

Neural networks

Logistic regression

Suppose $\mathcal{X} = \mathbb{R}^d$ and $\mathcal{Y} = \{0, 1\}$.

Logistic regression: statistical model for $Y \mid \mathbf{X} = \mathbf{x}$ for each $\mathbf{x} \in \mathbb{R}^d$:

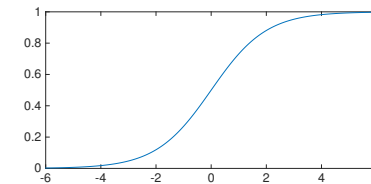
$$\mathcal{P} = \left\{ P_{\mathbf{w}} : \mathbf{w} \in \mathbb{R}^d \right\},$$

where

$$f_{\mathbf{w}}(\mathbf{x}) := P_{\mathbf{w}}(Y = 1 \mid \mathbf{X} = \mathbf{x}) = g(\langle \mathbf{w}, \mathbf{x} \rangle)$$

and

$$g(z) := \text{logistic}(z) = \frac{1}{1 + e^{-z}}.$$



(Bias term omitted for simplicity.)

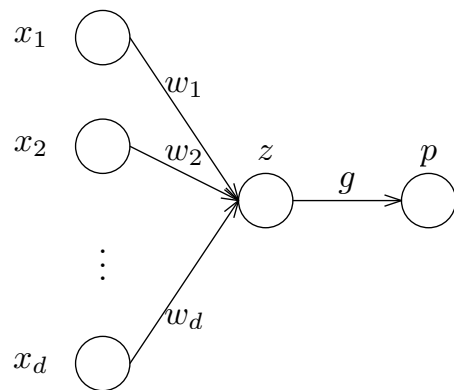
1 / 24

2 / 24

Logistic regression

Learning \mathbf{w} from data

Network diagram for $f_{\mathbf{w}}$:



$$z := \sum_{j=1}^d w_j x_j, \quad p := g(z) \quad (g = \text{logistic}).$$

Recall $f_{\mathbf{w}}(\mathbf{x}) = g(\langle \mathbf{w}, \mathbf{x} \rangle)$ where $g = \text{logistic}$ (called the *link function*).

Using log-loss function

$$\ell_{\log}(y, p) := y \log \frac{1}{p} + (1 - y) \log \frac{1}{1 - p}$$

in ERM framework with training data $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$:

$$\min_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1}^n \ell_{\log}(y_i, g(\langle \mathbf{w}, \mathbf{x}_i \rangle)).$$

3 / 24

4 / 24

Stochastic gradient method

Given: training data $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$.

- ▶ Start with some initial $\mathbf{w}^{(1)} \in \mathbb{R}^d$.
- ▶ For $t = 1, 2, \dots$ until some stopping condition is satisfied.
 - ▶ Pick $(\mathbf{x}, y) \in \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ uniformly at random.
 - ▶ Compute

$$\boldsymbol{\lambda}^{(t)} := \nabla_{\mathbf{w}} \ell_{\log}(y, g(\langle \mathbf{w}, \mathbf{x} \rangle)) \Big|_{\mathbf{w}=\mathbf{w}^{(t)}}.$$

- ▶ Update:

$$\mathbf{w}^{(t+1)} := \mathbf{w}^{(t)} - \eta_t \boldsymbol{\lambda}^{(t)}.$$

How to compute $\boldsymbol{\lambda}^{(t)}$?

5 / 24

Using the chain rule

Usual *chain rule* from calculus gives generic way to compute gradient of log-loss (w.r.t. \mathbf{w}) on example (\mathbf{x}, y) :

$$\nabla_{\mathbf{w}} \{ \ell_{\log}(y, f_{\mathbf{w}}(\mathbf{x})) \} = \frac{\partial}{\partial p} \{ \ell_{\log}(y, p) \} \Big|_{p=f_{\mathbf{w}}(\mathbf{x})} \cdot \nabla_{\mathbf{w}} \{ f_{\mathbf{w}}(\mathbf{x}) \}.$$

- ▶ First term:

$$\frac{\partial}{\partial p} \{ \ell_{\log}(y, p) \} = -\frac{y}{p} + \frac{1-y}{1-p},$$

and plug-in $f_{\mathbf{w}}(\mathbf{x})$ for p .

- ▶ Second term: since $f_{\mathbf{w}}(\mathbf{x}) = \text{logistic}(\langle \mathbf{w}, \mathbf{x} \rangle)$, applying chain rule again gives

$$\begin{aligned} \nabla_{\mathbf{w}} \{ f_{\mathbf{w}}(\mathbf{x}) \} &= \frac{\partial}{\partial z} \text{logistic}(z) \Big|_{z=\langle \mathbf{w}, \mathbf{x} \rangle} \cdot \mathbf{x} \\ &= \text{logistic}(\langle \mathbf{w}, \mathbf{x} \rangle) \cdot (1 - \text{logistic}(\langle \mathbf{w}, \mathbf{x} \rangle)) \cdot \mathbf{x}. \end{aligned}$$

6 / 24

Extensions

- ▶ Using log-loss and logistic link function, objective

$$\min_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1}^n \ell_{\log}(y_i, \text{logistic}(\langle \mathbf{w}, \mathbf{x}_i \rangle))$$

is same as negative log-likelihood for logistic regression.

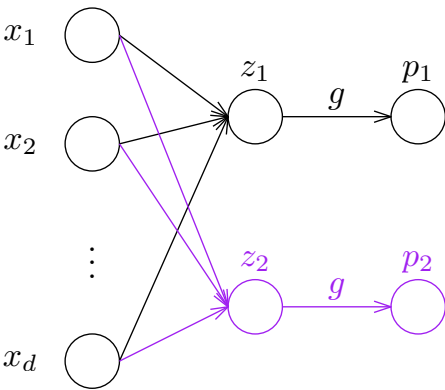
- ▶ Can also use other loss functions $\ell(y, p)$ and other link functions $g(z)$.
E.g., $\ell(y, p) = (y - p)^2$, $g(z) = \tanh(z)$.

Note: objective not necessarily convex, even if $p \rightarrow \ell(y, p)$ is convex.

Nevertheless, stochastic gradient method may still be effective.

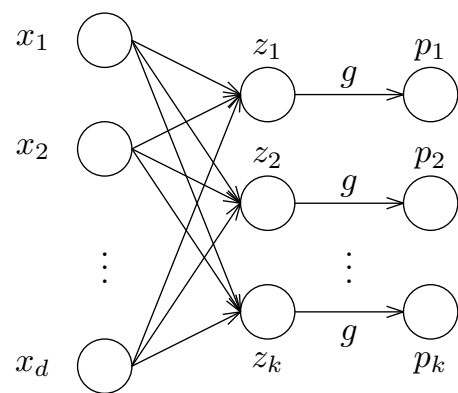
7 / 24

Two-output network



$$z_i := \sum_{j=1}^d W_{i,j} x_j, \quad y_i := g(z_i), \quad i \in \{1, 2\}.$$

8 / 24



$$z_i := \sum_{j=1}^d W_{i,j} x_j, \quad y_i := g(z_i), \quad i \in \{1, 2, \dots, k\}.$$

- Suppose we have two functions

$$\begin{aligned} f_{\mathbf{W}^{(1)}} : \mathbb{R}^d &\rightarrow \mathbb{R}^k, & (\mathbf{W}^{(1)} &\in \mathbb{R}^{k \times d}), \\ f_{\mathbf{W}^{(2)}} : \mathbb{R}^k &\rightarrow \mathbb{R}^\ell, & (\mathbf{W}^{(2)} &\in \mathbb{R}^{\ell \times k}), \end{aligned}$$

where

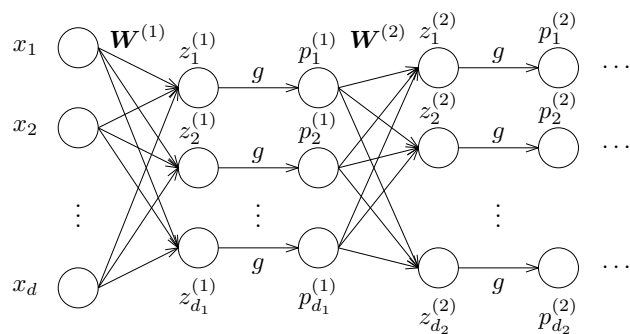
$$f_{\mathbf{W}}(\mathbf{x}) := G(\mathbf{W}\mathbf{x}),$$

and G applies the link function g coordinate-wise to a vector.

- Composition:** $f_{\mathbf{W}^{(1)}, \mathbf{W}^{(2)}} := f_{\mathbf{W}^{(2)}} \circ f_{\mathbf{W}^{(1)}}$ is defined by

$$f_{\mathbf{W}^{(1)}, \mathbf{W}^{(2)}}(\mathbf{x}) := f_{\mathbf{W}^{(2)}}(f_{\mathbf{W}^{(1)}}(\mathbf{x})).$$

This is a *two-layer neural network*.



- A one-layer neural network with a monotonic link function is (essentially) a *linear classifier*.

Cannot represent XOR function (Minsky and Papert, 1969).

- Any continuous function f can be approximated arbitrarily well by a two-layer neural network

$$f \approx \underbrace{f_{\mathbf{W}^{(1)}}}_{\mathbb{R}^k \rightarrow \mathbb{R}} \circ \underbrace{f_{\mathbf{W}^{(0)}}}_{\mathbb{R}^d \rightarrow \mathbb{R}^k}$$

(Cybenko, 1989; Barron, 1993).

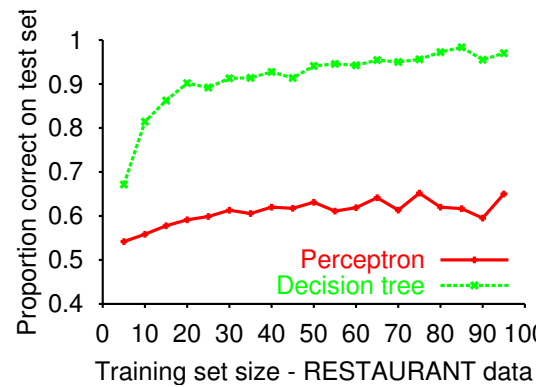
Caveat: may need a very large number of nodes.

- Some functions can be approximated with *exponentially fewer nodes* by using more than two layers (Eldan and Shamir, 2015; Telgarsky, 2015).

- $(z_i^{(l)}, p_i^{(l)})$ regarded as a single “unit”.
(Typically represented by a single node.)
- All except *input* (\mathbf{x}) and *output* units are considered “hidden”.

Limits of single-layer network

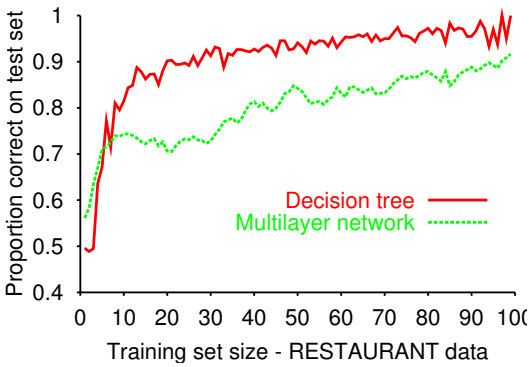
Linear classifier unable to represent complex functions.



(Figure from Stuart Russell aima.eecs.berkeley.edu/slides-pdf/chapter20b.pdf)

Benefits of depth

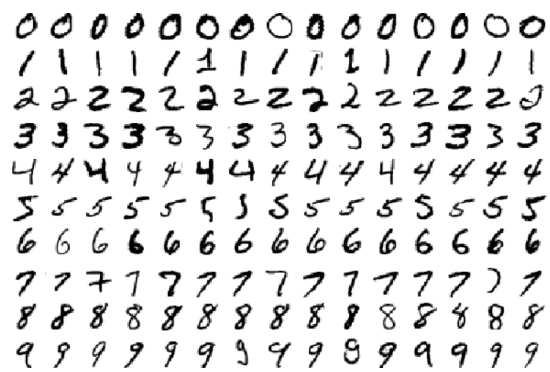
Neural network with four hidden units.



(Figure from Stuart Russell aima.eecs.berkeley.edu/slides-pdf/chapter20b.pdf)

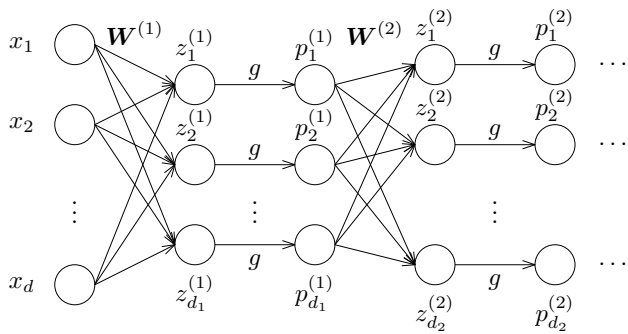
MNIST

OCR digits classification task



- ▶ 3-nearest-neighbor: 3.0% error rate
- ▶ Neural network with 700 hidden units: 1.6% error rate
- ▶ Current best convolutional network: 0.23% error rate

H-layer neural network



- ▶ Index layers by $l = 0, 1, 2, \dots, H$; d_l “units” in layer l .
- ▶ $p_j^{(l)}$ is output of unit j in layer l . ($\mathbf{p}^{(0)} = \mathbf{x} \in \mathbb{R}^{d_0}$ is the *input layer*.)
- ▶ $w_{i,j}^{(l)}$ is weight of edge from unit j in layer $l - 1$ to unit i in layer l :

$$z_i^{(l)} = \langle \mathbf{w}_i^{(l)}, \mathbf{p}^{(l-1)} \rangle = \sum_{j=1}^{d_{l-1}} w_{i,j}^{(l)} \cdot p_j^{(l-1)}.$$

- ▶ (Transpose of) $\mathbf{w}_i^{(l)}$ is i -th row of $\mathbf{W}^{(l)}$.

Evaluating a neural network

H -layer neural network f with weights $\{\mathbf{w}_i^{(l)} : 0 \leq l \leq H-1, 1 \leq i \leq d_l\}$ where $\mathbf{w}_i^{(l)} \in \mathbb{R}^{d_{l-1}}$.

Forward propagation to compute $f(\mathbf{x})$.

- ▶ Start with $\mathbf{p}^{(0)} := \mathbf{x}$.
- ▶ For $l = 1, 2, \dots, H$:
 - ▶ $z_i^{(l)} := \langle \mathbf{w}_i^{(l)}, \mathbf{p}^{(l-1)} \rangle$ for each $i = 1, 2, \dots, d_l$.
 - ▶ $p_i^{(l)} := g(z_i^{(l)})$ for each $i = 1, 2, \dots, d_l$.
- ▶ Return $\mathbf{p}^{(H)}$.

17 / 24

Training a neural network

How to compute gradient of loss $\ell(y, f(\mathbf{x}))$ w.r.t weights?

Direct calculation of gradients w.r.t. $\mathbf{w}_i^{(l)}$ is complicated.

Mid-to-late 20th century researchers discovered how to use *partial derivatives* to manage gradient computation: **backpropagation algorithm**.

Another view of backpropagation: introduce auxiliary variables $\mathbf{p}^{(l)}$ and constraints

$$\mathbf{p}^{(l)} = f_{\mathbf{w}^{(l)}}(\mathbf{p}^{(l-1)})$$

for $1 \leq l \leq H$, and apply “method of adjoints”.

18 / 24

Recipe for method of adjoints

- ▶ Introduce auxiliary variables $\mathbf{p}^{(l)}$ for $1 \leq l \leq H$, and constraints

$$\mathbf{p}^{(l)} = f_{\mathbf{w}^{(l)}}(\mathbf{p}^{(l-1)}), \quad \text{for } 1 \leq l \leq H.$$

(Also define $\mathbf{p}^{(0)} := \mathbf{x}$; not variable.)

- ▶ Introduce Lagrange variables $\alpha^{(l)} \in \mathbb{R}^{d_l}$ for $1 \leq l \leq H$, and Lagrangian

$$\mathcal{L} := \ell(y, f(\mathbf{x})) - \sum_{l=1}^H \left(\mathbf{p}^{(l)} - f_{\mathbf{w}^{(l)}}(\mathbf{p}^{(l-1)}) \right)^\top \alpha^{(l)}.$$

- ▶ Compute values of auxiliary variables and Lagrange variables that satisfy

$$\nabla_{\alpha^{(l)}} \mathcal{L} = \mathbf{0}, \quad \nabla_{\mathbf{p}^{(l)}} \mathcal{L} = \mathbf{0}, \quad \text{for } 1 \leq l \leq H.$$

- ▶ Compute gradient of \mathcal{L} w.r.t. $\mathbf{w}_i^{(l)}$ in terms of auxiliary and Lagrange variables; much simpler!

(This method is justified by some very beautiful optimization theory.)

19 / 24

Setting the auxiliary and Lagrange variables

$$\text{Lagrangian: } \mathcal{L} := \ell(y, f(\mathbf{x})) - \sum_{l=1}^H \left(\mathbf{p}^{(l)} - f_{\mathbf{w}^{(l)}}(\mathbf{p}^{(l-1)}) \right)^\top \alpha^{(l)}.$$

- ▶ For $1 \leq l \leq H$: $\nabla_{\alpha^{(l)}} \mathcal{L} = \mathbf{p}^{(l)} - f_{\mathbf{w}^{(l)}}(\mathbf{p}^{(l-1)})$.

This is zero if we set $\mathbf{p}^{(l)} := f_{\mathbf{w}^{(l)}}(\mathbf{p}^{(l-1)})$ (i.e., forward propagation).

- ▶ For $l = H$: $\nabla_{\mathbf{p}^{(H)}} \mathcal{L} = \nabla_{\mathbf{p}^{(H)}} \left\{ \ell(y, \mathbf{p}^{(H)}) \right\} - \alpha^{(H)}$.

This is zero if we set $\alpha^{(H)} := \nabla_{\mathbf{p}^{(H)}} \left\{ \ell(y, \mathbf{p}^{(H)}) \right\}$.

- ▶ For $1 \leq l < H$: $\nabla_{\mathbf{p}^{(l)}} \mathcal{L} = \nabla_{\mathbf{p}^{(l)}} \left\{ f_{\mathbf{w}^{(l+1)}}(\mathbf{p}^{(l)})^\top \alpha^{(l+1)} \right\} - \alpha^{(l)}$.

This is zero if we set

$$\alpha^{(l)} := \nabla_{\mathbf{p}^{(l)}} \left\{ f_{\mathbf{w}^{(l+1)}}(\mathbf{p}^{(l)})^\top \alpha^{(l+1)} \right\} = \sum_{i=1}^{d_{l+1}} \alpha_i^{(l+1)} \cdot g'(\langle \mathbf{w}_i^{(l+1)}, \mathbf{p}^{(l)} \rangle) \cdot \mathbf{w}_i^{(l+1)}.$$

- ▶ For $1 \leq l \leq H$: $\nabla_{\mathbf{w}_i^{(l)}} \mathcal{L} = \alpha_i^{(l)} \cdot g'(\langle \mathbf{w}_i^{(l)}, \mathbf{p}^{(l-1)} \rangle) \cdot \mathbf{p}^{(l-1)}$ (simple!).

20 / 24

Backward propagation

Backward propagation (“backprop”) to compute gradient of $\ell(y, f(\mathbf{x}))$ with respect to weights.

Input: labeled example (\mathbf{x}, y) , current weights $\mathbf{W}^{(l)}$ for $1 \leq l \leq H$.

- ▶ Run forward propagation, saving intermediate computations:
 $\mathbf{p}^{(l)} = f_{\mathbf{W}^{(l)}}(\mathbf{p}^{(l-1)})$ and $z_i^{(l)} = \langle \mathbf{w}_i^{(l)}, \mathbf{p}^{(l-1)} \rangle$ for $1 \leq l \leq H$.
- ▶ Set $\boldsymbol{\alpha}^{(H)} := \nabla_{\mathbf{p}^{(H)}} \ell(y, \mathbf{p}^{(H)})$.
- ▶ For $l = H, H-1, \dots, 1$:
 - ▶ Set $\nabla_{\mathbf{w}_i^{(l)}} \ell(y, f(\mathbf{x})) := \alpha_i^{(l)} \cdot g'(z_i^{(l)}) \cdot \mathbf{p}^{(l-1)}$ for $1 \leq i \leq d_l$.
 - ▶ Set $\boldsymbol{\alpha}^{(l-1)} := \sum_{i=1}^{d_l} \alpha_i^{(l)} \cdot g'(z_i^{(l)}) \cdot \mathbf{w}_i^{(l)}$ (unless $l = 1$).
- ▶ Return gradients $\nabla_{\mathbf{w}_i^{(l)}} \ell(y, f(\mathbf{x}))$ for $1 \leq l \leq H, 1 \leq i \leq d_l$.

21 / 24

Practical tips

- ▶ Apply stochastic gradient method to examples in random order.
- ▶ Standardize inputs (i.e., center and divide by standard deviation).
- ▶ (More expensive): completely de-correlate inputs (i.e., PCA).
- ▶ Can also use other loss functions, e.g., $\ell(y, p) = (y - p)^2$ (square loss).
- ▶ Link function: prefer “centered” range, e.g., \tanh (as opposed to logistic).
- ▶ Initialization: random? But take care so that weights are not too large or too small.
E.g., for node with d inputs, draw weights iid from $N(0, 1/d)$.

22 / 24

Wrap-up

- ▶ Frontier of experimental machine learning research.
- ▶ Basic ideas same as from the 1980s.
- ▶ What's new?
 - ▶ Designing complex “architectures” that work well for specific problems.
 - ▶ Fast algorithms / hardware for optimizing neural network weights on large data sets (e.g., using stream-based processing, distributed systems).
 - ▶ Applications: visual detection and recognition, speech recognition, general function fitting (e.g., learning “reward” functions of different actions of video games), etc.
 - ▶ ...

23 / 24

Key takeaways

1. Structure of simple multilayer neural networks; concept of link functions.
2. Optimization problems with neural networks; use of chain rule to compute gradients.
3. Arguments for depth.
4. High-level idea of backprop.

24 / 24