infty$. ■

SGD Variants

$$w_{k+1} = w_k - \alpha_k g(w_k, \xi_k)$$



- Generalization
- Work complexity
- Probabilistic results
- ...

Optional Readings

Published as a conference paper at ICLR 2017

ON LARGE-BATCH TRAINING FOR DEEP LEARNING: GENERALIZATION GAP AND SHARP MINIMA

Nitish Shirish Keskar*
Northwestern University
Evanston, IL 60208
keskar.nitish@u.northwestern.edu

Jorge Nocedal
Northwestern University
Evanston, IL 60208
j-nocedal@northwestern.edu

Ping Tak Peter Tang
Intel Corporation
Santa Clara, CA 95054
peter.tang@intel.com

Dheevatsa Mudigere
Intel Corporation
Bangalore, India
dheevatsa.mudigere@intel.com

Mikhail Smelyanskiy
Intel Corporation
Santa Clara, CA 95054
mikhail.smelyanskiy@intel.com

Optimization Methods for Large-Scale Machine Learning

Léon Bottou* Frank E. Curtis† Jorge Nocedal‡

June 23, 2018

4	Analyses of Stochastic Gradient Methods	21
4.1	Two Fundamental Lemmas	23
4.2	SG for Strongly Convex Objectives	25
4.3	SG for General Objectives	31
4.4	Work Complexity for Large-Scale Learning	34
4.5	Commentary	37
5	Noise Reduction Methods	40
5.1	Reducing Noise at a Geometric Rate	41
5.2	Dynamic Sample Size Methods	42
5.2.1	Practical Implementation	44
5.3	Gradient Aggregation	46
5.3.1	SVRG	46
5.3.2	SAGA	47
5.3.3	Commentary	49
5.4	Iterate Averaging Methods	49
6	Second-Order Methods	50
6.1	Hessian-Free Inexact Newton Methods	52
6.1.1	Subsampled Hessian-Free Newton Methods	53
6.1.2	Dealing with Nonconvexity	55
6.2	Stochastic Quasi-Newton Methods	55
6.2.1	Deterministic to Stochastic	56
6.2.2	Algorithms	57
6.3	Gauss-Newton Methods	59
6.4	Natural Gradient Method	61
6.5	Methods that Employ Diagonal Scalings	64
7	Other Popular Methods	68
7.1	Gradient Methods with Momentum	68
7.2	Accelerated Gradient Methods	70
7.3	Coordinate Descent Methods	70
8	Methods for Regularized Models	74
8.1	First-Order Methods for Generic Convex Regularizers	75
8.1.1	Iterative Soft-Thresholding Algorithm (ISTA)	77
8.1.2	Bound-Constrained Methods for ℓ_1 -norm Regularized Problems	77
8.2	Second-Order Methods	78
8.2.1	Proximal Newton Methods	79
8.2.2	Orthant-Based Methods	80

Network Training

Overfitting

- Cross validation
- **Dropout method:** randomly deleting some nodes and their connections from the network ([link](#))
- Regularization

$$\min_{\beta} E(\beta) + \lambda \sum_{s,t} \beta_{st}^2 \quad \text{or} \quad \min_{\beta} E(\beta) + \lambda \sum_{s,t} |\beta_{st}|$$

Vanishing/Exploding Gradients

- Gradient gets very small (large) moving from last layers to earlier layers
- Training takes a long time or simply fails

$$\frac{\partial E_i}{\partial \beta_{ts}} = \delta_t z_s$$

$$= \left(\sigma'(a_t) \sum_u \delta_u \beta_{ut} \right) z_s$$

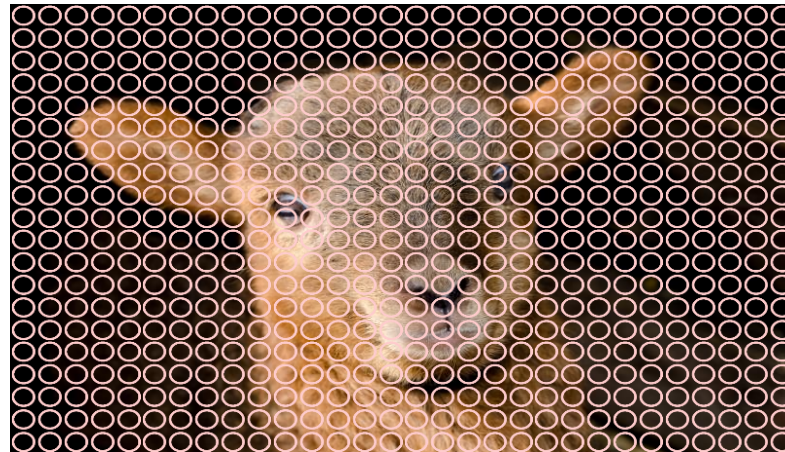
$$= \left(\underbrace{\sigma'(a_t)}_{< 1} \sum_u \left(\underbrace{\sigma'(a_u)}_{< 1} \sum_v \delta_v \beta_{vu} \right) \beta_{ut} \right) \underbrace{z_s}_{\leq 1}$$

$$\sigma(a) = \frac{1}{1 + e^{-a}} \implies \max_{a \in \mathbb{R}} \sigma'(a) = \frac{1}{4}$$

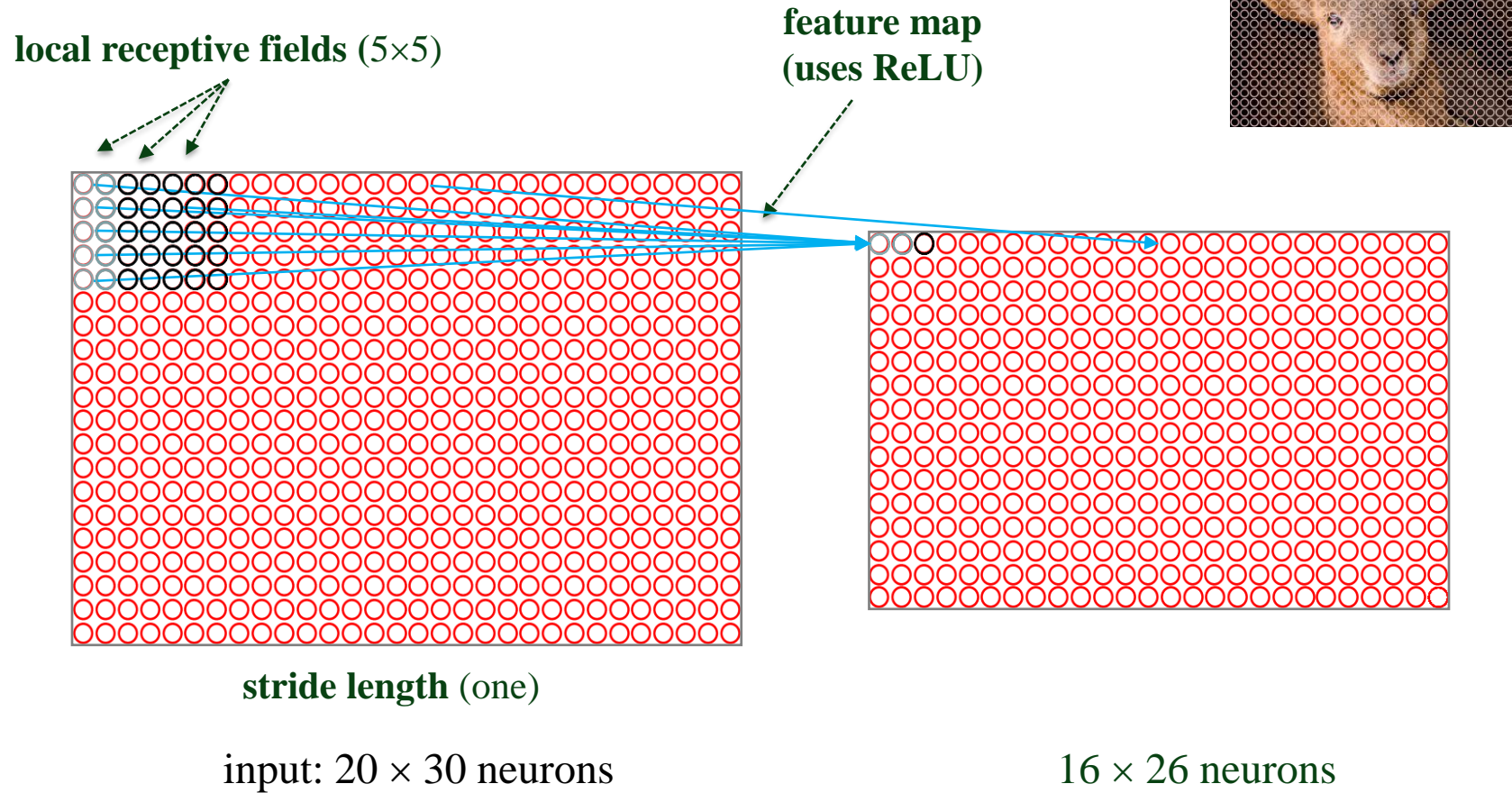
- Smaller initial weights lead to vanishing gradients
- Larger initial weights lead to explosion of gradients

Convolutional Neural Network (CNN)

- Carefully constructed multi-layer neural network
- Exploits the spatial structure of data (object recognition)
- First layers learn simple patterns, subsequent layers recombine them
- Hidden layers detect the same local structure
- ReLU is used for activation
- The number of weights and biases decreases significantly

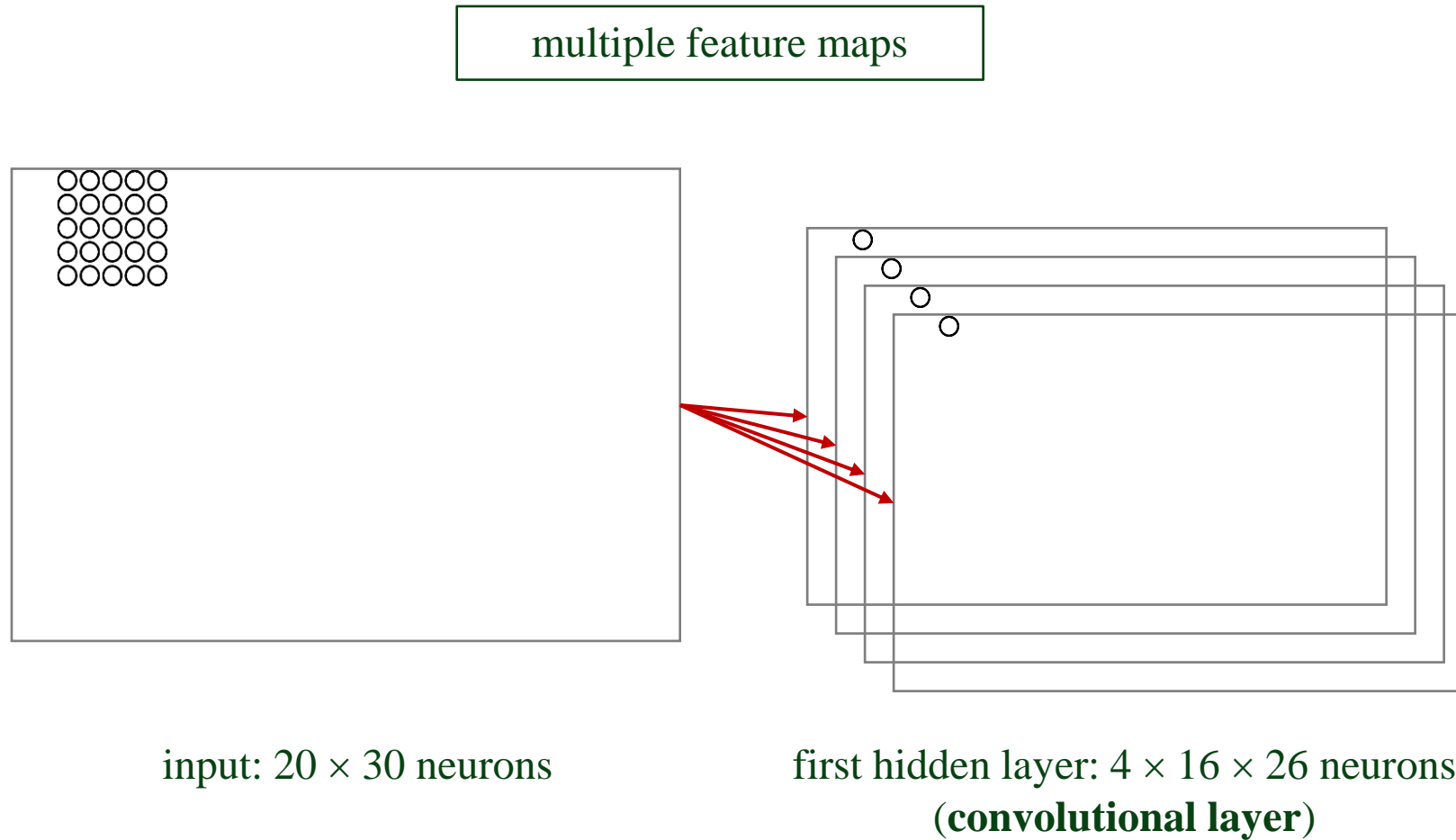


CNN - Construction



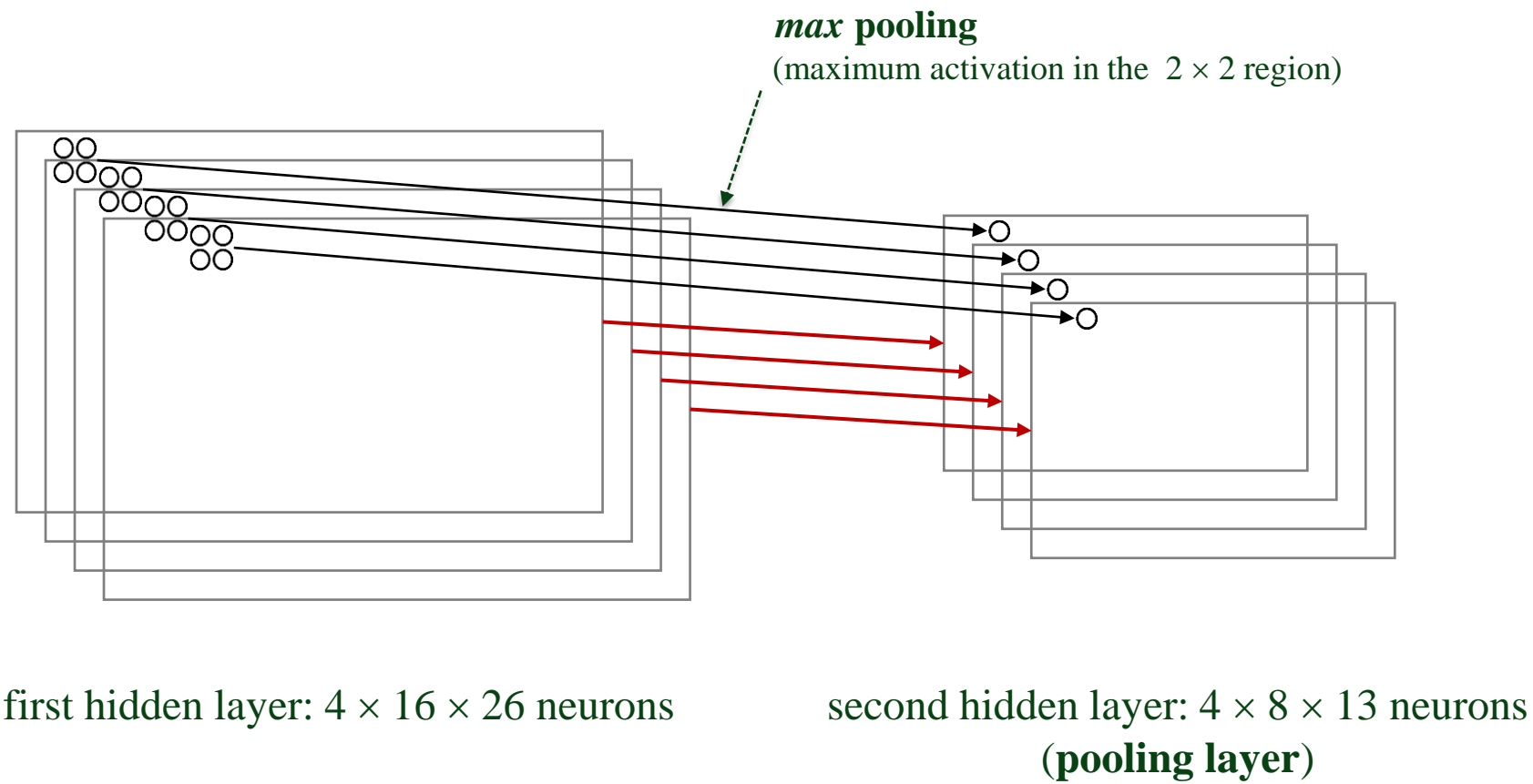
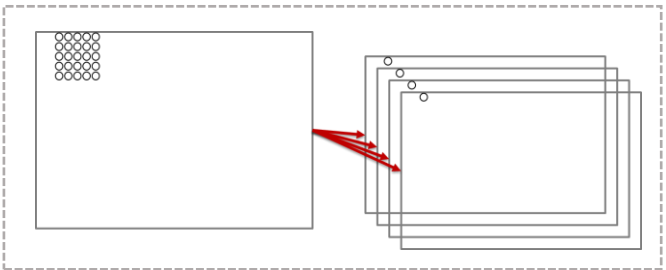
weights and bias are shared for each hidden neuron
(shared weights and bias define a **filter** or a **kernel**)

CNN – Hidden Layer

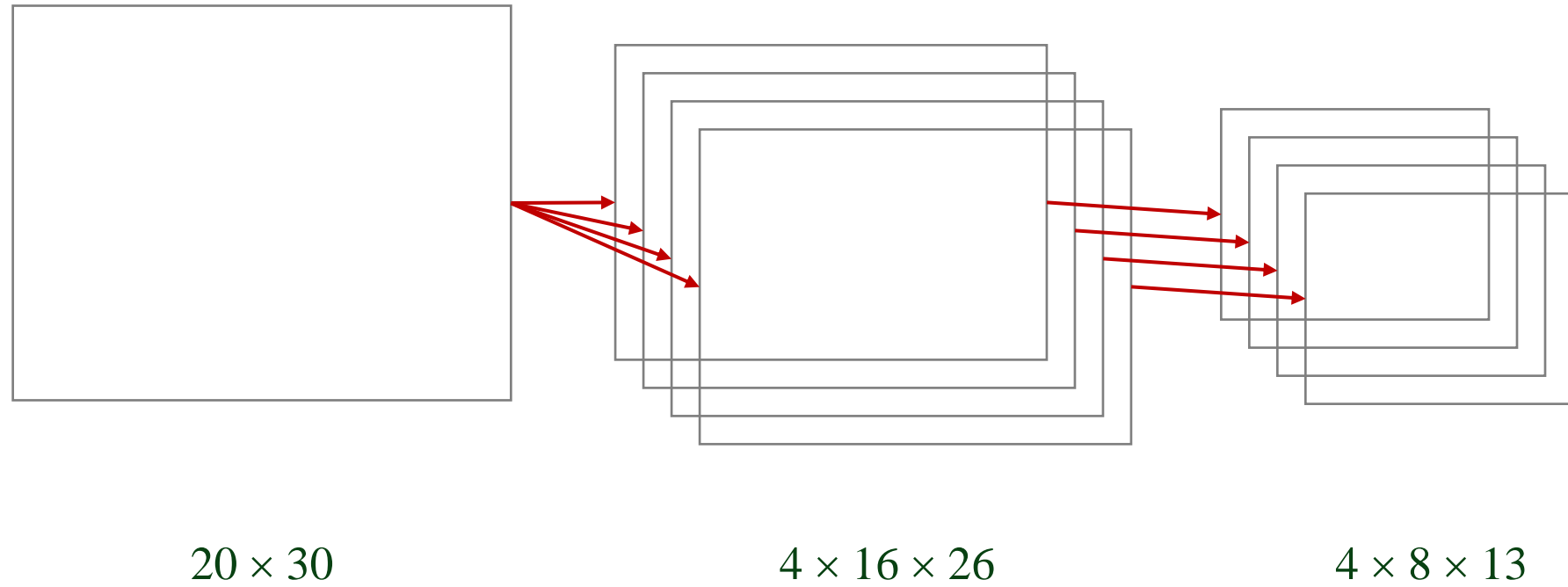


$$\text{number of weights and biases} = (25 + 1) \times 4 = 104$$

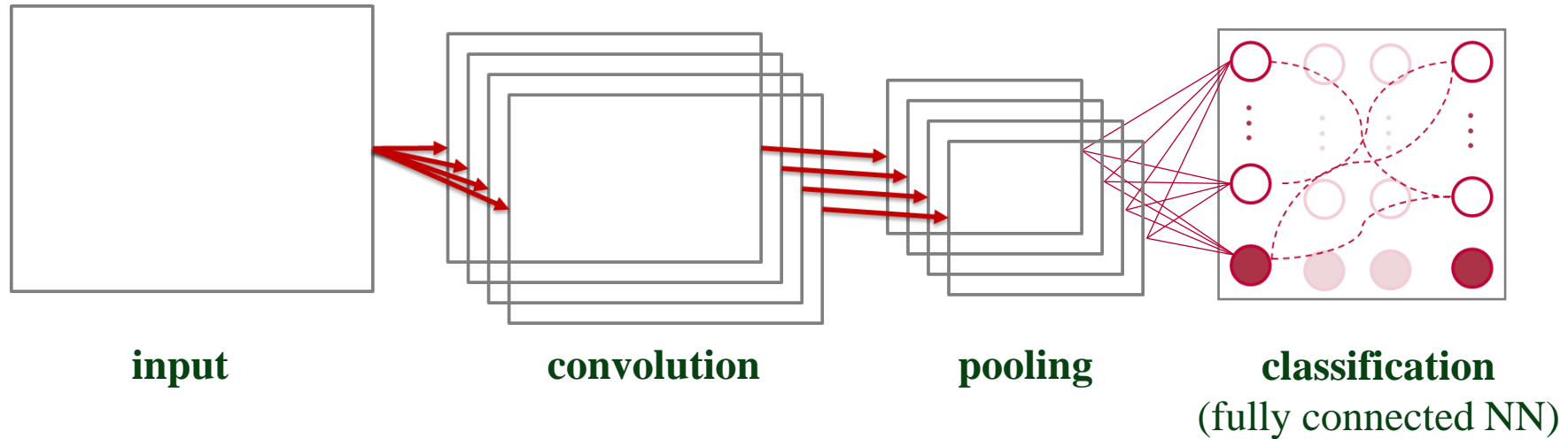
CNN – Pooling Layer



CNN – Merging Layers

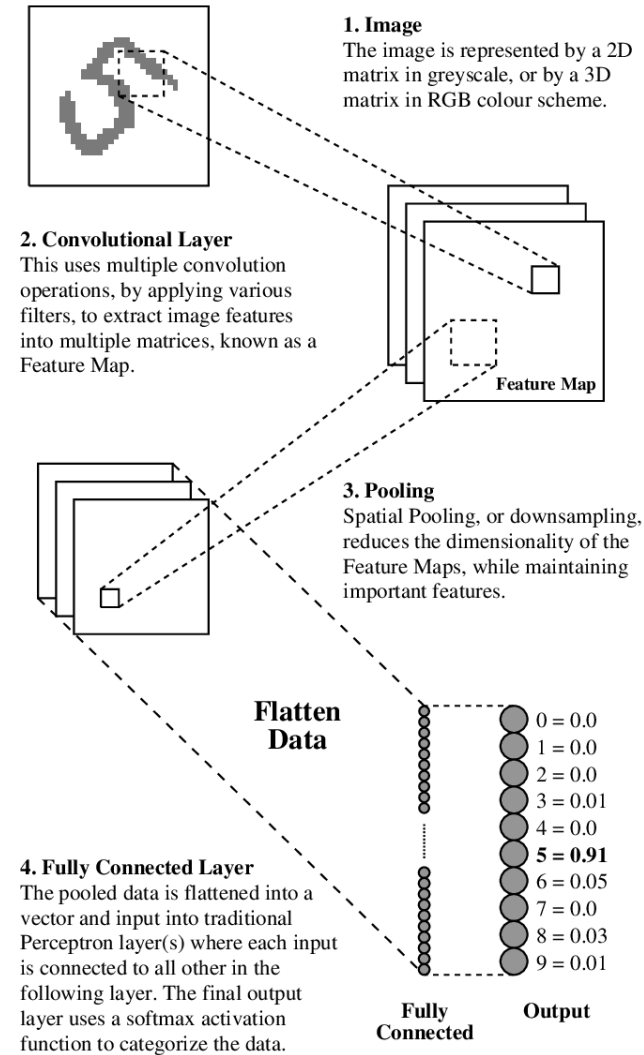


CNN – Complete Network



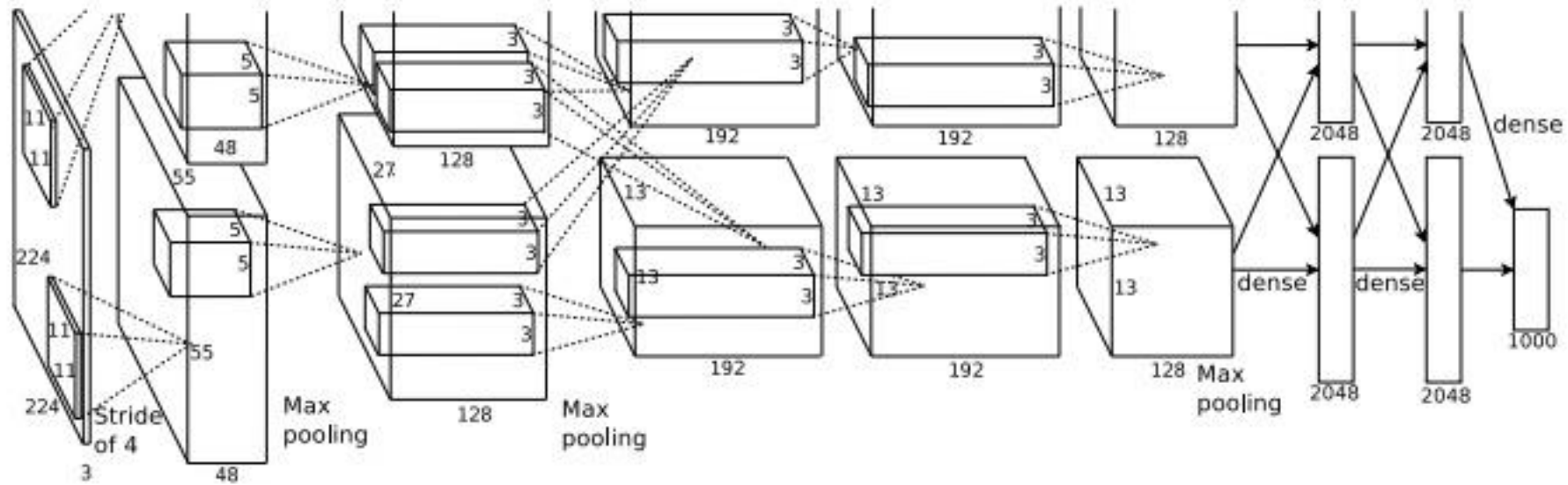
- Common to use one layer for classification in the last stage
- There could be multiple convolution layers
- Too many matrix or vector multiplications are required
- Vectorization along with GPUs are used (for all deep learning networks)

CNN – Object Recognition



([source](#))

CNN – ImageNet Classification

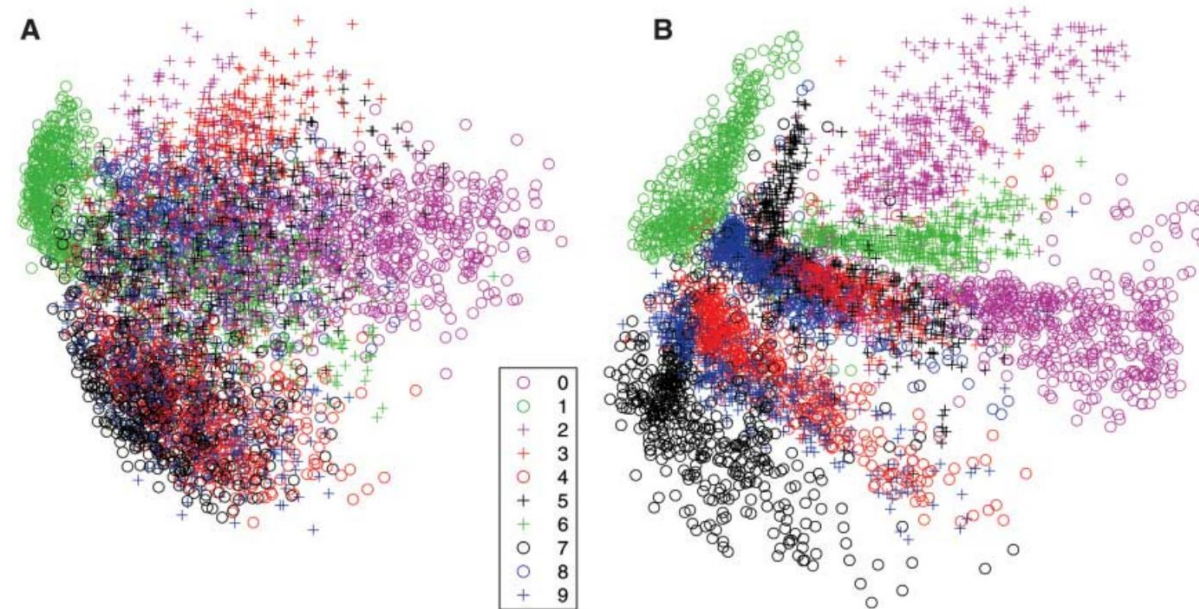


(source)

Autoencoders

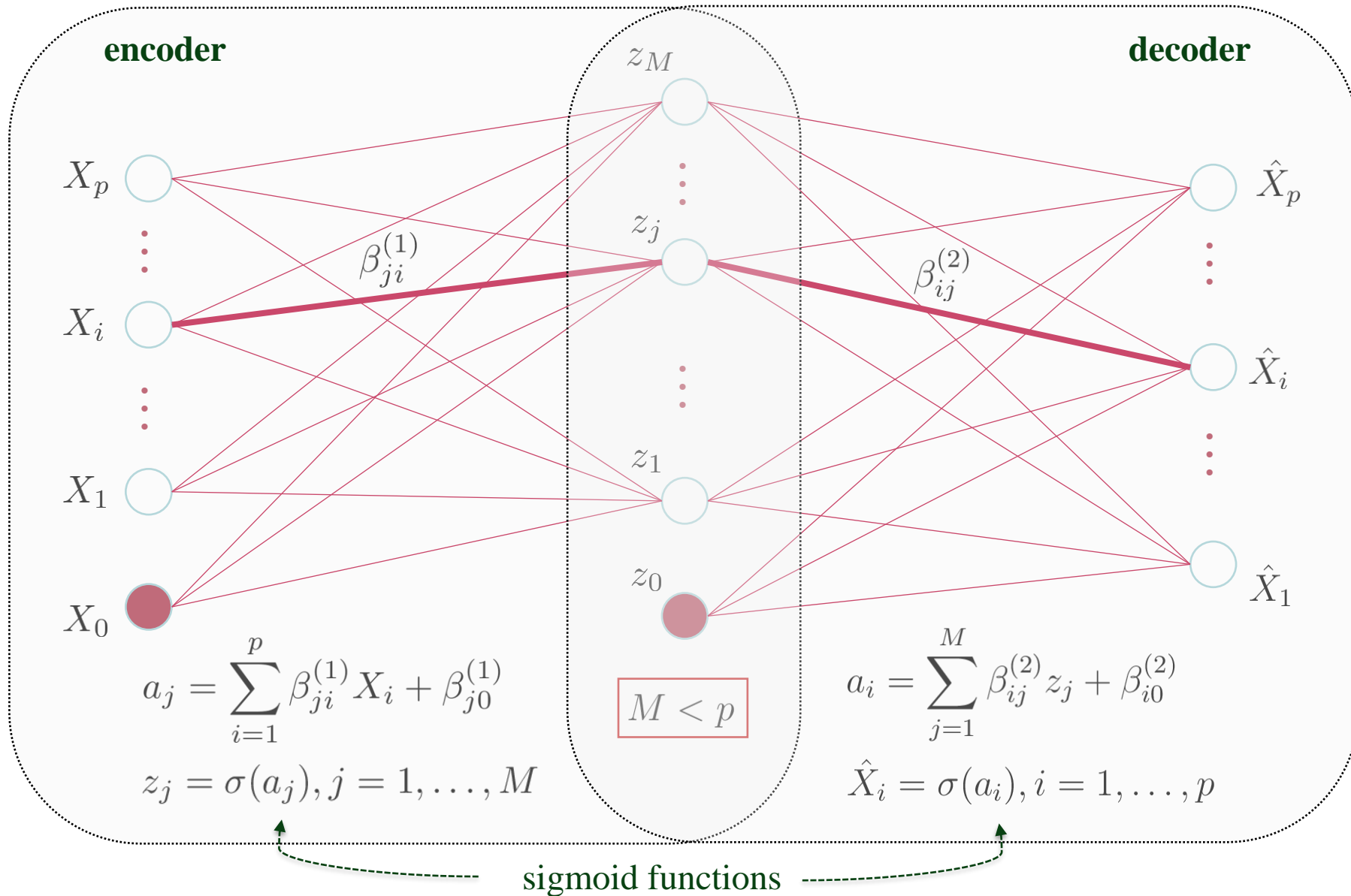
- Unsupervised learning
- Motto: Extract features, reconstruct input
- Dimension reduction technique

Fig. 3. (A) The two-dimensional codes for 500 digits of each class produced by taking the first two principal components of all 60,000 training images. (B) The two-dimensional codes found by a 784-1000-500-250-2 autoencoder. For an alternative visualization, see (8).



()

Autoencoders



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

Perceptron (P)



Feed Forward (FF)



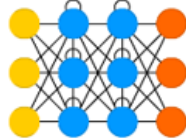
Radial Basis Network (RBF)



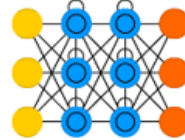
Deep Feed Forward (DFF)



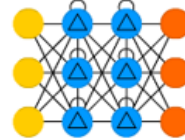
Recurrent Neural Network (RNN)



Long / Short Term Memory (LSTM)



Gated Recurrent Unit (GRU)



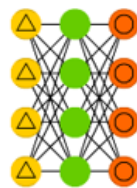
Auto Encoder (AE)



Variational AE (VAE)



Denoising AE (DAE)



Sparse AE (SAE)



Markov Chain (MC)



Hopfield Network (HN)



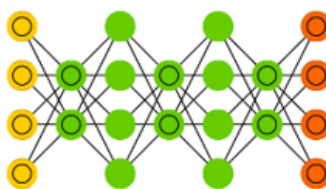
Boltzmann Machine (BM)



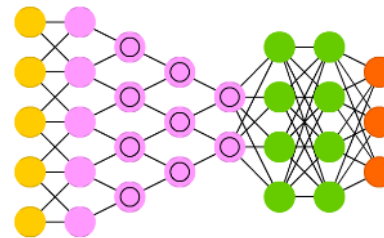
Restricted BM (RBM)



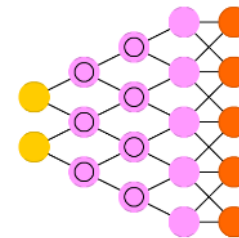
Deep Belief Network (DBN)



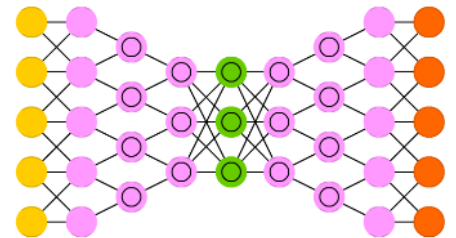
Deep Convolutional Network (DCN)



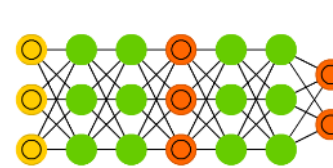
Deconvolutional Network (DN)



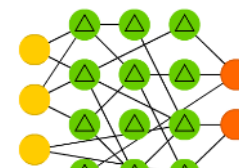
Deep Convolutional Inverse Graphics Network (DCIGN)



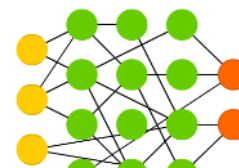
Generative Adversarial Network (GAN)



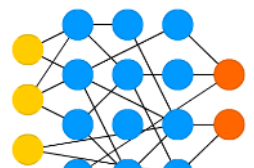
Liquid State Machine (LSM)



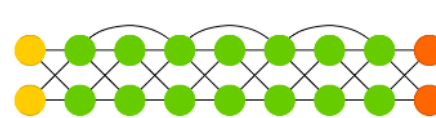
Extreme Learning Machine (ELM)



Echo State Network (ESN)



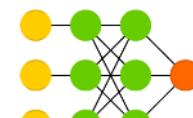
Deep Residual Network (DRN)



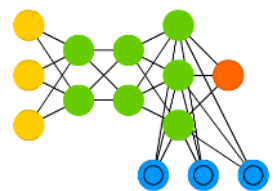
Kohonen Network (KN)



Support Vector Machine (SVM)



Neural Turing Machine (NTM)



(source - 2016)