

Multiclass Logistic Regression, Multilayer Perceptron

Milan Straka

 October 24, 2022



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



unless otherwise stated

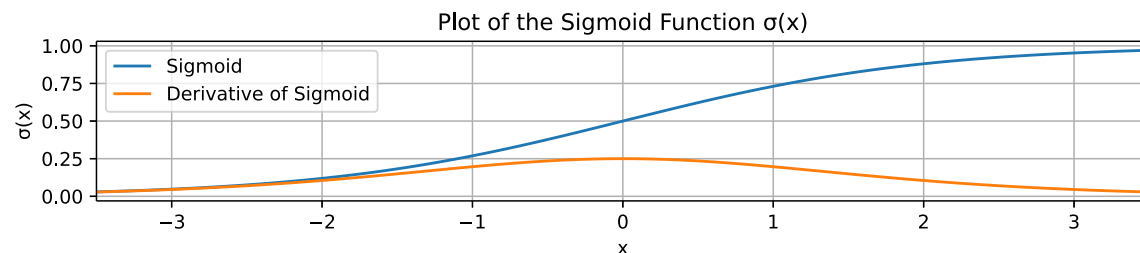
An extension of perceptron, which models the conditional probabilities of $p(C_0|\mathbf{x})$ and of $p(C_1|\mathbf{x})$. Logistic regression can in fact handle also more than two classes, which we will see shortly.

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 σ is a **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.

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} :
 - $\mathbf{g} \leftarrow \frac{1}{|\mathbb{B}|} \sum_{i \in \mathbb{B}} \nabla_{\mathbf{w}} \left(-\log(p(C_{t_i}|\mathbf{x}_i; \mathbf{w})) \right)$
 - $\mathbf{w} \leftarrow \mathbf{w} - \alpha \mathbf{g}$

Linearity in Logistic Regression

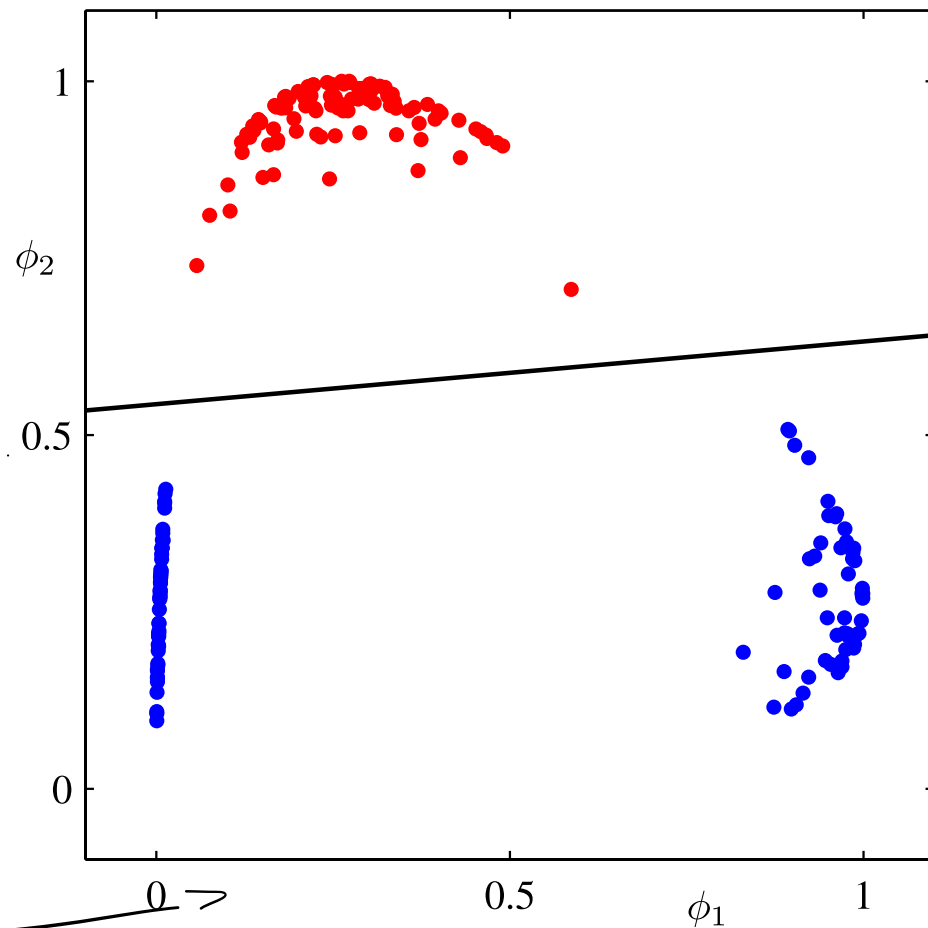
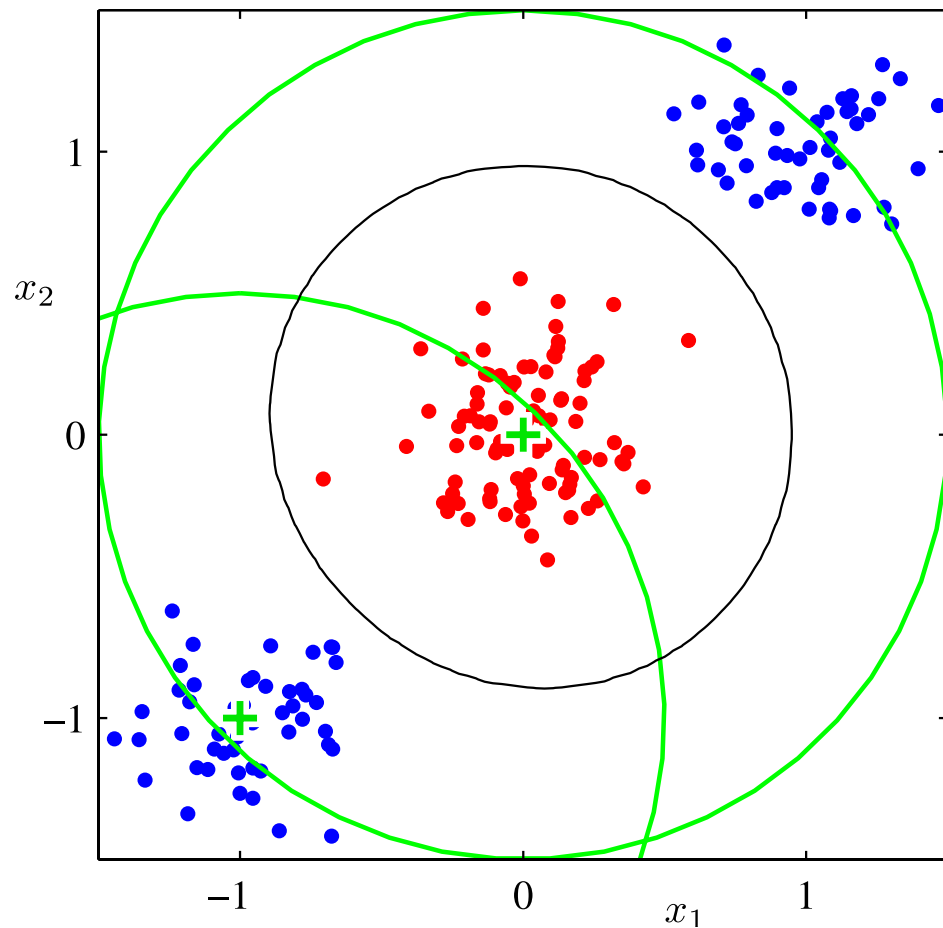


Figure 4.12 of Pattern Recognition and Machine Learning.

Tady jsou si vytvořili „mapping“ na vzdálenost 2 bodů.

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 | $\text{NLL} \propto \mathbb{E} - \log(p(t \mathbf{x}))$ | ? |

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

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

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

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:

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

*> to říct to, že někdy máš definiční $p(C_1) \Rightarrow t=1$
nebo $p(C_0) \Rightarrow t=0$ a vždy jedná
o členů toho výrazu je 1. - tedy nic měnit.*

$$\nabla_{\mathbf{w}} - \log p(t|\mathbf{x}; \mathbf{w}) =$$

$$= \nabla_{\mathbf{w}} \left(-t \cdot \log(\sigma(\bar{y})) - (1-t) \cdot \log(1 - \sigma(\bar{y})) \right)$$

$$\frac{\partial}{\partial x} \log g(x) = \frac{1}{g(x)} \cdot \frac{\partial}{\partial x} g(x)$$

$$= -t \cdot \frac{1}{\sigma(\bar{y})} \cdot \nabla_{\mathbf{w}} \sigma(\bar{y}) - (1-t) \cdot \frac{1}{1 - \sigma(\bar{y})} \cdot \nabla_{\mathbf{w}} (1 - \sigma(\bar{y}))$$

$$\frac{\partial}{\partial x} f(g(x)) = \frac{\partial}{\partial g(x)} f(g(x)) \cdot \frac{\partial}{\partial x} g(x) = \frac{\partial}{\partial z} f(z) \cdot \frac{\partial}{\partial x} g(x)$$

$$\nabla_{\mathbf{w}} \sigma(\bar{y}) = \frac{\partial}{\partial \bar{y}} \sigma(\bar{y}) \cdot \nabla_{\mathbf{w}} \bar{y}$$

$$= -t \cdot \frac{1}{\sigma(\bar{y})} \cdot \sigma(\bar{y}) \cdot (1 - \sigma(\bar{y})) \cdot \nabla_{\mathbf{w}} \bar{y} + (1-t) \cdot \frac{1}{1 - \sigma(\bar{y})} \cdot \sigma(\bar{y}) \cdot (1 - \sigma(\bar{y})) \cdot \nabla_{\mathbf{w}} \bar{y}$$

$$= (-t + t\sigma(\bar{y}) + \sigma(\bar{y}) - t\sigma(\bar{y})) \mathbf{x}$$

$$= (y(\mathbf{x}; \mathbf{w}) - t) \mathbf{x}$$

$$\nabla_{\mathbf{w}} \bar{y} = \nabla_{\mathbf{w}} \mathbf{x}^T \mathbf{w} = \mathbf{x}$$

Generalized Linear Models

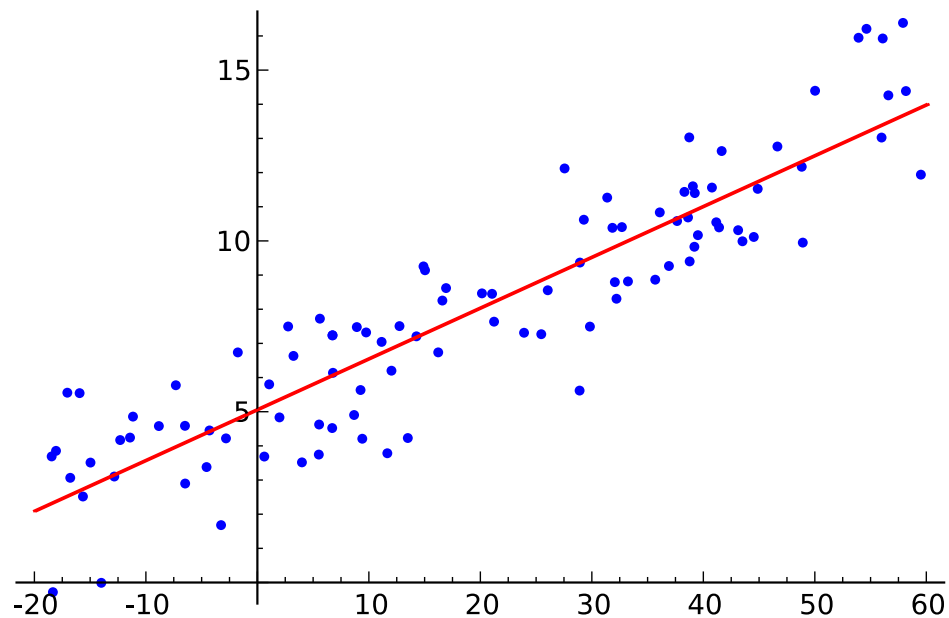
The logistic regression is in fact an extended linear regression. A linear regression model, which is followed by some **activation function** a , is called **generalized linear model**:

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

| 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

During regression, we predict a number, not a real probability distribution. In order to generate a distribution, we might consider a distribution with the mean of the predicted value and a fixed variance σ^2 – the most general such a distribution is the normal distribution.



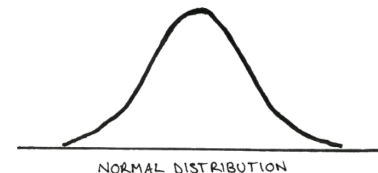
https://upload.wikimedia.org/wikipedia/commons/3/3a/Linear_regression.svg

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 | $\text{NLL} \propto \text{MSE}$ | $(y(\boldsymbol{x}) - t)\boldsymbol{x}$ |
| logistic regression | $\sigma(\bar{y})$ | Bernoulli | $\text{NLL} \propto \mathbb{E} - \log(p(t \boldsymbol{x}))$ | $(y(\boldsymbol{x}) - t)\boldsymbol{x}$ |

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
dan di softmax na 1.
resulaty lin. modelu

$$\text{softmax}(\mathbf{z})_i = \frac{e^{z_i}}{\sum_j e^{z_j}} \cdot 1$$

— maximum viny vedle sebe

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**.

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:

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

The difference between softmax and sigmoid output can be compared on the binary case, where the binary logistic regression outputs of the linear part of the model are

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

while the outputs of the softmax variant with two outputs can be interpreted as $\bar{\mathbf{y}}(\mathbf{x}; \mathbf{W})_0 = \log(p(C_0|\mathbf{x}; \mathbf{W})) + c$ and $\bar{\mathbf{y}}(\mathbf{x}; \mathbf{W})_1 = \log(p(C_1|\mathbf{x}; \mathbf{W})) + c$.

If we consider $\bar{\mathbf{y}}(\mathbf{x}; \mathbf{W})_0$ to be zero, the model can then predict only the probability $p(C_1|\mathbf{x})$, and the constant c is fixed to $-\log(p(C_0|\mathbf{x}; \mathbf{W}))$, recovering the original interpretation.

Generalizing to a K -class classification, we could produce only $K - 1$ outputs and define $\bar{\mathbf{y}}_0 = 0$, resulting in the interpretation of the linear part outputs analogous to the binary case:

$$\bar{\mathbf{y}}(\mathbf{x}; \mathbf{W})_i = \log \left(\frac{p(C_i|\mathbf{x}; \mathbf{W})}{p(C_0|\mathbf{x}; \mathbf{W})} \right).$$

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 $\bar{\mathbf{y}}(\mathbf{x})$.

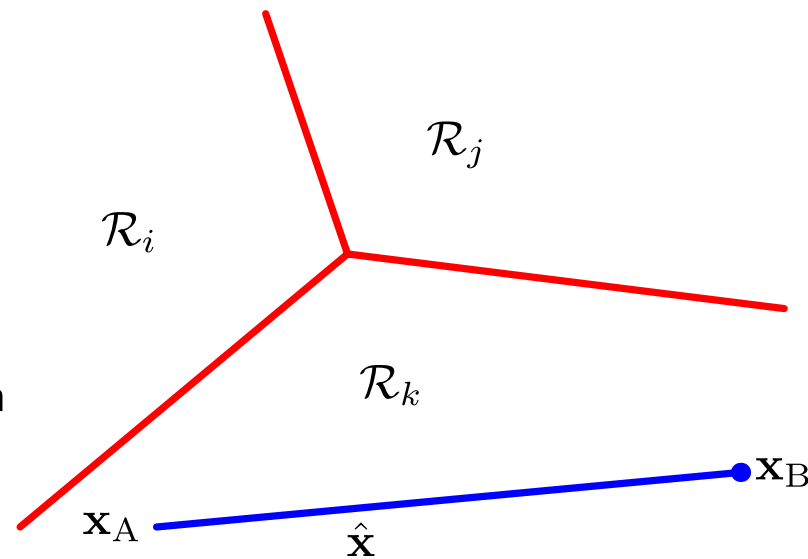


Figure 4.3 of Pattern Recognition and Machine Learning.

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 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 $\mathbf{x}(\mathbf{y}(\mathbf{x}) - \mathbf{1}_t)^T$.

Several other GLMs exist, we now describe a final one, this time for regression and not for classification. Compared to regular linear regression, where we assume the output distribution is normal, we turn our attention to **Poisson distribution**.

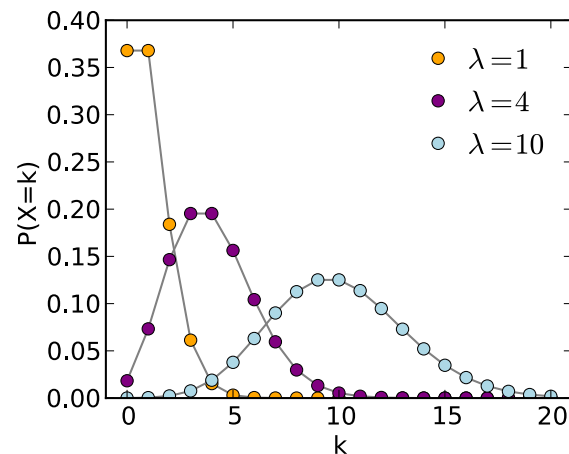
Poisson Distribution

Poisson distribution is a discrete distribution suitable for modeling the probability of a given number of events occurring in a fixed time interval, if these events occur at a known rate and independently of each other.

$$P(x = k; \lambda) = \frac{\lambda^k e^{-\lambda}}{k!}$$

It is easy to show that if x has Poisson distribution,

$$\begin{aligned}\mathbb{E}[x] &= \lambda \\ \text{Var}(x) &= \lambda\end{aligned}$$

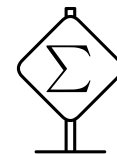


https://upload.wikimedia.org/wikipedia/commons/1/16/Poisson_pmf.svg

Poisson Distribution Derivation

The Poisson distribution can be obtained as a limit of the binomial distribution.

Assume we are considering n independent events, each with probability p_n , and that np_n converges to λ . Then



$$\lim_{n \rightarrow \infty} \binom{n}{k} p_n^k (1 - p_n)^{n-k} = e^{-\lambda} \frac{\lambda^k}{k!}.$$

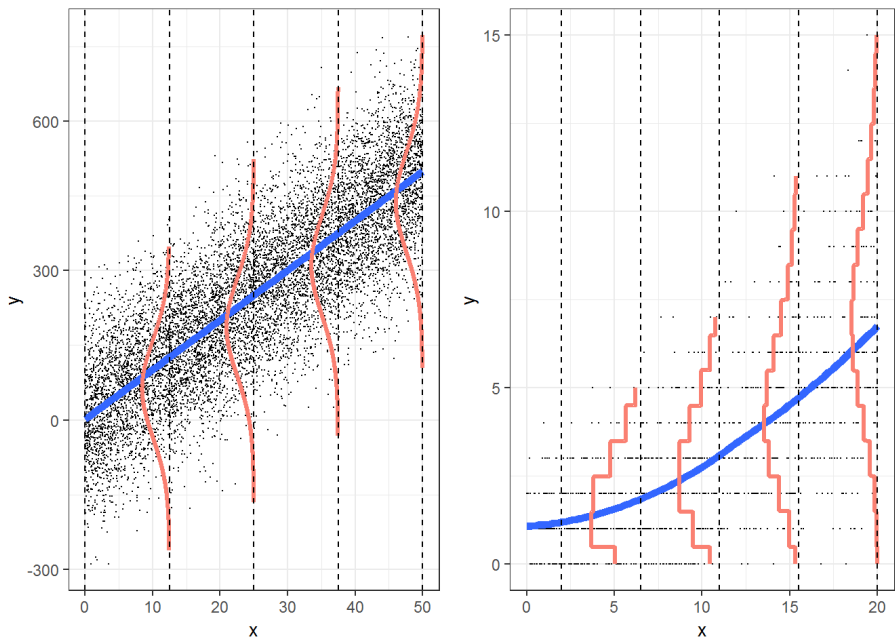
$$\begin{aligned} \lim_{n \rightarrow \infty} \binom{n}{k} p_n^k (1 - p_n)^{n-k} &= \lim_{n \rightarrow \infty} \frac{n(n-1)(n-2) \cdots (n-k+1)}{k!} \left(\frac{\lambda}{n}\right)^k \left(1 - \frac{\lambda}{n}\right)^{n-k} \\ &= \lim_{n \rightarrow \infty} \frac{n^k + \mathcal{O}(n^{k-1})}{k!} \frac{\lambda^k}{n^k} \left(1 - \frac{\lambda}{n}\right)^{n-k} = \lim_{n \rightarrow \infty} \frac{\lambda^k}{k!} \left(1 - \frac{\lambda}{n}\right)^{n-k} \end{aligned}$$

and the result follows, since $\lim_{n \rightarrow \infty} \left(1 - \frac{\lambda}{n}\right)^n = e^{-\lambda}$ and $\lim_{n \rightarrow \infty} \left(1 - \frac{\lambda}{n}\right)^{-k} = 1$.

Poisson Distribution

An important difference compared to the normal distribution is that the latter assumes that the variance does not depend on the mean, i.e., that the model “makes errors of the same magnitude everywhere”.

On the other hand, the variance of a Poisson distribution increases with the mean. It is useful if we want to measure error relatively, not as an absolute difference.



https://bookdown.org/roback/bookdown-bysh/bookdown-bysh_files/figure-html/OLSpois-1.png

Poisson Regression

Poisson regression is a generalized linear model producing a Poisson distribution (i.e., the mean rate λ).

Again, we use NLL as the loss. To choose a suitable activation, we might be interested in obtaining the same gradient as for other GLMs – solving for an activation function while requiring the gradient to be $(a(\bar{y}(\mathbf{x})) - t) \cdot \mathbf{x}$ yields $a(\bar{y}) = \exp(\bar{y})$, which means the linear part of the model is predicting $\log(\lambda)$.

| 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 | $\text{softmax}(\bar{\mathbf{y}})$ | categorical | $\text{NLL} \propto \mathbb{E} - \log(p(t \mathbf{x}))$ | $((\mathbf{y}(\mathbf{x}) - \mathbf{1}_t)\mathbf{x}^T)^T$ |
| Poisson regression | $\exp(\bar{y})$ | Poisson | $\text{NLL} \propto \mathbb{E} - \log(p(t \mathbf{x}))$ | $(y(\mathbf{x}) - t)\mathbf{x}$ |

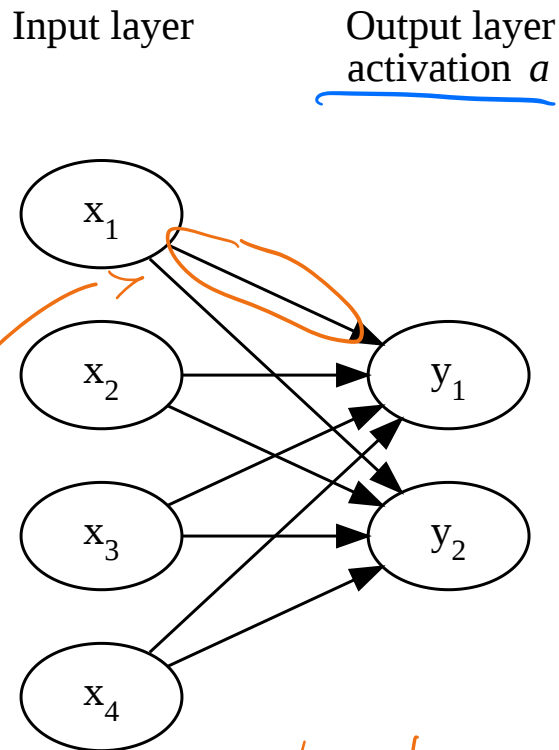
Multilayer Perceptron

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 :

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

vyjádření hodnoty hmoty výstupu



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

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

- The computation is performed analogously:

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

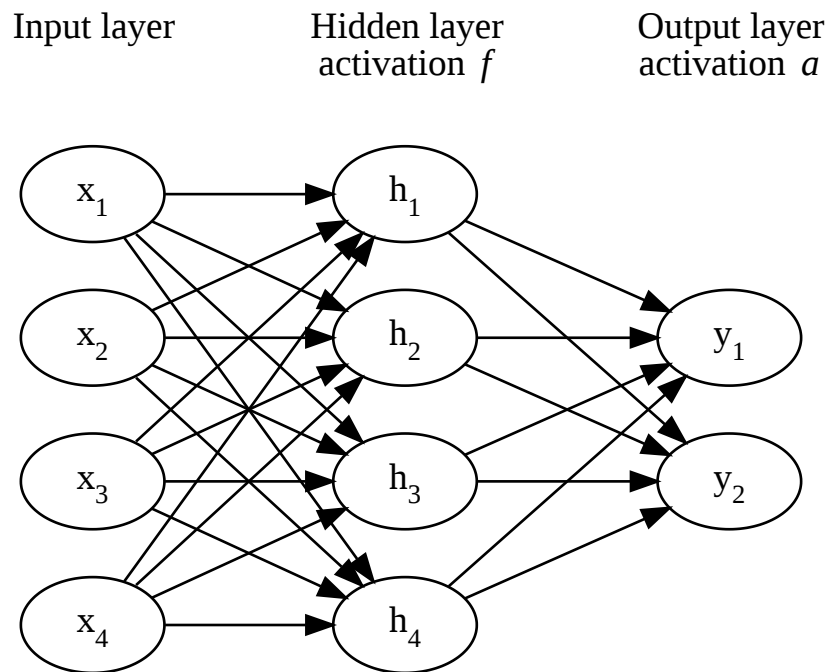
$$y_i = a \left(\sum_j h_j w_{j,i}^{(y)} + b_i^{(y)} \right),$$

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.

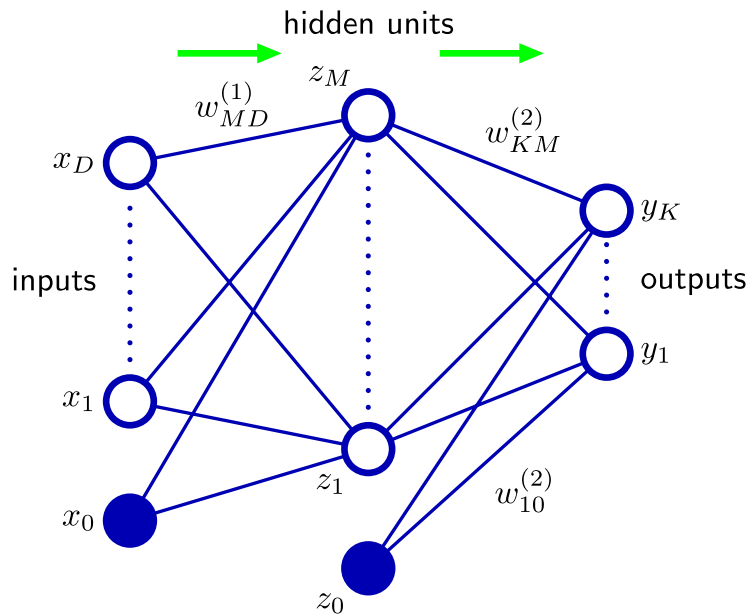


Figure 5.1 of Pattern Recognition and Machine Learning.

Output Layer Activation Functions

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

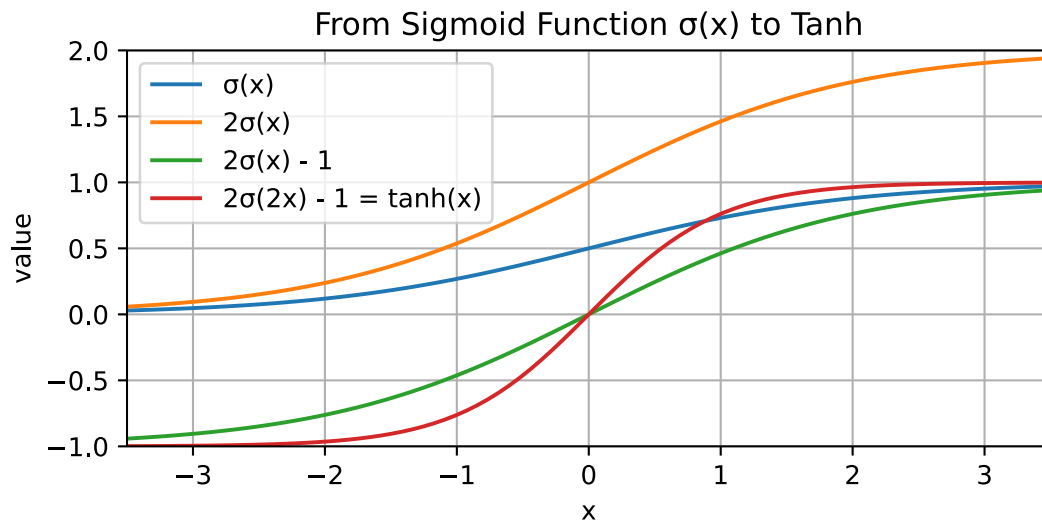
$$\sigma(x) \stackrel{\text{def}}{=} \frac{1}{1 + e^{-x}}$$

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

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

Hidden Layer Activation Functions

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



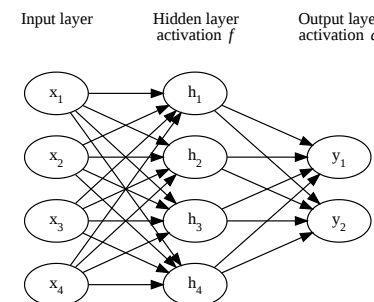
The multilayer perceptron can be trained using again a minibatch SGD algorithm:

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} :
 - $\mathbf{g} \leftarrow \frac{1}{|\mathbb{B}|} \sum_{i \in \mathbb{B}} \nabla_{\mathbf{w}} \left(-\log(p(t_i | \mathbf{x}_i; \mathbf{w})) \right)$
 - $\mathbf{w} \leftarrow \mathbf{w} - \alpha \mathbf{g}$

Assume we have an MLP with input of size D , weights $\mathbf{W}^{(h)} \in \mathbb{R}^{D \times H}$, $\mathbf{b}^{(h)} \in \mathbb{R}^H$, hidden layer of size H and activation f with weights $\mathbf{W}^{(y)} \in \mathbb{R}^{H \times K}$, $\mathbf{b}^{(y)} \in \mathbb{R}^K$, and finally an output layer of size K with activation a .



In order to compute the gradient of the loss L with respect to all weights, you should proceed gradually:

- first compute $\frac{\partial L}{\partial \mathbf{y}}$,
- then compute $\frac{\partial \mathbf{y}}{\partial \mathbf{y}^{(in)}}$, where $\mathbf{y}^{(in)}$ are the inputs to the output layer (i.e., before applying activation function a ; in other words, $\mathbf{y} = a(\mathbf{y}^{(in)})$),
- then compute $\frac{\partial \mathbf{y}^{(in)}}{\partial \mathbf{W}^{(y)}}$ and $\frac{\partial \mathbf{y}^{(in)}}{\partial \mathbf{b}^{(y)}}$, which allows us to obtain $\frac{\partial L}{\partial \mathbf{W}^{(y)}} = \frac{\partial L}{\partial \mathbf{y}} \cdot \frac{\partial \mathbf{y}}{\partial \mathbf{y}^{(in)}} \cdot \frac{\partial \mathbf{y}^{(in)}}{\partial \mathbf{W}^{(y)}}$ and analogously $\frac{\partial L}{\partial \mathbf{b}^{(y)}}$,
- followed by $\frac{\partial \mathbf{y}^{(in)}}{\partial \mathbf{h}}$ and $\frac{\partial \mathbf{h}}{\partial \mathbf{h}^{(in)}}$,
- and finally using $\frac{\partial \mathbf{h}^{(in)}}{\partial \mathbf{W}^{(h)}}$ and $\frac{\partial \mathbf{h}^{(in)}}{\partial \mathbf{b}^{(h)}}$ to compute $\frac{\partial L}{\partial \mathbf{W}^{(h)}}$ and $\frac{\partial L}{\partial \mathbf{b}^{(h)}}$.

One way how to interpret the hidden layer is:

- the part from the hidden layer to the output layer is the previously used generalized linear model (linear regression, logistic regression, ...);
- the part from the inputs to the hidden layer can be considered automatically constructed features.

The features are a linear mapping of the input values followed by a nonlinearity, and the theorem on the next slide proves they can always be constructed to achieve as good a fit of the training data as is required.

Note that the weights in an MLP must be initialized randomly. If we used just zeros, all the constructed features (hidden layer nodes) would behave identically and we would never distinguish them.

Using random weights corresponds to starting with random features, which allows the SGD to make progress (improve the individual features).

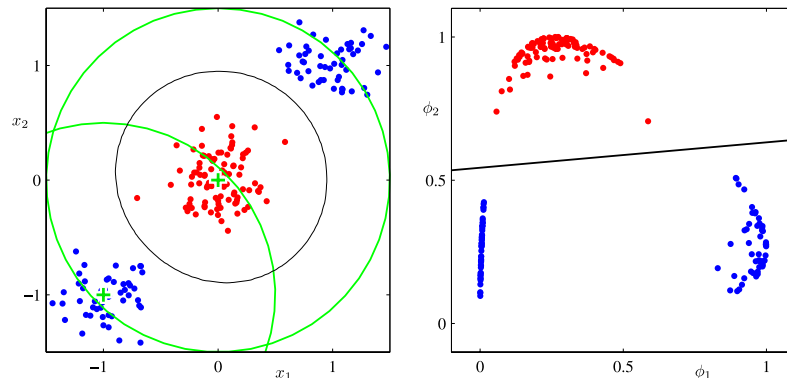


Figure 4.12 of *Pattern Recognition and Machine Learning*.

Let $\varphi(x) : \mathbb{R} \rightarrow \mathbb{R}$ be a nonconstant, bounded and nondecreasing continuous function.
(Later a proof was given also for $\varphi = \text{ReLU}$ and even for any nonpolynomial function.)

For any $\varepsilon > 0$ and any continuous function $f : [0, 1]^D \rightarrow \mathbb{R}$, there exists $H \in \mathbb{N}$, $\mathbf{v} \in \mathbb{R}^H$, $\mathbf{b} \in \mathbb{R}^H$ and $\mathbf{W} \in \mathbb{R}^{D \times H}$, such that if we denote

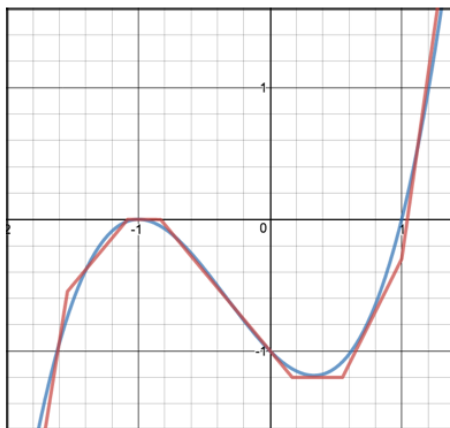
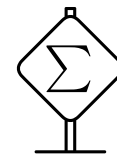
$$F(\mathbf{x}) = \mathbf{v}^T \varphi(\mathbf{x}^T \mathbf{W} + \mathbf{b}) = \sum_{i=1}^H v_i \varphi(\mathbf{x}^T \mathbf{W}_{*,i} + b_i),$$

where φ is applied elementwise, then for all $\mathbf{x} \in [0, 1]^D$:

$$|F(\mathbf{x}) - f(\mathbf{x})| < \varepsilon.$$

Sketch of the proof:

- If a function is continuous on a closed interval, it can be approximated by a sequence of lines to arbitrary precision.



https://miro.medium.com/max/844/1*lihbPNQgl7oKjpCsmzPDKw.png

$$n_1(x) = \text{Relu}(-5x - 7.7)$$

$$n_2(x) = \text{Relu}(-1.2x - 1.3)$$

$$n_3(x) = \text{Relu}(1.2x + 1)$$

$$n_4(x) = \text{Relu}(1.2x - .2)$$

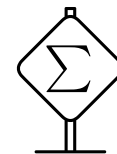
$$n_5(x) = \text{Relu}(2x - 1.1)$$

$$n_6(x) = \text{Relu}(5x - 5)$$

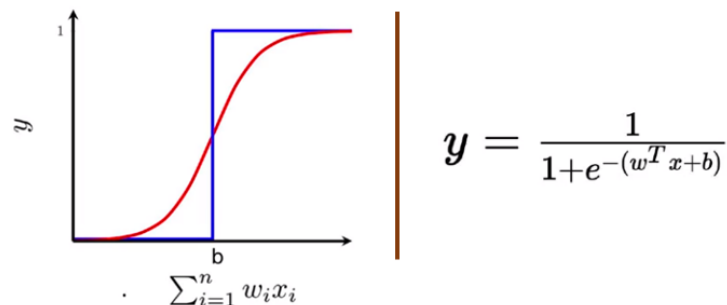
$$Z(x) = -n_1(x) - n_2(x) - n_3(x) \\ + n_4(x) + n_5(x) + n_6(x)$$

- However, we can create a sequence of k linear segments as a sum of k ReLU units – on every endpoint a new ReLU starts (i.e., the input ReLU value is zero at the endpoint), with a tangent which is the difference between the target tangent and the tangent of the approximation until this point.

Sketch of the proof for a squashing function $\varphi(x)$ (i.e., nonconstant, bounded and nondecreasing continuous function like sigmoid):

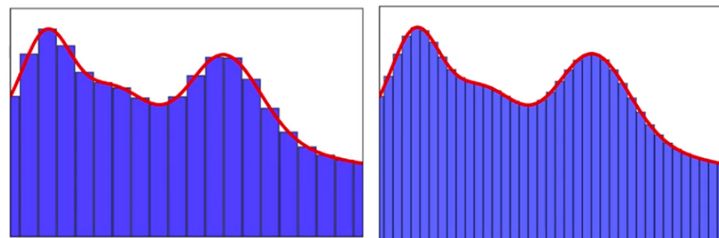


- We can prove φ can be arbitrarily close to a hard threshold by compressing it horizontally.



https://hackernoon.com/hn-images/1*N7dfPwbiXC-Kk4TCbfRerA.png

- Then we approximate the original function using a series of straight line segments



https://hackernoon.com/hn-images/1*hVuJgUTLUFWTMmJhl_fomg.png