

# Multiclass Logistic Regression, Multilayer Perceptron

Jindřich Libovický (reusing materials by Milan Straka)

 October 24, 2023



Charles University in Prague  
Faculty of Mathematics and Physics  
Institute of Formal and Applied Linguistics



unless otherwise stated

- Implement **multiclass classification** with softmax.
- Reason about linear regression, logistic regression and softmax classification in a **single probabilistic framework**: with different target distributions, activation functions and training using maximum likelihood estimate.
- Explain **multi-layer perceptron** as a further generalization of linear models.

# Logistic Regression

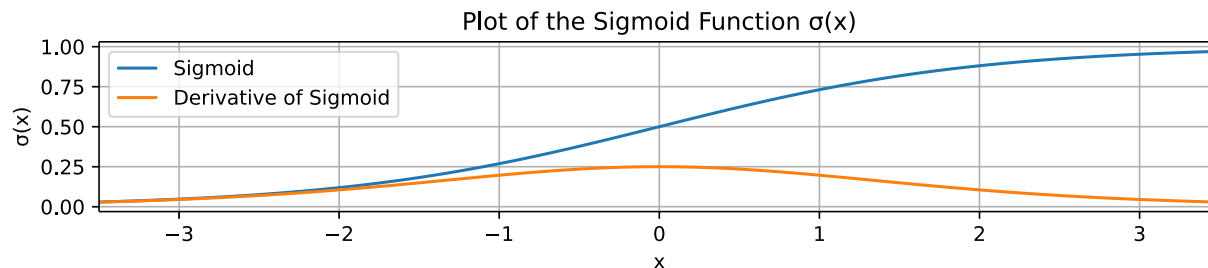
An extension of perceptron, which models the conditional probabilities of  $p(C_0|\mathbf{x})$  and of  $p(C_1|\mathbf{x})$ . (It can in fact handle also more than two classes, which we will see shortly.)

Logistic regression employs the following parametrization of the conditional class probabilities:

$$p(C_1|\mathbf{x}) = \sigma(\mathbf{x}^T \mathbf{w} + b)$$
$$p(C_0|\mathbf{x}) = 1 - p(C_1|\mathbf{x}),$$

where  $\sigma$  is the **sigmoid function**

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



It can be trained using the SGD algorithm.

We denote the output of the “linear part” of the logistic regression as  $\bar{y}(\mathbf{x}; \mathbf{w}) = \mathbf{x}^T \mathbf{w}$  and the overall prediction as  $y(\mathbf{x}; \mathbf{w}) = \sigma(\bar{y}(\mathbf{x}; \mathbf{w})) = \sigma(\mathbf{x}^T \mathbf{w})$ .

The logistic regression output  $y(\mathbf{x}; \mathbf{w})$  models the probability of class  $C_1$ ,  $p(C_1|\mathbf{x})$ .

To give some meaning to the output of the linear part  $\bar{y}(\mathbf{x}; \mathbf{w})$ , starting with

$$p(C_1|\mathbf{x}) = \sigma(\bar{y}(\mathbf{x}; \mathbf{w})) = \frac{1}{1 + e^{-\bar{y}(\mathbf{x}; \mathbf{w})}},$$

we arrive at

$$\bar{y}(\mathbf{x}; \mathbf{w}) = \log \left( \frac{p(C_1|\mathbf{x})}{1 - p(C_1|\mathbf{x})} \right) = \log \left( \frac{p(C_1|\mathbf{x})}{p(C_0|\mathbf{x})} \right),$$

which is called a **logit** and it is a logarithm of odds of the probabilities of the two classes.

# Logistic Regression

To train the logistic regression, we use MLE (the maximum likelihood estimation). Its application is straightforward, given that  $p(C_1|\mathbf{x}; \mathbf{w})$  is directly the model output  $y(\mathbf{x}; \mathbf{w})$ .

Therefore, the loss for a minibatch  $\mathbb{X} = \{(\mathbf{x}_1, t_1), (\mathbf{x}_2, t_2), \dots, (\mathbf{x}_N, t_N)\}$  is

$$E(\mathbf{w}) = \frac{1}{N} \sum_i -\log(p(C_{t_i}|\mathbf{x}_i; \mathbf{w})).$$

**Input:** Input dataset  $(\mathbf{X} \in \mathbb{R}^{N \times D}, \mathbf{t} \in \{0, +1\}^N)$ , learning rate  $\alpha \in \mathbb{R}^+$ .

- $\mathbf{w} \leftarrow \mathbf{0}$  or we initialize  $\mathbf{w}$  randomly
- until convergence (or patience runs out), process a minibatch of examples  $\mathbb{B}$ :
  - $\mathbf{g} \leftarrow \frac{1}{|\mathbb{B}|} \sum_{i \in \mathbb{B}} \nabla_{\mathbf{w}} \left( -\log(p(C_{t_i}|\mathbf{x}_i; \mathbf{w})) \right)$
  - $\mathbf{w} \leftarrow \mathbf{w} - \alpha \mathbf{g}$

# Generalized Linear Models

The logistic regression is in fact an extended linear regression. A linear regression model, which is followed by an **activation function**  $a$ , is called **generalized linear model**:

$$p(t|\mathbf{x}; \mathbf{w}, b) = y(\mathbf{x}; \mathbf{w}, b) = a(\bar{y}(\mathbf{x}; \mathbf{w}, b)) = a(\mathbf{x}^T \mathbf{w} + b).$$

| Name                | Activation        | Distribution | Loss  | Gradient                        |
|---------------------|-------------------|--------------|---|---------------------------------|
| linear regression   | identity          | ?            | $\text{MSE} \propto \mathbb{E}(y(\mathbf{x}) - t)^2$    | $(y(\mathbf{x}) - t)\mathbf{x}$ |
| logistic regression | $\sigma(\bar{y})$ | Bernoulli    | $\text{NLL} \propto \mathbb{E} - \log(p(t \mathbf{x}))$ | ?                               |

# Logistic Regression Gradient

We start by computing the gradient of the  $\sigma(x)$ .

$$\begin{aligned}
 \frac{\partial}{\partial x} \sigma(x) &= \frac{\partial}{\partial x} \frac{1}{1 + e^{-x}} \\
 &= \frac{\frac{\partial}{\partial x} (1 + e^{-x})}{(1 + e^{-x})^2} \\
 &= \frac{1}{1 + e^{-x}} \cdot \frac{e^{-x}}{1 + e^{-x}} \\
 &= \sigma(x) \cdot \frac{e^{-x} + 1 - 1}{1 + e^{-x}} \\
 &= \sigma(x) \cdot (1 - \sigma(x))
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial}{\partial x} \frac{1}{g(x)} &= -\frac{\frac{\partial}{\partial x} g(x)}{g(x)^2} \\
 \frac{\partial}{\partial x} e^{g(x)} &= e^{g(x)} \cdot \frac{\partial}{\partial x} g(x)
 \end{aligned}$$

# Logistic Regression Gradient

Consider the log-likelihood of logistic regression  $\log p(t|\mathbf{x}; \mathbf{w})$ . For brevity, we denote  $\bar{y}(\mathbf{x}; \mathbf{w}) = \mathbf{x}^T \mathbf{w}$  just as  $\bar{y}$  in the following computation.

Remembering that for  $t \sim \text{Ber}(\varphi)$  we have  $p(t) = \varphi^t (1 - \varphi)^{1-t}$ , we can rewrite the log-likelihood to:

$$\begin{aligned}\log p(t|\mathbf{x}; \mathbf{w}) &= \log \sigma(\bar{y})^t (1 - \sigma(\bar{y}))^{1-t} \\ &= t \cdot \log(\sigma(\bar{y})) + (1 - t) \cdot \log(1 - \sigma(\bar{y}))\end{aligned}$$



# Logistic Regression Gradient

$$\begin{aligned}
 \nabla_{\mathbf{w}} - \log p(t|\mathbf{x}; \mathbf{w}) &= \\
 &= \nabla_{\mathbf{w}} \left( -t \cdot \log(\sigma(\bar{y})) - (1-t) \cdot \log(1 - \sigma(\bar{y})) \right) \\
 & \qquad \qquad \qquad \frac{\partial}{\partial x} \log g(x) = \frac{1}{g(x)} \cdot \frac{\partial}{\partial x} g(x) \\
 &= -t \cdot \frac{1}{\sigma(\bar{y})} \cdot \nabla_{\mathbf{w}} \sigma(\bar{y}) - (1-t) \cdot \frac{1}{1 - \sigma(\bar{y})} \cdot \nabla_{\mathbf{w}} (1 - \sigma(\bar{y})) \\
 & \qquad \qquad \qquad \frac{\partial}{\partial x} f(g(x)) = \frac{\partial}{\partial g(x)} f(g(x)) \cdot \frac{\partial}{\partial x} g(x) = \frac{\partial}{\partial z} f(z) \cdot \frac{\partial}{\partial x} g(x) \\
 & \qquad \qquad \qquad \nabla_{\mathbf{w}} \sigma(\bar{y}) = \frac{\partial}{\partial \bar{y}} \sigma(\bar{y}) \cdot \nabla_{\mathbf{w}} \bar{y} \\
 &= -t \cdot \frac{1}{\sigma(\bar{y})} \cdot \sigma(\bar{y}) \cdot (1 - \sigma(\bar{y})) \cdot \nabla_{\mathbf{w}} \bar{y} + (1-t) \cdot \frac{1}{1 - \sigma(\bar{y})} \cdot \sigma(\bar{y}) \cdot (1 - \sigma(\bar{y})) \cdot \nabla_{\mathbf{w}} \bar{y} \\
 &= (-t + t\sigma(\bar{y}) + \sigma(\bar{y}) - t\sigma(\bar{y})) \mathbf{x} \\
 &= (y(\mathbf{x}; \mathbf{w}) - t) \mathbf{x}
 \end{aligned}$$

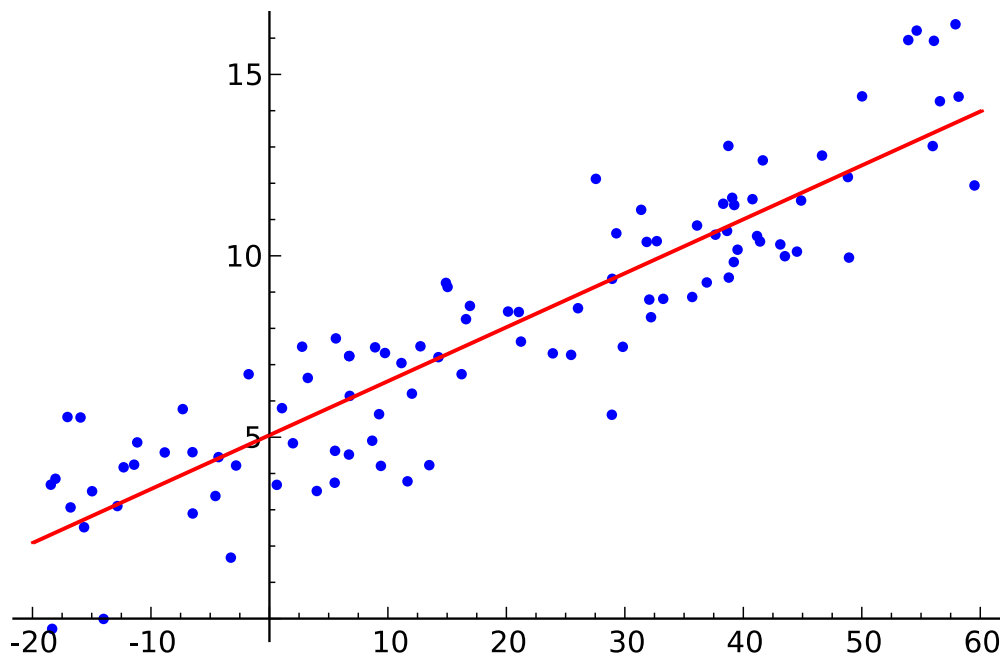
# Generalized Linear Models

The logistic regression is in fact an extended linear regression. A linear regression model, which is followed by some **activation function**  $a$ , is called **generalized linear model**:

$$p(t|\mathbf{x}; \mathbf{w}, b) = y(\mathbf{x}; \mathbf{w}, b) = a(\bar{y}(\mathbf{x}; \mathbf{w}, b)) = a(\mathbf{x}^T \mathbf{w} + b).$$

| Name                | Activation        | Distribution | Loss  | Gradient                        |
|---------------------|-------------------|--------------|---|---------------------------------|
| linear regression   | identity          | ?            | $\text{MSE} \propto \mathbb{E}(y(\mathbf{x}) - t)^2$    | $(y(\mathbf{x}) - t)\mathbf{x}$ |
| logistic regression | $\sigma(\bar{y})$ | Bernoulli    | $\text{NLL} \propto \mathbb{E} - \log(p(t \mathbf{x}))$ | $(y(\mathbf{x}) - t)\mathbf{x}$ |

# Mean Square Error as MLE



[https://upload.wikimedia.org/wikipedia/commons/3/3a/Linear\\_regression.svg](https://upload.wikimedia.org/wikipedia/commons/3/3a/Linear_regression.svg)

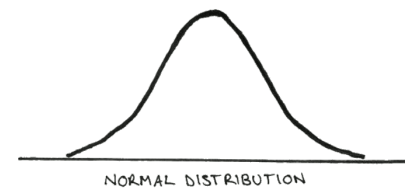
During regression, we predict a number, not a probability distribution. To generate a distribution, we might consider a distribution with the mean of the predicted value and a fixed variance  $\sigma^2$  – the most general such a distribution is the normal distribution.

# Mean Square Error as MLE

Therefore, assume our model generates a distribution  $p(t|\mathbf{x}; \mathbf{w}) = \mathcal{N}(t; y(\mathbf{x}; \mathbf{w}), \sigma^2)$ .

Now we can apply the maximum likelihood estimation and get

$$\begin{aligned}
 \arg \max_{\mathbf{w}} p(\mathbf{t}|\mathbf{X}; \mathbf{w}) &= \arg \min_{\mathbf{w}} \sum_{i=1}^N -\log p(t_i|\mathbf{x}_i; \mathbf{w}) \\
 &= \arg \min_{\mathbf{w}} - \sum_{i=1}^N \log \sqrt{\frac{1}{2\pi\sigma^2}} e^{-\frac{(t_i - y(\mathbf{x}_i; \mathbf{w}))^2}{2\sigma^2}} \\
 &= \arg \min_{\mathbf{w}} -N \log(2\pi\sigma^2)^{-1/2} - \sum_{i=1}^N -\frac{(t_i - y(\mathbf{x}_i; \mathbf{w}))^2}{2\sigma^2} \\
 &= \arg \min_{\mathbf{w}} \sum_{i=1}^N \frac{(t_i - y(\mathbf{x}_i; \mathbf{w}))^2}{2\sigma^2} = \arg \min_{\mathbf{w}} \frac{1}{N} \sum_{i=1}^N (y(\mathbf{x}_i; \mathbf{w}) - t_i)^2.
 \end{aligned}$$



*Freeman*  
<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2465539/>

We have therefore extended the GLM table to

| Name                | Activation        | Distribution | Loss  | Gradient                        |
|---------------------|-------------------|--------------|---|---------------------------------|
| linear regression   | identity          | Normal       | $\text{NLL} \propto \text{MSE}$                         | $(y(\mathbf{x}) - t)\mathbf{x}$ |
| logistic regression | $\sigma(\bar{y})$ | Bernoulli    | $\text{NLL} \propto \mathbb{E} - \log(p(t \mathbf{x}))$ | $(y(\mathbf{x}) - t)\mathbf{x}$ |

# Multiclass Logistic Regression

To extend the binary logistic regression to a multiclass case with  $K$  classes, we:

- generate  $K$  outputs, each with its own set of weights, so that for  $\mathbf{W} \in \mathbb{R}^{D \times K}$ ,

$$\bar{\mathbf{y}}(\mathbf{x}; \mathbf{W}) = \mathbf{x}^T \mathbf{W}, \quad \text{or in other words, } \bar{\mathbf{y}}(\mathbf{x}; \mathbf{W})_i = \mathbf{x}^T (\mathbf{W}_{*,i})$$

- generalize the sigmoid function to a softmax function, such that

$$\text{softmax}(\mathbf{z})_i = \frac{e^{z_i}}{\sum_j e^{z_j}}.$$

Note that the original sigmoid function can be written as

$$\sigma(x) = \text{softmax} \left( \begin{bmatrix} x & 0 \end{bmatrix} \right)_0 = \frac{e^x}{e^x + e^0} = \frac{1}{1 + e^{-x}}.$$

The resulting classifier is also known as **multinomial logistic regression**, **maximum entropy classifier** or **softmax regression**.

# Multiclass Logistic Regression

Using the softmax function, we naturally define that

$$p(C_i|\mathbf{x}; \mathbf{W}) = \mathbf{y}(\mathbf{x}; \mathbf{W})_i = \text{softmax}(\bar{\mathbf{y}}(\mathbf{x}; \mathbf{W}))_i = \text{softmax}(\mathbf{x}^T \mathbf{W})_i = \frac{e^{(\mathbf{x}^T \mathbf{W})_i}}{\sum_j e^{(\mathbf{x}^T \mathbf{W})_j}}.$$

Considering the definition of the softmax function, it is natural to obtain the interpretation of the linear part of the model  $\bar{\mathbf{y}}(\mathbf{x}; \mathbf{W})$  as **logits** by computing a logarithm of the above:

$$\bar{\mathbf{y}}(\mathbf{x}; \mathbf{W})_i = \log(p(C_i|\mathbf{x}; \mathbf{W})) + c.$$

The constant  $c$  is present, because the output of the model is *overparametrized* (for example, the probability of the last class could be computed from the remaining ones). This is connected to the fact that softmax is invariant to addition of a constant:

$$\text{softmax}(\mathbf{z} + c)_i = \frac{e^{z_i + c}}{\sum_j e^{z_j + c}} = \frac{e^{z_i}}{\sum_j e^{z_j}} \cdot \frac{e^c}{e^c} = \text{softmax}(\mathbf{z})_i.$$

# Multiclass Logistic Regression

To train  $K$ -class classification, analogously to the binary logistic regression we can use MLE and train the model using minibatch stochastic gradient descent:

**Input:** Input dataset  $(\mathbf{X} \in \mathbb{R}^{N \times D}, \mathbf{t} \in \{0, 1, \dots, K - 1\}^N)$ , learning rate  $\alpha \in \mathbb{R}^+$ .

**Model:** Let  $\mathbf{w}$  denote all parameters of the model (in our case, the parameters are a weight matrix  $\mathbf{W}$  and maybe a bias vector  $\mathbf{b}$ ).

- $\mathbf{w} \leftarrow \mathbf{0}$  or we initialize  $\mathbf{w}$  randomly
- until convergence (or patience runs out), process a minibatch of examples  $\mathbb{B}$ :
  - $\mathbf{g} \leftarrow \frac{1}{|\mathbb{B}|} \sum_{i \in \mathbb{B}} \nabla_{\mathbf{w}} \left( -\log(p(C_{t_i} | \mathbf{x}_i; \mathbf{w})) \right)$
  - $\mathbf{w} \leftarrow \mathbf{w} - \alpha \mathbf{g}$



# Multiclass Logistic Regression

Note that the decision regions of the binary/multiclass logistic regression are convex (and therefore connected).

To see this, consider  $\mathbf{x}_A$  and  $\mathbf{x}_B$  in the same decision region  $\mathcal{R}_k$ .

Any point  $\mathbf{x}$  lying on the line connecting them is their convex combination,  $\mathbf{x} = \lambda \mathbf{x}_A + (1 - \lambda) \mathbf{x}_B$ , and from the linearity of  $\bar{\mathbf{y}}(\mathbf{x}) = \mathbf{x}^T \mathbf{W}$  it follows that

$$\bar{\mathbf{y}}(\mathbf{x}) = \lambda \bar{\mathbf{y}}(\mathbf{x}_A) + (1 - \lambda) \bar{\mathbf{y}}(\mathbf{x}_B).$$

Given that  $\bar{\mathbf{y}}(\mathbf{x}_A)_k$  was the largest among  $\bar{\mathbf{y}}(\mathbf{x}_A)$  and also given that  $\bar{\mathbf{y}}(\mathbf{x}_B)_k$  was the largest among  $\bar{\mathbf{y}}(\mathbf{x}_B)$ , it must be the case that  $\bar{\mathbf{y}}(\mathbf{x})_k$  is the largest among all  $\bar{\mathbf{y}}(\mathbf{x})$ .

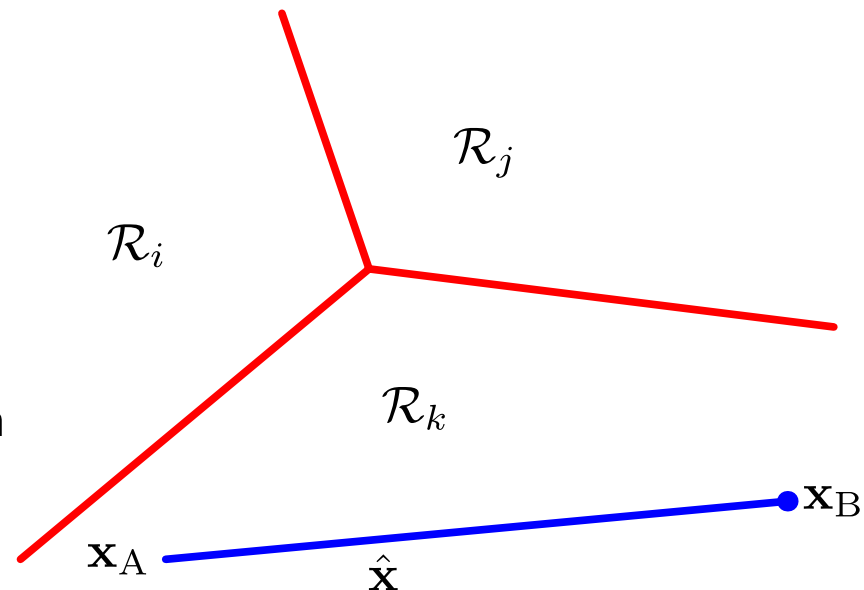
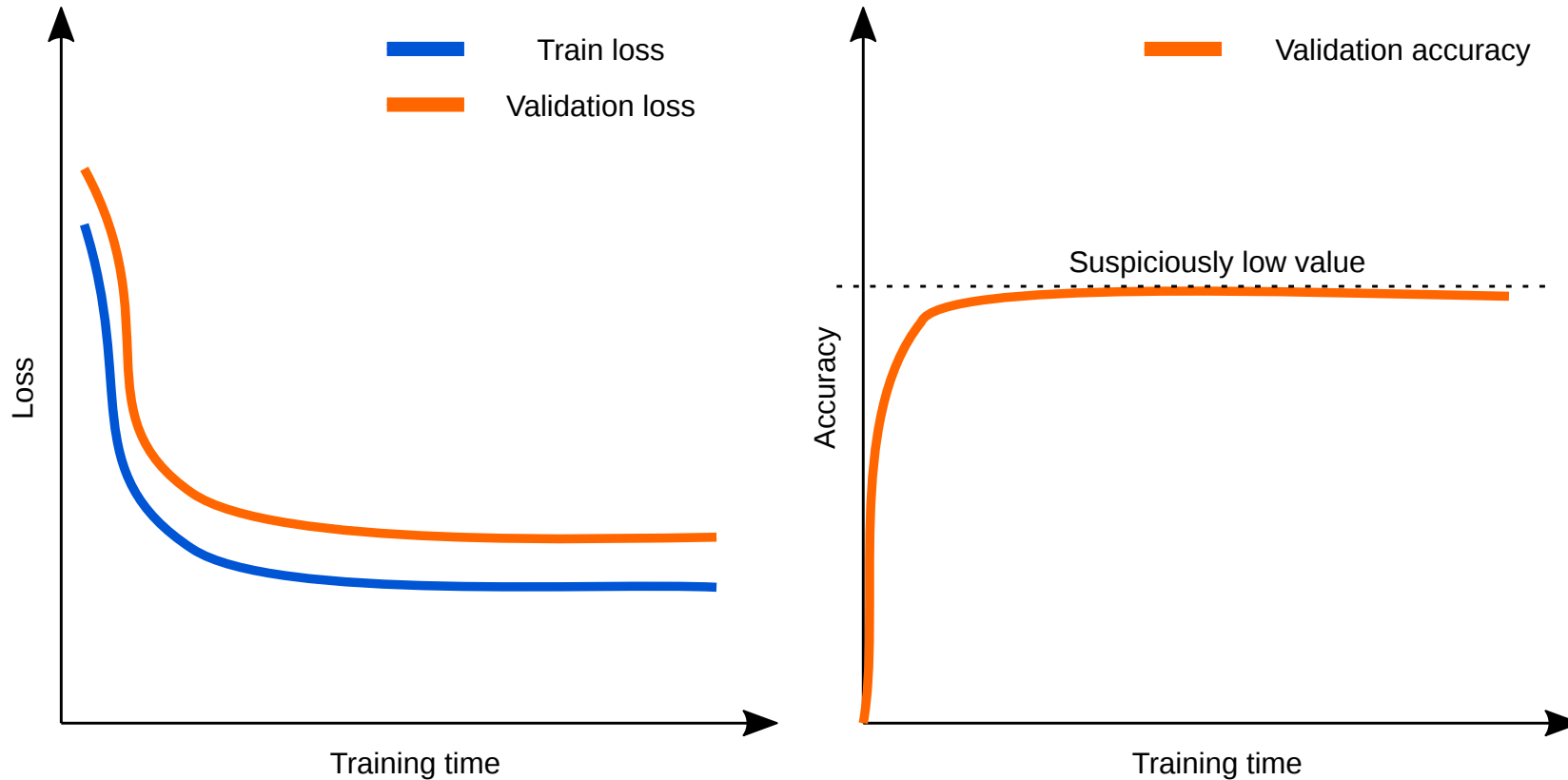


Figure 4.3 of Pattern Recognition and Machine Learning.

# What went wrong?



The model only predicts the majority class. Insufficient features, too high learning rate.

# Generalized Linear Models

The multiclass logistic regression can now be added to the GLM table:

| Name                              | Activation                         | Distribution | Loss  | Gradient  |
|-----------------------------------|------------------------------------|--------------|---|---|
| linear regression                 | identity                           | Normal       | $\text{NLL} \propto \text{MSE}$                         | $(y(\mathbf{x}) - t)\mathbf{x}$                           |
| logistic regression               | $\sigma(\bar{y})$                  | Bernoulli    | $\text{NLL} \propto \mathbb{E} - \log(p(t \mathbf{x}))$ | $(y(\mathbf{x}) - t)\mathbf{x}$                           |
| multiclass<br>logistic regression | $\text{softmax}(\bar{\mathbf{y}})$ | categorical  | $\text{NLL} \propto \mathbb{E} - \log(p(t \mathbf{x}))$ | $((\mathbf{y}(\mathbf{x}) - \mathbf{1}_t)\mathbf{x}^T)^T$ |

Recall that  $\mathbf{1}_t = ([i = t])_{i=0}^{K-1}$  is one-hot representation of target  $t \in \{0, 1, \dots, K - 1\}$ .

The gradient  $((\mathbf{y}(\mathbf{x}) - \mathbf{1}_t)\mathbf{x}^T)^T$  can be of course also computed as  $\mathbf{x}(\mathbf{y}(\mathbf{x}) - \mathbf{1}_t)^T$ .

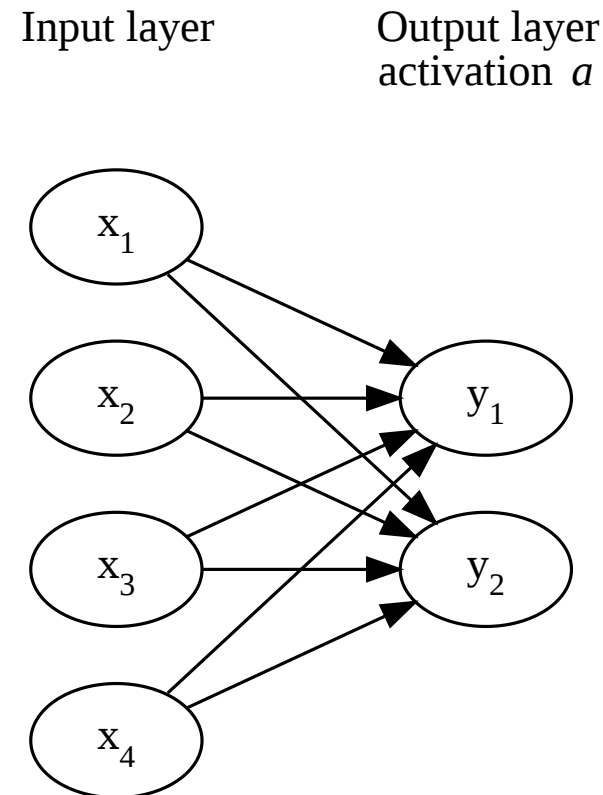
# Multilayer Perceptron

We can reformulate the generalized linear models in the following framework.

- Assume we have an input node for every input feature.
- Additionally, we have an output node for every model output (one for linear regression or binary classification,  $K$  for classification in  $K$  classes).
- Every input node and output node are connected with a directed edge, and every edge has an associated weight.
- Value of every (output) node is computed by summing the values of predecessors multiplied by the corresponding weights, added to a bias of this node, and finally passed through an activation function  $a$ :

$$y_i = a \left( \sum_j x_j w_{j,i} + b_i \right)$$

or in matrix form  $\mathbf{y} = a(\mathbf{x}^T \mathbf{W} + \mathbf{b})$ , or for a batch of examples  $\mathbf{X}, \mathbf{Y} = a(\mathbf{XW} + \mathbf{b})$ .



# Multilayer Perceptron

We now extend the model by adding a **hidden layer** with activation  $f$ .

- The computation is performed analogously:

$$h_i = f \left( \sum_j x_j w_{j,i}^{(h)} + b_i^{(h)} \right),$$

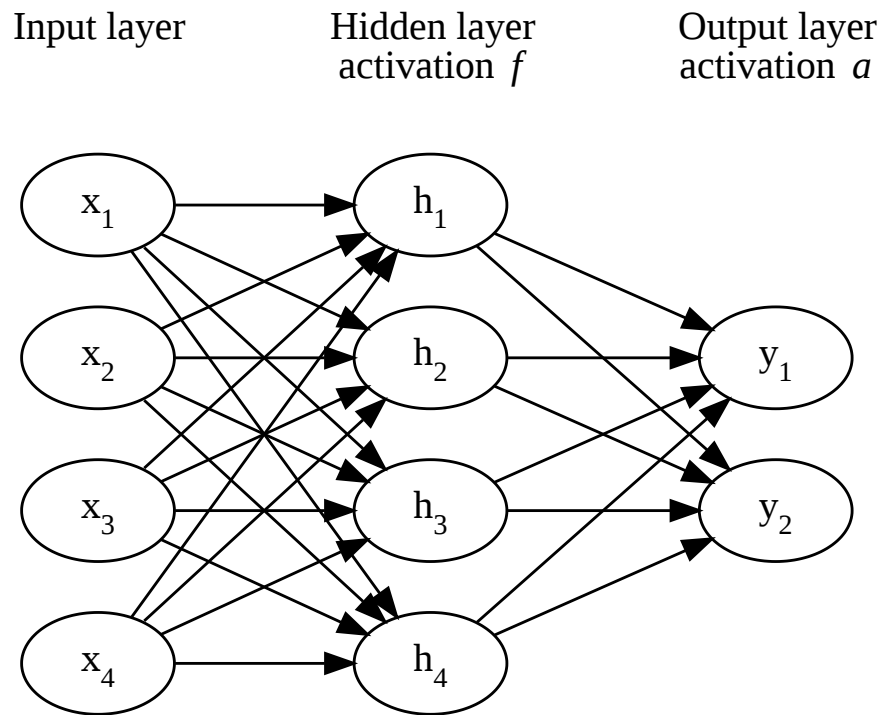
$$y_i = a \left( \sum_j h_j w_{j,i}^{(y)} + b_i^{(y)} \right),$$

or in matrix form

$$\mathbf{h} = f \left( \mathbf{x}^T \mathbf{W}^{(h)} + \mathbf{b}^{(h)} \right),$$

$$\mathbf{y} = a \left( \mathbf{h}^T \mathbf{W}^{(y)} + \mathbf{b}^{(y)} \right),$$

and for batch of inputs  $\mathbf{H} = f \left( \mathbf{XW}^{(h)} + \mathbf{b}^{(h)} \right)$  and  $\mathbf{Y} = a \left( \mathbf{HW}^{(y)} + \mathbf{b}^{(y)} \right)$ .



# Multilayer Perceptron

Note that:

- the structure of the *input* layer depends on the input features;
- the structure and the *activation* function of the *output* layer depends on the target data;
- however, the *hidden* layer has no pre-image in the data and is completely arbitrary – which is the reason why it is called a *hidden* layer.

Also note that we can absorb biases into weights analogously to the generalized linear models.

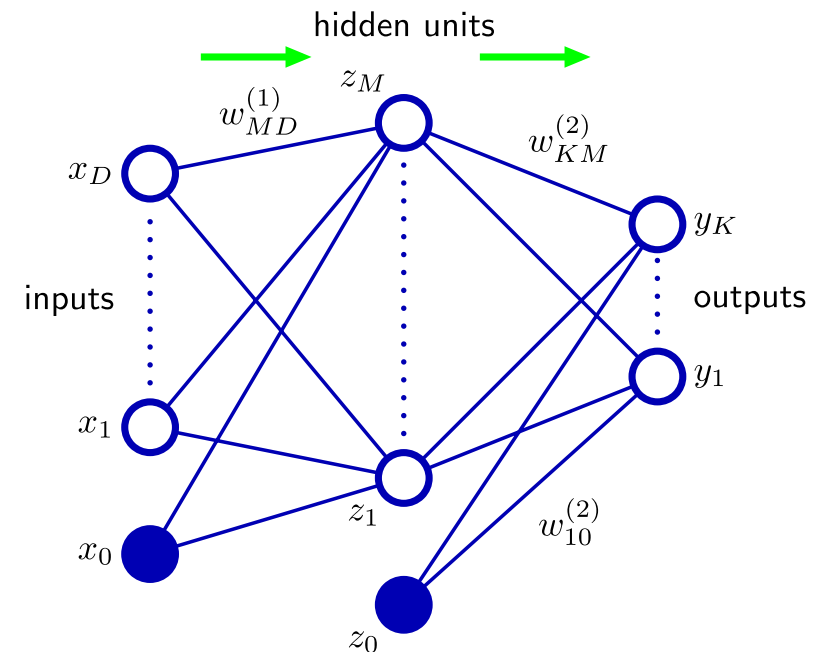


Figure 5.1 of Pattern Recognition and Machine Learning.

## Output Layer Activation Functions

- regression:
  - identity activation: we model normal distribution on output (linear regression)
- binary classification:
  - $\sigma(\boldsymbol{x})$ : we model the Bernoulli distribution (the model predicts a probability)

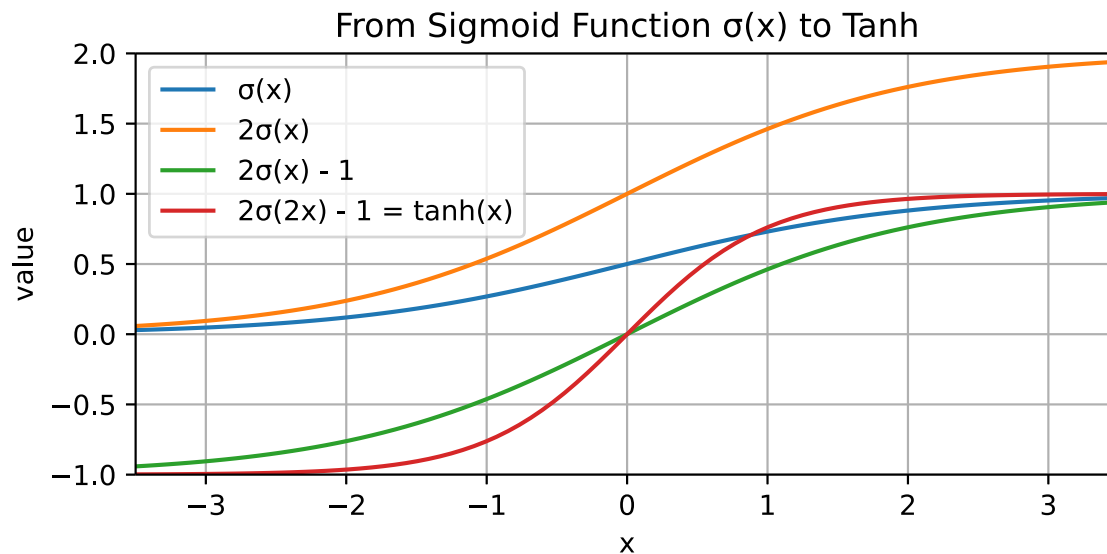
$$\sigma(x) \stackrel{\text{def}}{=} \frac{1}{1 + e^{-x}}$$

- $K$ -class classification:
  - $\text{softmax}(\boldsymbol{x})$ : we model the (usually overparametrized) categorical distribution

$$\text{softmax}(\boldsymbol{x}) \propto e^{\boldsymbol{x}}, \quad \text{softmax}(\boldsymbol{x})_i \stackrel{\text{def}}{=} \frac{e^{x_i}}{\sum_j e^{x_j}}$$

## Hidden Layer Activation Functions

- no activation (identity): does not help, composition of linear mapping is a linear mapping
- $\sigma$  (but works suboptimally – nonsymmetrical,  $\frac{d\sigma}{dx}(0) = 1/4$ )
- $\tanh$ 
  - result of making  $\sigma$  symmetrical and making derivation in zero 1
  - $\tanh(x) = 2\sigma(2x) - 1$
- ReLU
  - $\max(0, x)$
  - the most common nonlinear activation used nowadays





The multilayer perceptron can be trained using again a minibatch SGD algorithm:

**Input:** Input dataset ( $\mathbf{X} \in \mathbb{R}^{N \times D}$ ,  $\mathbf{t}$  targets), learning rate  $\alpha \in \mathbb{R}^+$ .

**Model:** Let  $\mathbf{w}$  denote all parameters of the model (all weight matrices and bias vectors).

- initialize  $\mathbf{w}$ 
  - set weights randomly
    - for a weight matrix processing a layer of size  $M$  to a layer of size  $O$ , we can sample its elements uniformly for example from the  $\left[-\frac{1}{\sqrt{M}}, \frac{1}{\sqrt{M}}\right]$  range
    - the exact range becomes more important for networks with many hidden layers
  - set biases to 0
- until convergence (or patience runs out), process a minibatch of examples  $\mathbb{B}$ :
  - $\mathbf{g} \leftarrow \frac{1}{|\mathbb{B}|} \sum_{i \in \mathbb{B}} \nabla_{\mathbf{w}} \left( -\log(p(t_i | \mathbf{x}_i; \mathbf{w})) \right)$
  - $\mathbf{w} \leftarrow \mathbf{w} - \alpha \mathbf{g}$