

## Support Vector Machines for Regression

- Consider regressing an unidimensional continuous random variable on a  $D$ -dimensional continuous random variable.
- Consider a training set  $\{(\mathbf{x}_n, t_n)\}$ . Consider using the linear model 
$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b$$
- To get a sparse solution, instead of minimizing the classical regularized error function

$$\frac{1}{2} \sum_n (y(\mathbf{x}_n) - t_n)^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

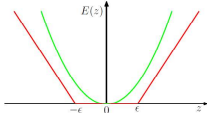
consider minimizing the  $\epsilon$ -insensitive regularized error function

$$C \sum_n E_\epsilon(y(\mathbf{x}_n) - t_n) + \frac{1}{2} \|\mathbf{w}\|^2$$

where  $C > 0$  controls regularization and

$$E_\epsilon(y(\mathbf{x}) - t) = \begin{cases} 0 & \text{if } |y(\mathbf{x}) - t| \leq \epsilon \\ |y(\mathbf{x}) - t| - \epsilon & \text{otherwise} \end{cases}$$

**Figure 7.6** Plot of an  $\epsilon$ -insensitive error function (in red) in which the error increases linearly with distance beyond the insensitive region. Also shown for comparison is the quadratic error function (in green).



13/18

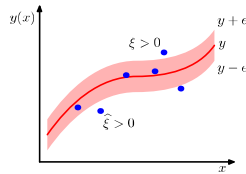
## Support Vector Machines for Regression

- The values of  $C$  and  $\epsilon$  can be decided by cross-validation.
- Consider the slack variables  $\xi_n \geq 0$  and  $\widehat{\xi}_n \geq 0$  such that

$$\xi_n = \begin{cases} t_n - y(\mathbf{x}_n) - \epsilon & \text{if } t_n > y(\mathbf{x}_n) + \epsilon \\ 0 & \text{otherwise} \end{cases}$$

and

$$\widehat{\xi}_n = \begin{cases} y(\mathbf{x}_n) - \epsilon - t_n & \text{if } t_n < y(\mathbf{x}_n) - \epsilon \\ 0 & \text{otherwise} \end{cases}$$



14/18

## Support Vector Machines for Regression

- The optimal regression curve is given by

$$\arg \min_{\mathbf{w}, b, \{\xi_n\}, \{\widehat{\xi}_n\}} C \sum_n (\xi_n + \widehat{\xi}_n) + \frac{1}{2} \|\mathbf{w}\|^2$$

subject to  $\xi \geq 0$ ,  $\widehat{\xi}_n \geq 0$ ,  $t_n \leq y(\mathbf{x}_n) + \epsilon + \xi_n$  and  $t_n \geq y(\mathbf{x}_n) - \epsilon - \widehat{\xi}_n$ .

- To minimize the previous expression, we minimize

$$C \sum_n (\xi_n + \widehat{\xi}_n) + \frac{1}{2} \|\mathbf{w}\|^2 - \sum_n (\mu_n \xi_n + \widehat{\mu}_n \widehat{\xi}_n)$$

$$- \sum_n (y(\mathbf{x}_n) + \epsilon + \xi_n - t_n) - \sum_n (\widehat{\mu}_n (t_n - y(\mathbf{x}_n) + \epsilon + \widehat{\xi}_n))$$

where  $\mu_n \geq 0$ ,  $\widehat{\mu}_n \geq 0$ ,  $a_n \geq 0$  and  $\widehat{a}_n \geq 0$  are Lagrange multipliers.

- Setting its derivatives with respect to  $\mathbf{w}$ ,  $b$ ,  $\xi_n$  and  $\widehat{\xi}_n$  to zero gives

$$\mathbf{w} = \sum_n (\widehat{a}_n - \widehat{a}_n) \phi(\mathbf{x}_n)$$

$$0 = \sum_n (\widehat{a}_n - \widehat{a}_n)$$

$$C = a_n + \mu_n$$

$$C = \widehat{a}_n + \widehat{\mu}_n$$

15/18

## Support Vector Machines for Regression

- Replacing these in the Lagrangian function gives the dual representation of the problem, in which we maximize

$$\frac{1}{2} \sum_n \sum_m (\widehat{a}_n - \widehat{a}_m) k(\mathbf{x}_n, \mathbf{x}_m) - \epsilon \sum_n (\widehat{a}_n + \widehat{a}_n) + \sum_n (\widehat{a}_n - \widehat{a}_n) t_n$$

subject to  $a_n \geq 0$  and  $a_n \leq C$  for all  $n$ , because  $\mu_n \geq 0$ . Similarly for  $\widehat{a}_n$ .

- When the Lagrangian function is maximized, the Karush-Kuhn-Tucker conditions hold for all  $n$ :

$$a_n (y(\mathbf{x}_n) + \epsilon + \xi_n - t_n) = 0$$

$$\widehat{a}_n (t_n - y(\mathbf{x}_n) + \epsilon + \widehat{\xi}_n) = 0$$

$$\mu_n \xi_n = 0$$

$$\widehat{\mu}_n \widehat{\xi}_n = 0$$

- Then,  $a_n > 0$  if and only if  $y(\mathbf{x}_n) + \epsilon + \xi_n - t_n = 0$ , which implies that  $\mathbf{x}_n$  lies on or above the upper margin of the  $\epsilon$ -tube. Similarly for  $\widehat{a}_n > 0$ .

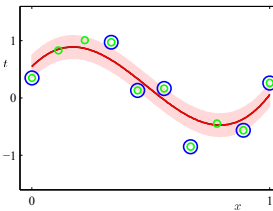
16/18

## Support Vector Machines for Regression

- The prediction for a new point  $\mathbf{x}$  is made according to

$$y(\mathbf{x}) = \sum_{m \in S} (\widehat{a}_m - \widehat{a}_m) k(\mathbf{x}, \mathbf{x}_m) + b$$

where  $S$  are the indexes of the support vectors. Sparse solution!



- To find  $b$ , consider any support vector  $\mathbf{x}_n$  with  $0 < a_n < C$ . Then,  $\mu_n > 0$  and thus  $\xi_n = 0$  and thus  $0 = t_n - \epsilon - y(\mathbf{x}_n)$ . Then,

$$b = t_n - \epsilon - \sum_{m \in S} (\widehat{a}_m - \widehat{a}_m) k(\mathbf{x}_n, \mathbf{x}_m)$$

17/18

## Neural Networks

- Consider binary classification with input space  $\mathbb{R}^D$ . Consider a training set  $\{(\mathbf{x}_n, t_n)\}$  where  $t_n \in \{-1, +1\}$ .

- SVMs classify a new point  $\mathbf{x}$  according to

$$y(\mathbf{x}) = \text{sgn} \left( \sum_{m \in S} a_m t_m k(\mathbf{x}, \mathbf{x}_m) + b \right)$$

- Consider regressing an unidimensional continuous random variable on a  $D$ -dimensional continuous random variable. Consider a training set  $\{(\mathbf{x}_n, t_n)\}$

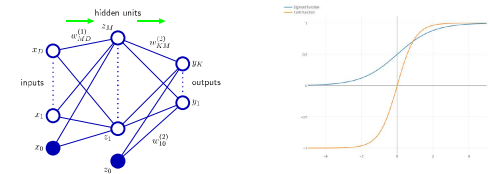
- For a new point  $\mathbf{x}$ , SVMs predict

$$y(\mathbf{x}) = \sum_{m \in S} (\widehat{a}_n - \widehat{a}_n) k(\mathbf{x}, \mathbf{x}_m) + b$$

- SVMs imply **data-selected user-defined** basis functions.
- NNs imply a **user-defined** number of **data-selected** basis functions.

4/16

## Neural Networks



- Activations:  $a_j = \sum_i w_{ji}^{(1)} x_i + w_{j0}^{(1)}$
- Hidden units and activation function:  $z_j = h(a_j)$
- Output activations:  $a_k = \sum_j w_{kj}^{(2)} z_j + w_{k0}^{(2)}$
- Output activation function for regression:  $y_k(\mathbf{x}) = a_k$
- Output activation function for classification:  $y_k(\mathbf{x}) = \sigma(a_k)$
- Sigmoid function:  $\sigma(a) = \frac{1}{1 + \exp(-a)}$
- Two-layer NN:

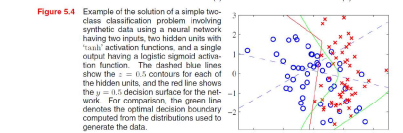
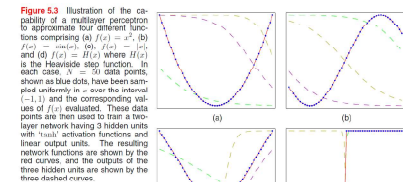
$$y_k(\mathbf{x}) = \sigma \left( \sum_j w_{kj}^{(2)} h \left( \sum_i w_{ji}^{(1)} x_i + w_{j0}^{(1)} \right) + w_{k0}^{(2)} \right)$$

- Evaluating the previous expression is known as forward propagation. The NN is said to have a feed-forward architecture.
- All the previous is, of course, generalizable to more layers.

5/16

## Neural Networks

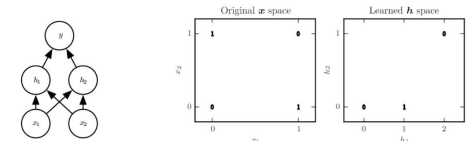
- For a large variety of activation functions, the two-layer NN can uniformly approximate any continuous function to arbitrary accuracy provided enough hidden units. Easy to fit the parameters? Overfitting?!



6/16

## Neural Networks

- Solving the XOR problem with NNs.
- No line shatters the points in the original space.
- The NN represents a mapping of the input space to an alternative space where a line can shatter the points. Note that the points (0,1) and (1,0) are mapped both to the point (1,0).
- It resembles SVMs.



$$\begin{aligned} w_{11}^{(1)} &= w_{12}^{(1)} = w_{21}^{(1)} = w_{22}^{(1)} = 1 \\ w_{10}^{(1)} &= 0, \quad w_{20}^{(1)} = -1 \\ h_j &= z_j = h(a_j) = \max\{0, a_j\} \\ w_{11}^{(2)} &= 1, \quad w_{12}^{(2)} = -2 \\ w_{10}^{(2)} &= 0 \\ y &= y_k = a_k \end{aligned}$$

7/16

## Backpropagation Algorithm

- Consider regressing an  $K$ -dimensional continuous random variable on a  $D$ -dimensional continuous random variable.
- Consider a training set  $\{(\mathbf{x}_n, \mathbf{t}_n)\}$ . Consider minimizing the sum-of-squares error function

$$E(\mathbf{w}) = \sum_n E_n(\mathbf{w}) = \sum_n \frac{1}{2} \|\mathbf{y}(\mathbf{x}_n) - \mathbf{t}_n\|^2 = \sum_n \sum_k \frac{1}{2} (y_k(\mathbf{x}_n) - t_{nk})^2$$

- This error function can be justified from a maximum likelihood approach to learning  $\mathbf{w}$ . To see it, **assume** that

$$p(t_k | \mathbf{x}, \mathbf{w}, \sigma) = \mathcal{N}(t_k | y_k(\mathbf{x}), \sigma)$$

- Then, the likelihood function is

$$p(\{\mathbf{t}_n\} | \{\mathbf{x}_n\}, \mathbf{w}, \sigma) = \prod_n \prod_k \mathcal{N}(t_{nk} | y_k(\mathbf{x}_n), \sigma) = \prod_n \prod_k \frac{1}{(2\pi\sigma^2)^{1/2}} e^{-\frac{1}{2\sigma^2} (t_{nk} - y_k(\mathbf{x}_n))^2}$$

and thus

$$-\ln p(\{\mathbf{t}_n\} | \{\mathbf{x}_n\}, \mathbf{w}, \sigma) = \sum_n \sum_k \frac{1}{2\sigma^2} (t_{nk} - y_k(\mathbf{x}_n))^2 + \frac{N}{2} \ln \sigma^2 + \frac{N}{2} \ln 2\pi$$

which is equivalent to the sum-of-squares error function for a given  $\sigma$ .

- If  $\sigma$  is not given, then we can find the ML estimates of  $\mathbf{w}$ , plug them into the log likelihood function, and maximize it with respect to  $\sigma$ .

9/16

## Backpropagation Algorithm

- The weight space is highly multimodal and, thus, we have to resort to approximate iterative methods to minimize the previous expression.
- Batch gradient descent

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \eta_t \nabla E(\mathbf{w}^t)$$

where  $\eta_t > 0$  is the learning rate ( $\sum_t \eta_t = \infty$  and  $\sum_t \eta_t^2 < \infty$  to ensure convergence, e.g.  $\eta_t = 1/t$ ).

- Sequential, stochastic or online gradient descent

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \eta_t \nabla E_n(\mathbf{w}^t)$$

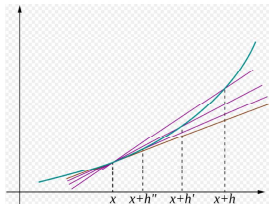
where  $n$  is chosen randomly or sequentially.

- Sequential gradient descent is less affected by the multimodality problem, as a local minimum of the whole data will not be generally a local minimum of each individual point.

9/16

## Backpropagation Algorithm

- Recall that  $f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$



- Recall that  $\nabla E_n(\mathbf{w}^t)$  is a vector whose components are the partial derivatives of  $E_n(\mathbf{w}^t)$ .

10/16

## Backpropagation Algorithm

- Since  $E_n$  depends on  $w_{ji}$  only via  $a_j$ , and  $a_j = \sum_i w_{ji} x_i$ , then

$$\frac{\partial E_n}{\partial w_{ji}} = \frac{\partial E_n}{\partial a_j} \frac{\partial a_j}{\partial w_{ji}} = \frac{\partial E_n}{\partial a_j} x_i = \delta_j x_i$$

- Since  $E_n$  depends on  $a_j$  only via  $a_k$ , then

$$\delta_j = \frac{\partial E_n}{\partial a_j} = \sum_k \frac{\partial E_n}{\partial a_k} \frac{\partial a_k}{\partial a_j} = \sum_k \delta_k \frac{\partial a_k}{\partial a_j}$$

- Since  $a_k = \sum_j w_{kj} z_j$  and  $z_j = h(a_j)$ , then

$$\frac{\partial a_k}{\partial a_j} = h'(a_j) w_{kj}$$

- Putting all together, we have that

$$\delta_j = h'(a_j) \sum_k \delta_k w_{kj}$$

- Since  $y_k = a_k$  for regression and  $a_k = \sum_j w_{kj} z_j$ , then

$$\frac{\partial E_n}{\partial w_{kj}} = \frac{\partial E_n}{\partial a_k} \frac{\partial a_k}{\partial w_{kj}} = \delta_k z_j \text{ and } \delta_k = \frac{\partial E_n}{\partial a_k} = y_k - t_k$$

- Backpropagation algorithm:

- Forward propagate to compute activations, and hidden and output units.
- Compute  $\delta_k$  for the output units.
- Backpropagate the  $\delta$ 's, i.e. evaluate  $\delta_j$  for the hidden units recursively.
- Compute the required derivatives.

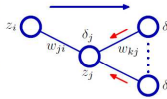
11/16

## Backpropagation Algorithm

- Backpropagation algorithm:

- Forward propagate to compute activations, and hidden and output units.
- Compute  $\delta_k$  for the output units.
- Backpropagate the  $\delta$ 's, i.e. evaluate  $\delta_j$  for the hidden units recursively.
- Compute the required derivatives.

**Figure 5.7** Illustration of the calculation of  $\delta_j$  for hidden unit  $j$  by backpropagation of the  $\delta$ 's from those units  $k$  to which unit  $j$  sends connections. The blue arrow denotes the direction of information flow during forward propagation, and the red arrows indicate the backward propagation of error information.



- For classification, we minimize the negative log likelihood function, a.k.a. cross-entropy error function:

$$E_n(\mathbf{w}) = - \sum_k [t_{nk} \ln y_k(\mathbf{x}_n) + (1 - t_{nk}) \ln(1 - y_k(\mathbf{x}_n))]$$

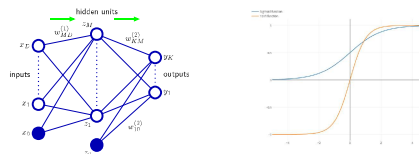
with  $t_{nk} \in \{0, 1\}$  and  $y_k(\mathbf{x}_n) = \sigma(a_k)$ . Then, again

$$\frac{\partial E_n}{\partial w_{kj}} = \delta_k z_j \text{ and } \delta_k = \frac{\partial E_n}{\partial a_k} = y_k - t_k$$

- This is an example of embarrassingly parallel algorithm.

12/16

## Backpropagation Algorithm



- Example:  $y_k = a_k$ , and  $z_j = h(a_j) = \tanh(a_j)$  where  $\tanh(a) = \frac{\exp(a) - \exp(-a)}{\exp(a) + \exp(-a)}$ .

- Note that  $h'(a) = 1 - h(a)^2$ .

- Backpropagation:

- Forward propagation, i.e. compute

$$a_j = \sum_i w_{ji} x_i \text{ and } z_j = h(a_j) \text{ and } y_k = \sum_j w_{kj} z_j$$

- Compute

$$\delta_k = y_k - t_k$$

- Backpropagate, i.e. compute

$$\delta_j = (1 - z_j^2) \sum_k w_{kj} \delta_k$$

- Compute

$$\frac{\partial E_n}{\partial w_{kj}} = \delta_k z_j \text{ and } \frac{\partial E_n}{\partial w_{ji}} = \delta_j x_i$$

13/16

## Kernel Classification

- The moving window rule gives equal weight to all the points in the ball, which may be counterintuitive. Then,

$$y_k(\mathbf{x}) = \begin{cases} 0 & \text{if } \sum_n \mathbf{1}_{\{t_n=1\}} k\left(\frac{\mathbf{x}-\mathbf{x}_n}{h}\right) \leq \sum_n \mathbf{1}_{\{t_n=0\}} k\left(\frac{\mathbf{x}-\mathbf{x}_n}{h}\right) \\ 1 & \text{otherwise} \end{cases}$$

where  $k: \mathbb{R}^D \rightarrow \mathbb{R}$  is a kernel function, which is usually non-negative and monotone decreasing along rays starting from the origin. The parameter  $h$  is called smoothing factor or width.

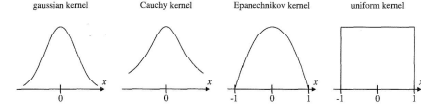
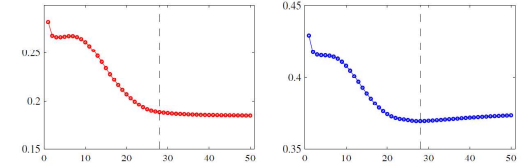


FIGURE 10.3. Various kernels on  $\mathbb{R}$ .

- Gaussian kernel:  $k(u) = \exp(-\|u\|^2)$  where  $\|\cdot\|$  is the Euclidean norm.
- Cauchy kernel:  $k(u) = 1/(1 + \|u\|^{2+1})$
- Epanechnikov kernel:  $k(u) = (1 - \|u\|^2) \mathbf{1}_{\{\|u\| \leq 1\}}$
- Moving window kernel:  $k(u) = \mathbf{1}_{\{u \in S(0,1)\}}$

6/16

## Regularization



**Figure 5.12** An illustration of the behaviour of training set error (left) and validation set error (right) during a typical training session, as a function of the iteration step, for the sinusoidal data set. The goal of achieving the best generalization performance suggests that training should be stopped at the point shown by the vertical dashed lines, corresponding to the minimum of the validation set error.

- Regularization when learning the parameters: Early stopping the backpropagation algorithm according to the error on some validation data.
- Regularization when learning the structure:
  - Cross-validation.
  - Penalizing complexity according to

$$E(\mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|^2 \text{ or } E(\mathbf{w}) + \frac{\lambda_1}{2} \|\mathbf{w}^{(1)}\|^2 + \frac{\lambda_2}{2} \|\mathbf{w}^{(2)}\|^2$$

and choose  $\lambda$ , or  $\lambda_1$  and  $\lambda_2$  by cross-validation. Note that the effect of the penalty is simply to add  $\lambda w_{ji}$  and  $\lambda w_{kj}$ , or  $\lambda_1 w_{ji}$  and  $\lambda_2 w_{kj}$  to the appropriate derivatives.

15/16

## Limitations of Neural Networks

### Theorem (Universal approximation theorem)

For every continuous function  $f: [a, b]^D \rightarrow \mathbb{R}$  and for every  $\epsilon > 0$ , there exists a NN with one hidden layer such that

$$\sup_{\mathbf{x} \in [a, b]^D} |f(\mathbf{x}) - y(\mathbf{x})| < \epsilon$$

### Theorem (Universal classification theorem)

Let  $\mathcal{C}^{(k)}$  contain all classifiers defined by NNs of one hidden layer with  $k$  hidden units and the sigmoid activation function. Then, for any distribution  $p(\mathbf{x}, t)$ ,

$$\lim_{k \rightarrow \infty} \inf_{y \in \mathcal{C}^{(k)}} L(y(\mathbf{x})) - L(p(t|\mathbf{x})) = 0$$

where  $L(\cdot)$  is the 0/1 loss function.

- How many hidden units has such a NN?
- How much data do we need to learn such a NN (and avoid overfitting) via the backpropagation algorithm?
- How fast does the backpropagation algorithm converge to such a NN? Assuming that it does not get trapped in a local minimum...
- The answer to the last two questions depends on the first: More hidden units implies more training time and higher generalization error.

4/16

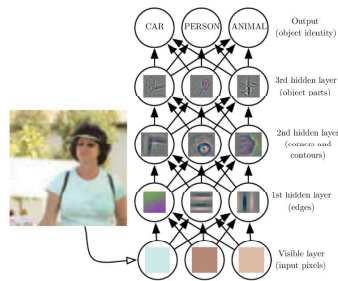
## Limitations of Neural Networks

- ▶ How many hidden units does the NN need ?
- ▶ Any Boolean function can be written in disjunctive normal form (OR of ANDs) or conjunctive normal form (AND of ORs). This is a depth-two logical circuit.
- ▶ For most Boolean functions, the size of the circuit is exponential in the size of the input.
- ▶ However, there are Boolean functions that have a polynomial-size circuit of depth  $k$  and an exponential-size circuit of depth  $k-1$ .
- ▶ Then, there is no universally right depth. Ideally, we should let the data determine the right depth.

### Theorem (No free lunch theorem)

For any algorithm, good performance on some problems comes at the expense of bad performance on some others.

## Deep Neural Networks



- ▶ A deep NN is a function that maps input to output.
- ▶ The mapping is formed by composing many simpler functions.
- ▶ Each layer provides a new representation of the input, i.e. complex concepts are built from simpler ones.
- ▶ The representation is learned automatically from data.

## Deep Neural Networks

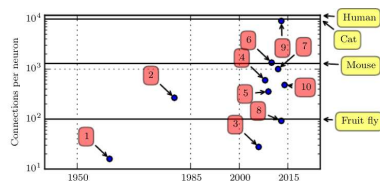


Figure 1.10: Initially, the number of connections between neurons in artificial neural networks was limited by hardware capabilities. Today, the number of connections between neurons is mostly a design consideration. Some artificial neural networks have nearly as many connections per neuron as a cat, and it is quite common for other neural networks to have as many connections per neuron as smaller mammals like mice. Even the human brain does not have an exorbitant amount of connections per neuron. Biological neural network sizes from Wikipedia (2015).

1. Adaptive linear element (Widrow and Hoff, 1960)
2. Neocognitron (Fukushima, 1980)
3. GPU-accelerated convolutional network (Chellapilla et al., 2006)
4. Deep Boltzmann machine (Salakhutdinov and Hinton, 2009a)
5. Unsupervised convolutional network (Jarrett et al., 2009)
6. GPU-accelerated multilayer perceptron (Ciresan et al., 2010)
7. Distributed autoencoder (Lee et al., 2012)
8. Multi-GPU convolutional network (Krizhevsky et al., 2013)
9. COTS HPC unsupervised convolutional network (Coston et al., 2013)
10. GoogLeNet (Szegedy et al., 2014a)

5/16

## Deep Neural Networks

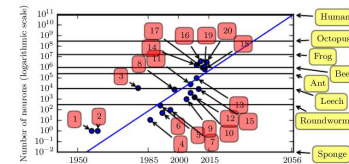


Figure 1.11: Since the introduction of hidden units, artificial neural networks have doubled in size roughly every 2.4 years. Biological neural network sizes from Wikipedia (2015).

1. Perceptron (Rosenblatt, 1958, 1962)
2. Adaptive linear element (Widrow and Hoff, 1960)
3. Neocognitron (Fukushima, 1980)
4. Early back-propagation network (Rumelhart et al., 1986)
5. Recurrent neural network for speech recognition (Bollmann and Fahlke, 1991)
6. Multilayer perceptron for speech recognition (Bengio et al., 1991)
7. Mean field sigmoid belief network (Neal et al., 1996)
8. LeNet-5 (LeCun et al., 1998)
9. Echo state network (Lange and Haas, 2004)
10. Deep belief network (Salakhutdinov et al., 2006)
11. GPU-accelerated convolutional network (Chellapilla et al., 2006)
12. Deep Boltzmann machine (Salakhutdinov and Hinton, 2009a)
13. GPU-accelerated deep belief network (Storn et al., 2010)
14. Unsupervised convolutional network (Jarrett et al., 2009)
15. GPU-accelerated multilayer perceptron (Ciresan et al., 2010)
16. DMF-1 network (Cristian and Ng, 2011)
17. Distributed autoencoder (Lee et al., 2012)
18. Multi-GPU convolutional network (Krizhevsky et al., 2013)
19. COTS HPC unsupervised convolutional network (Coston et al., 2013)
20. GoogLeNet (Szegedy et al., 2014a)

22 layers DNN, but  
12 times fewer weights  
than DNN 19

6/16

## Deep Neural Networks

- ▶ Training DNNs is difficult:
  - ▶ Typically, poorer generalization than (shallow) NNs.
  - ▶ The gradient may vanish/explode as we move away from the output layer, due to multiplying small/big quantities. E.g. the gradient of  $\sigma$  and  $\tanh$  is in  $[0, 1]$ . So, they may only suffer the gradient vanishing problem. Other activations functions may suffer the gradient exploding problem.
  - ▶ There may be larger plateaus and many more local minima than with NNs.
- ▶ Training DNNs is doable:
  - ▶ Convolutional networks, particularly suitable for image processing.
  - ▶ Rectifier activation function, a new activation function.
  - ▶ Layer-wise pre-training, to find a good starting point for training.
- ▶ In addition to performance, the computational demands of the training must be considered, e.g. CPU, GPU, memory, parallelism, etc.
  - ▶ The authors state that GoogLeNet was trained "using modest amount of model and data-parallelism. Although we used a CPU based implementation only, a rough estimate suggests that the GoogLeNet network could be trained to convergence using few high-end GPUs within a week, the main limitation being the memory usage".

5/16

6/16

7/16