

Multiclass Logistic Regression, Multiplayer Perceptron

Milan Straka

 November 11, 2019



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:

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

where σ is a *sigmoid function*

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

Can be trained using an SGD algorithm.

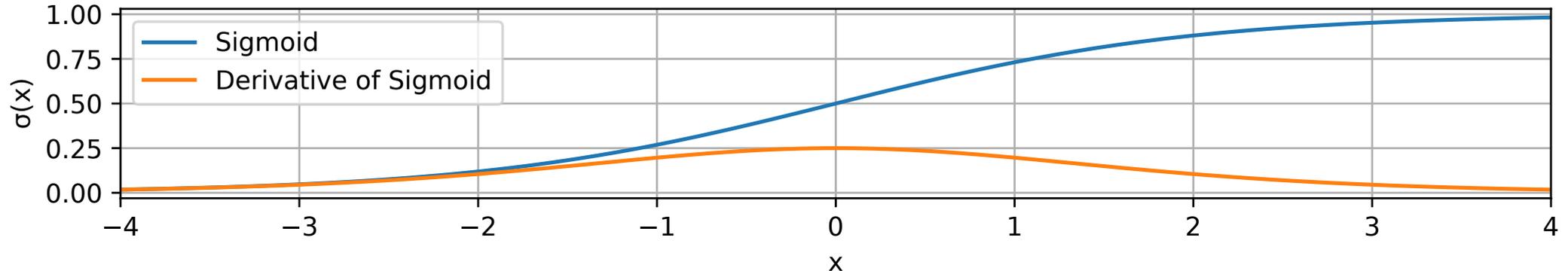
Sigmoid Function

The sigmoid function has values in range $(0, 1)$, it is monotonically increasing and it has a derivative of $\frac{1}{4}$ at $x = 0$.

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

$$\sigma'(x) = \sigma(x)(1 - \sigma(x))$$

Plot of the Sigmoid Function $\sigma(x)$



To give some meaning to the sigmoid function, starting with

$$P(C_1|\mathbf{x}) = \sigma(f(\mathbf{x}; \mathbf{w})) = \frac{1}{1 + e^{-f(\mathbf{x}; \mathbf{w})}}$$

we can arrive at

$$f(\mathbf{x}; \mathbf{w}) = \log \left(\frac{P(C_1|\mathbf{x})}{P(C_0|\mathbf{x})} \right),$$

where the prediction of the model $f(\mathbf{x}; \mathbf{w})$ is called a *logit* and it is a logarithm of odds of the two classes probabilities.

Logistic Regression

To train the logistic regression $y(\mathbf{x}; \mathbf{w}) = \mathbf{x}^T \mathbf{w}$, we use MLE (the maximum likelihood estimation). Note that $P(C_1 | \mathbf{x}; \mathbf{w}) = \sigma(y(\mathbf{x}; \mathbf{w}))$.

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

$$\mathcal{L}(\mathbb{X}) = \frac{1}{N} \sum_i -\log(P(C_{t_i} | \mathbf{x}_i; \mathbf{w})).$$

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

- $\mathbf{w} \leftarrow \mathbf{0}$
- until convergence (or until patience is over), process batch of N examples:
 - $\mathbf{g} \leftarrow -\frac{1}{N} \sum_i \nabla_{\mathbf{w}} \log(P(C_{t_i} | \mathbf{x}_i; \mathbf{w}))$
 - $\mathbf{w} \leftarrow \mathbf{w} - \alpha \mathbf{g}$

Linearity in Logistic Regression

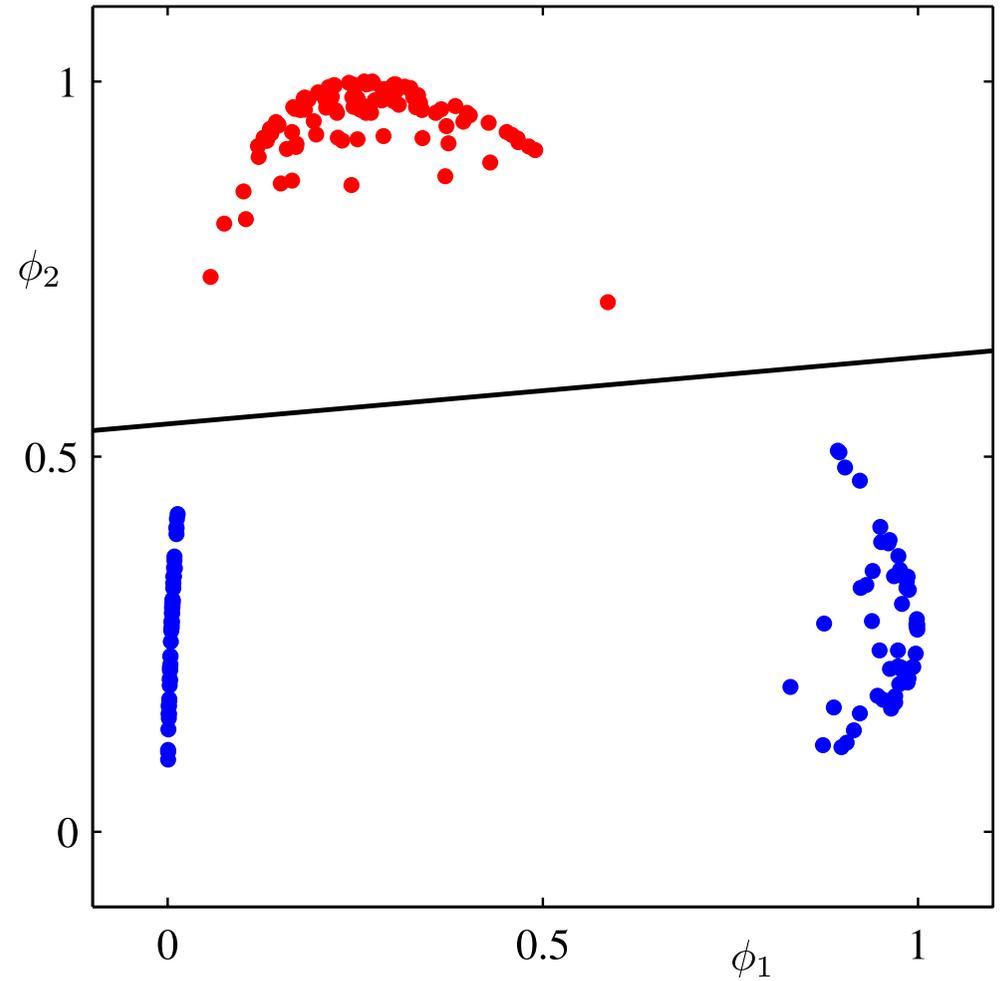
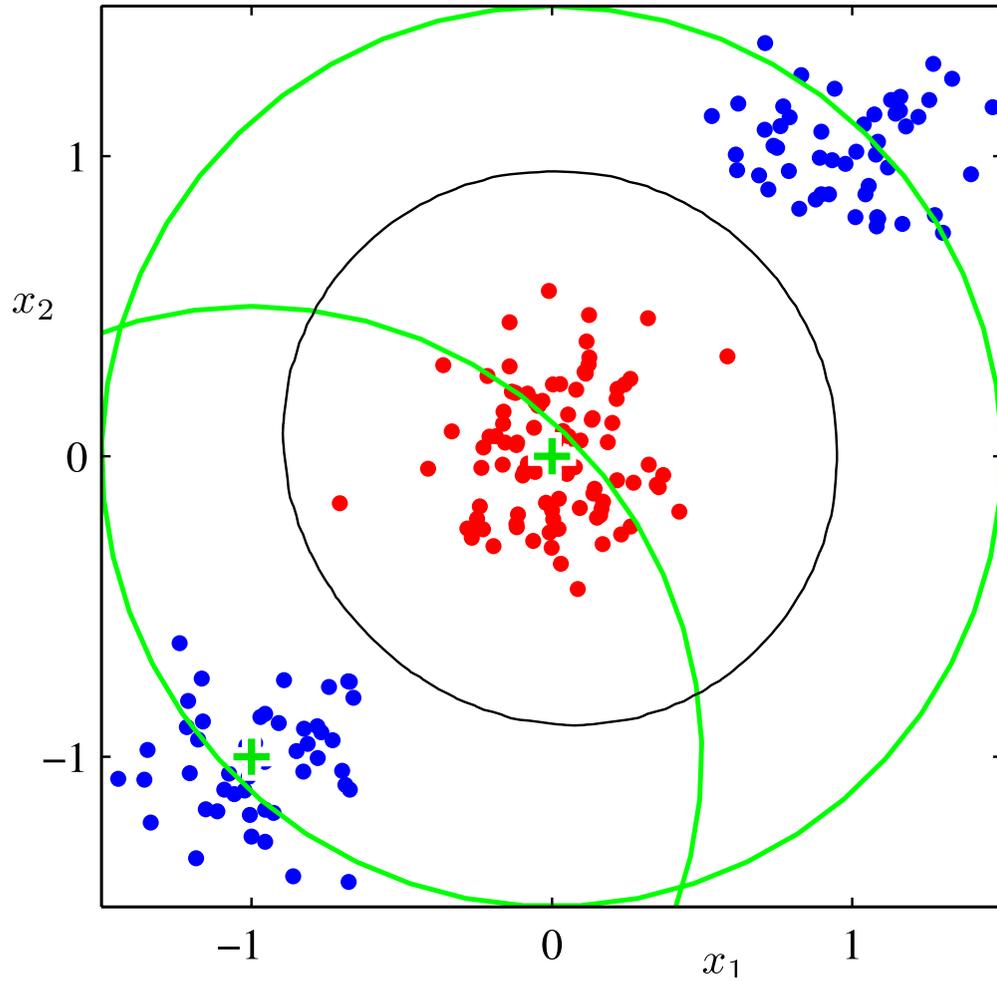


Figure 4.12 of Pattern Recognition and Machine Learning.

Multiclass Logistic Regression

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

- Generate multiple outputs, notably K outputs, each with its own set of weights, so that

$$y(\mathbf{x}; \mathbf{W})_i = \mathbf{W}_i \mathbf{x}.$$

- 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

$$\sigma(x) = \text{softmax}([\mathbf{x} \ 0])_0 = \frac{e^x}{e^x + e^0} = \frac{1}{1 + e^{-x}}.$$

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

Multiclass Logistic Regression

Note that as defined, the multiclass logistic regression is overparametrized. It is possible to generate only $K - 1$ outputs and define $z_K = 0$, which is the approach used in binary logistic regression.

In this settings, analogously to binary logistic regression, we can recover the interpretation of the model outputs $\mathbf{y}(\mathbf{x}; \mathbf{W})$ (i.e., the softmax inputs) as *logits*:

$$y(\mathbf{x}; \mathbf{W})_i = \log \left(\frac{P(C_i | \mathbf{x}; \mathbf{w})}{P(C_K | \mathbf{x}; \mathbf{w})} \right).$$

However, in all our implementations, we will use weights for all K outputs.

Multiclass Logistic Regression

Using the softmax function, we naturally define that

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

We can then use MLE and train the model using stochastic gradient descent.

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

- $\mathbf{w} \leftarrow \mathbf{0}$
- until convergence (or until patience is over), process batch of N examples:
 - $\mathbf{g} \leftarrow -\frac{1}{N} \sum_i \nabla_{\mathbf{w}} \log(P(C_{t_i} | \mathbf{x}_i; \mathbf{w}))$
 - $\mathbf{w} \leftarrow \mathbf{w} - \alpha \mathbf{g}$

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 linear combination, $\mathbf{x} = \lambda\mathbf{x}_A + (1 - \lambda)\mathbf{x}_B$, and from the linearity of $\mathbf{y}(\mathbf{x}) = \mathbf{W}\mathbf{x}$ it follows that

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

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

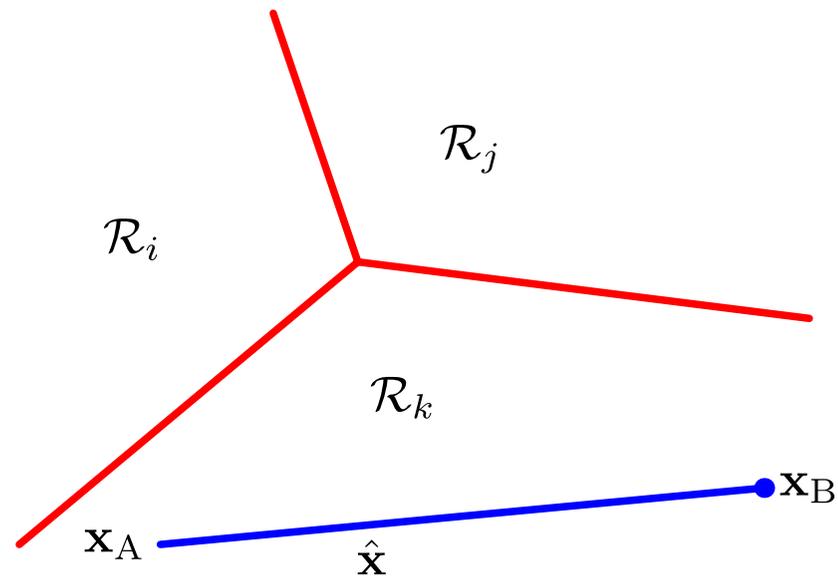


Figure 4.3 of Pattern Recognition and Machine Learning.

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.

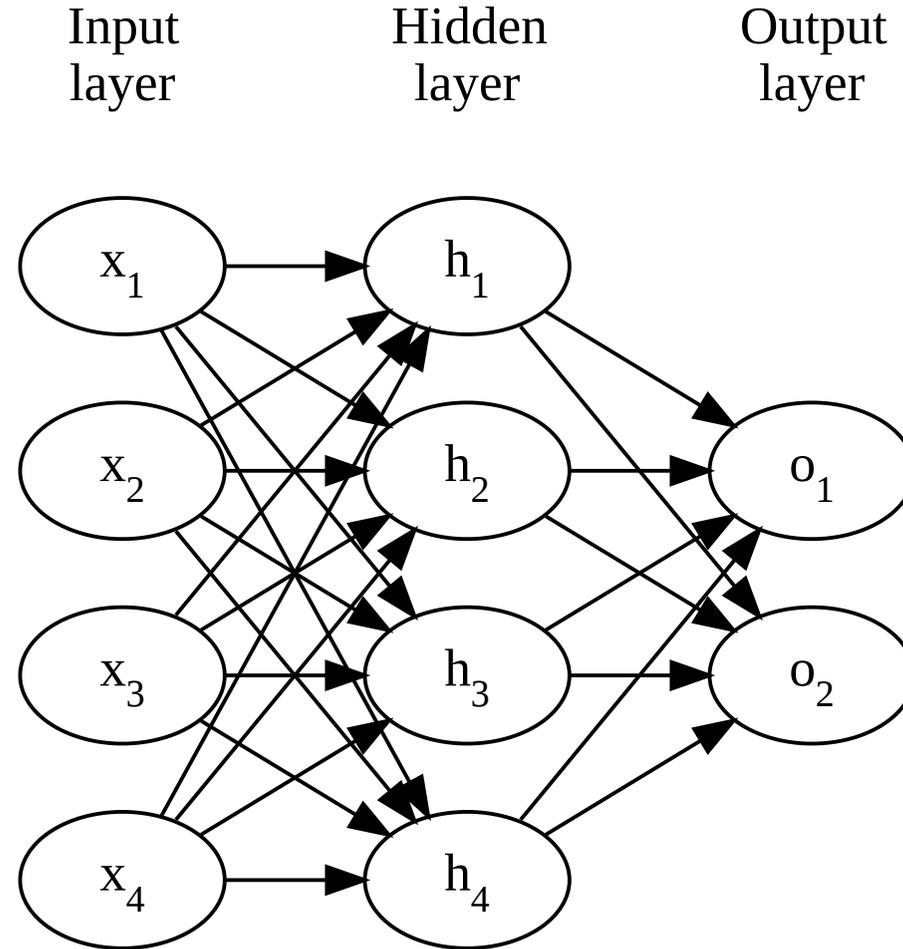
Mean Square Error as MLE

Therefore, assume our model generates a distribution

$$P(y|\mathbf{x}; \mathbf{w}) = \mathcal{N}(y; f(\mathbf{x}; \mathbf{w}), \sigma^2).$$

Now we can apply MLE and get

$$\begin{aligned} \arg \max_{\mathbf{w}} P(\mathbb{X}; \mathbf{w}) &= \arg \min_{\mathbf{w}} \sum_{i=1}^m -\log P(y_i | \mathbf{x}_i; \mathbf{w}) \\ &= - \arg \min_{\mathbf{w}} \sum_{i=1}^m \log \sqrt{\frac{1}{2\pi\sigma^2}} e^{-\frac{(y_i - f(\mathbf{x}_i; \mathbf{w}))^2}{2\sigma^2}} \\ &= - \arg \min_{\mathbf{w}} m \log(2\pi\sigma^2)^{-1/2} + \sum_{i=1}^m -\frac{(y_i - f(\mathbf{x}_i; \mathbf{w}))^2}{2\sigma^2} \\ &= \arg \min_{\mathbf{w}} \sum_{i=1}^m \frac{(y_i - f(\mathbf{x}_i; \mathbf{w}))^2}{2\sigma^2} = \arg \min_{\mathbf{w}} \sum_{i=1}^m (y_i - f(\mathbf{x}_i; \mathbf{w}))^2. \end{aligned}$$



There is a weight on each edge, and an activation function f is performed on the hidden layers, and optionally also on the output layer.

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

If the network is composed of layers, we can use matrix notation and write:

$$\mathbf{h} = f(\mathbf{W}\mathbf{x} + \mathbf{b})$$

Multilayer Perceptron and Biases

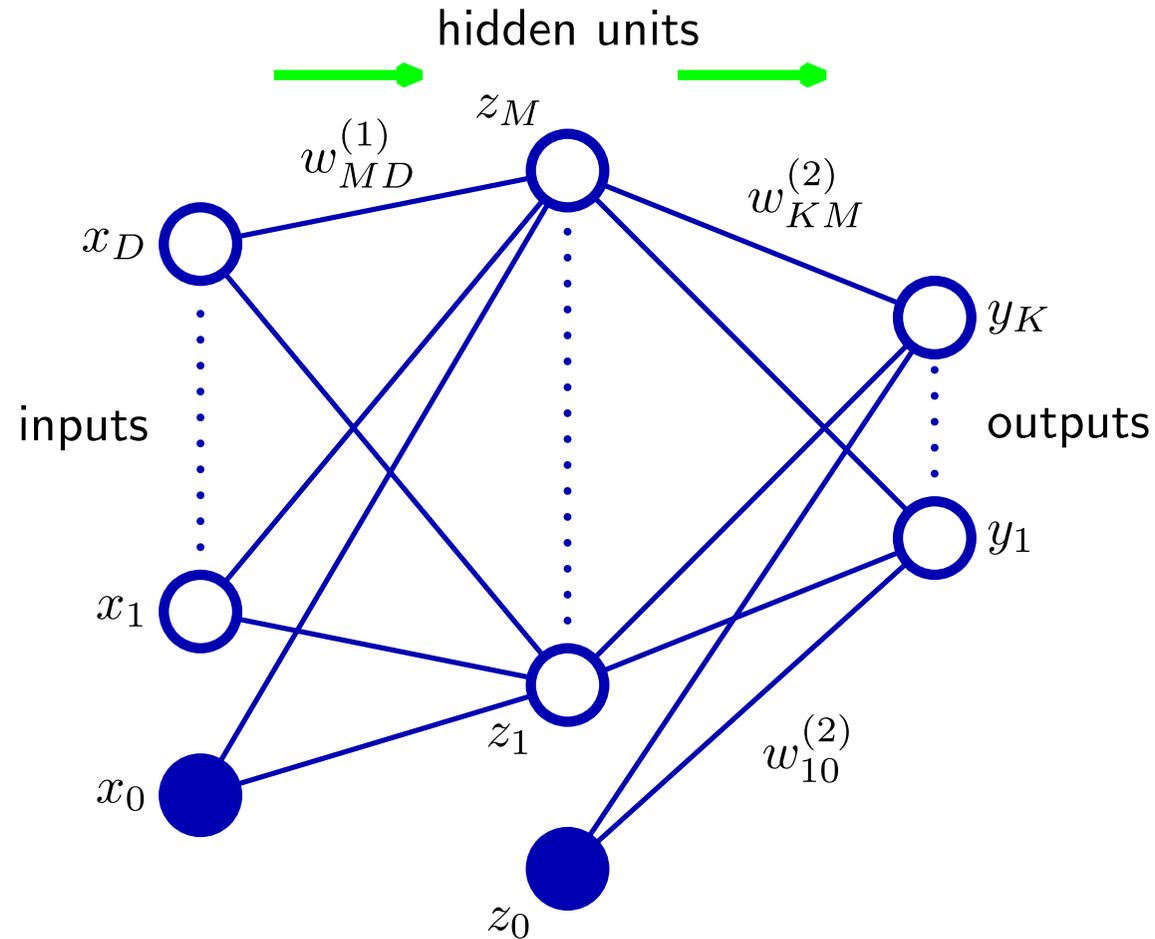


Figure 5.1 of Pattern Recognition and Machine Learning.

Output Layers

- none (linear regression if there are no hidden layers)
- sigmoid (logistic regression model if there are no hidden layers)

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

- softmax (maximum entropy model if there are no hidden layers)

$$\text{softmax}(\mathbf{x}) \propto e^{\mathbf{x}}$$

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

Hidden Layers

- none (does not help, composition of linear mapping is a linear mapping)
- σ (but works badly – 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
 - $\max(0, x)$

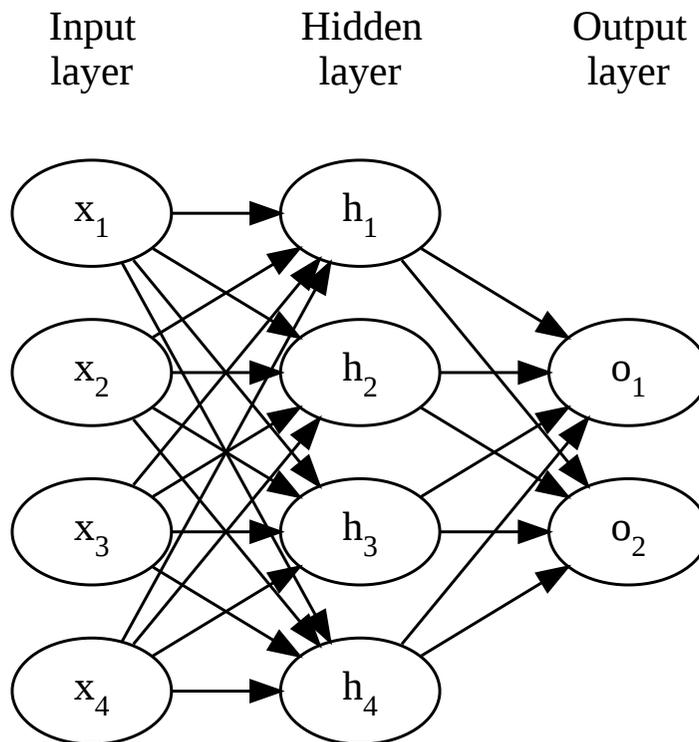
The multilayer perceptron can be trained using an SGD algorithm:

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

- $\mathbf{w} \leftarrow 0$
- until convergence (or until patience is over), process batch of N examples:
 - $\mathbf{g} \leftarrow \nabla_{\mathbf{w}} \frac{1}{N} \sum_j -\log p(y_j | \mathbf{x}_j; \mathbf{w})$
 - $\mathbf{w} \leftarrow \mathbf{w} - \alpha \mathbf{g}$

Training MLP – Computing the Derivatives

Assume a network with an input of size N_1 , then weights $\mathbf{U} \in \mathbb{R}^{N_1 \times N_2}$, hidden layer with size N_2 and activation h , weights $\mathbf{V} \in \mathbb{R}^{N_2 \times N_3}$, and finally an output layer of size N_3 with activation o .



(to be finished later)

Universal Approximation Theorem '89

Let $\varphi(x)$ be a nonconstant, bounded and nondecreasing continuous function.

(Later a proof was given also for $\varphi = \text{ReLU}$.)

Then for any $\varepsilon > 0$ and any continuous function f on $[0, 1]^m$ there exists an $N \in \mathbb{N}$, $v_i \in \mathbb{R}$, $b_i \in \mathbb{R}$ and $\mathbf{w}_i \in \mathbb{R}^m$, such that if we denote

$$F(\mathbf{x}) = \sum_{i=1}^N v_i \varphi(\mathbf{w}_i \cdot \mathbf{x} + b_i)$$

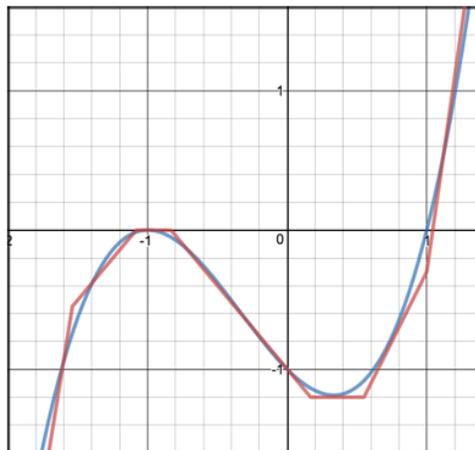
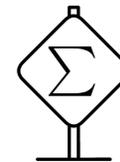
then for all $\mathbf{x} \in [0, 1]^m$

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

Universal Approximation Theorem for ReLUs

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

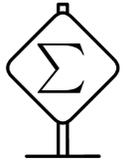
$$\begin{aligned}n_1(x) &= \text{Relu}(-5x - 7.7) \\n_2(x) &= \text{Relu}(-1.2x - 1.3) \\n_3(x) &= \text{Relu}(1.2x + 1) \\n_4(x) &= \text{Relu}(1.2x - .2) \\n_5(x) &= \text{Relu}(2x - 1.1) \\n_6(x) &= \text{Relu}(5x - 5)\end{aligned}$$

$$\begin{aligned}Z(x) &= -n_1(x) - n_2(x) - n_3(x) \\&\quad + n_4(x) + n_5(x) + n_6(x)\end{aligned}$$

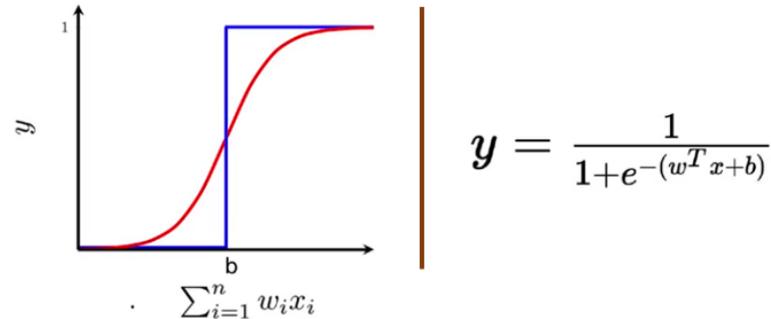
- However, we can create a sequence of k linear segments as a sum of k ReLU units – on every endpoint a new ReLU starts (i.e., the input ReLU value is zero at the endpoint), with a tangent which is the difference between the target target and the tangent of the approximation until this point.

Universal Approximation Theorem for Squashes

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

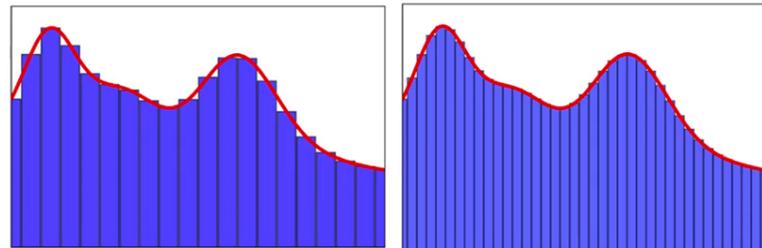


- We can prove φ can be arbitrarily close to a hard threshold by compressing it horizontally.



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

Lagrange Multipliers – Equality Constraints

Given a function $J(\mathbf{x})$, we can find a maximum with respect to a vector $\mathbf{x} \in \mathbb{R}^d$, by investigating the critical points $\nabla_{\mathbf{x}} J(\mathbf{x}) = 0$.

Consider now finding maximum subject to a constraint $g(\mathbf{x}) = 0$.

- Note that $\nabla_{\mathbf{x}} g(\mathbf{x})$ is orthogonal to the surface of the constraint, because if \mathbf{x} and a nearby point $\mathbf{x} + \boldsymbol{\varepsilon}$ lie on the surface, from the Taylor expansion $g(\mathbf{x} + \boldsymbol{\varepsilon}) \approx g(\mathbf{x}) + \boldsymbol{\varepsilon}^T \nabla_{\mathbf{x}} g(\mathbf{x})$ we get $\boldsymbol{\varepsilon}^T \nabla_{\mathbf{x}} g(\mathbf{x}) \approx 0$.
- In the sought maximum, $\nabla_{\mathbf{x}} f(\mathbf{x})$ must also be orthogonal to the constraint surface (or else moving in the direction of the derivative would increase the value).
- Therefore, there must exist λ such that $\nabla_{\mathbf{x}} f + \lambda \nabla_{\mathbf{x}} g = 0$.

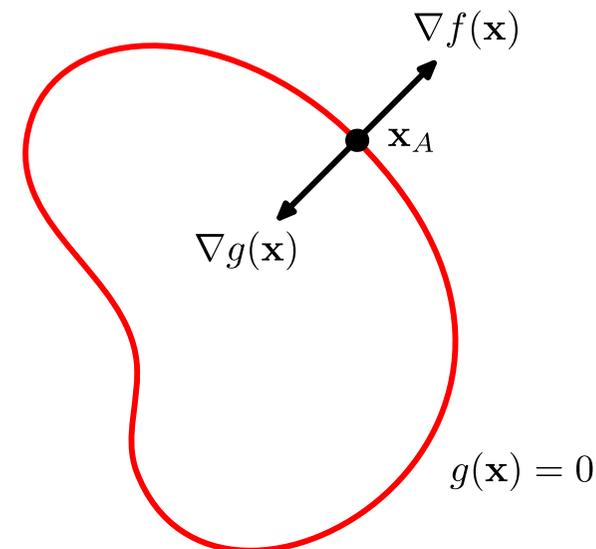


Figure E.1 of Pattern Recognition and Machine Learning.

Lagrange Multipliers – Equality Constraints

We therefore introduce the *Lagrangian function*

$$L(\mathbf{x}, \lambda) \stackrel{\text{def}}{=} f(\mathbf{x}) + \lambda g(\mathbf{x}).$$

We can then find the maximum under the constraint by inspecting critical points of $L(\mathbf{x}, \lambda)$ with respect to both \mathbf{x} and λ :

- $\frac{\partial L}{\partial \lambda} = 0$ leads to $g(\mathbf{x}) = 0$;
- $\frac{\partial L}{\partial \mathbf{x}} = 0$ is the previously derived $\nabla_{\mathbf{x}} f + \lambda \nabla_{\mathbf{x}} g = 0$.

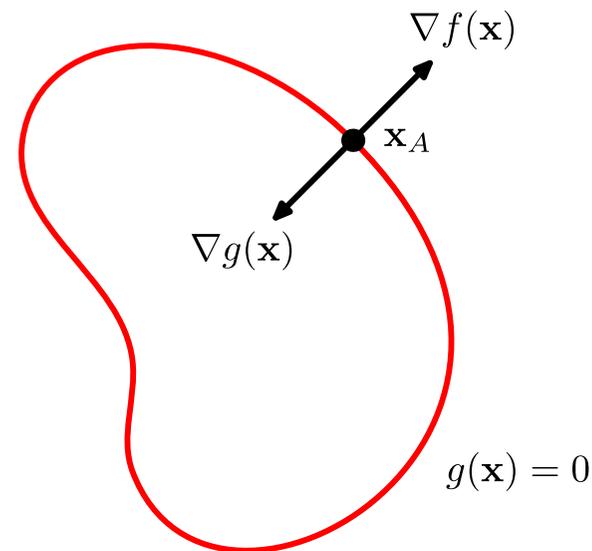
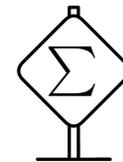


Figure E.1 of *Pattern Recognition and Machine Learning*.

Many optimization techniques depend on minimizing a function $J(\mathbf{w})$ with respect to a vector $\mathbf{w} \in \mathbb{R}^d$, by investigating the critical points $\nabla_{\mathbf{w}} J(\mathbf{w}) = 0$.



A function of a function, $J[f]$, is known as a **functional**, for example entropy $H[\cdot]$.

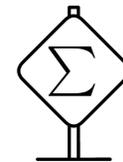
Similarly to partial derivatives, we can take **functional derivatives** of a functional $J[f]$ with respect to individual values $f(\mathbf{x})$ for all points \mathbf{x} . The functional derivative of J with respect to a function f in a point \mathbf{x} is denoted as

$$\frac{\partial}{\partial f(\mathbf{x})} J.$$

For this class, we will use only the following theorem, which states that for all differentiable functions f and differentiable functions $g(y = f(\mathbf{x}), \mathbf{x})$ with continuous derivatives, it holds that

$$\frac{\partial}{\partial f(\mathbf{x})} \int g(f(\mathbf{x}), \mathbf{x}) d\mathbf{x} = \frac{\partial}{\partial y} g(y, \mathbf{x}).$$

An intuitive view is to think about $f(\mathbf{x})$ as a vector of uncountably many elements (for every value \mathbf{x}). In this interpretation the result is analogous to computing partial derivatives of a vector $\mathbf{w} \in \mathbb{R}^d$:



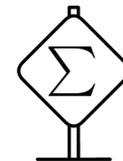
$$\frac{\partial}{\partial w_i} \sum_j g(w_j, \mathbf{x}) = \frac{\partial}{\partial w_i} g(w_i, \mathbf{x}).$$

$$\frac{\partial}{\partial f(\mathbf{x})} \int g(f(\mathbf{x}), \mathbf{x}) d\mathbf{x} = \frac{\partial}{\partial y} g(y, \mathbf{x}).$$

Function with Maximum Entropy

What distribution over \mathbb{R} maximizes entropy $H[p] = -\mathbb{E}_x \log p(x)$?

For continuous values, the entropy is an integral $H[p] = -\int p(x) \log p(x) dx$.



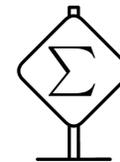
We cannot just maximize H with respect to a function p , because:

- the result might not be a probability distribution – we need to add a constraint that $\int p(x) dx = 1$;
- the problem is unspecified because a distribution can be shifted without changing entropy – we add a constraint $\mathbb{E}[x] = \mu$;
- because entropy increases as variance increases, we ask which distribution with a *fixed* variance σ^2 has maximum entropy – adding a constraint $\text{Var}(x) = \sigma^2$.

Function with Maximum Entropy

Lagrangian of all the constraints and the entropy function is

$$L(p; \mu, \sigma^2) = \lambda_1 \left(\int p(x) dx - 1 \right) + \lambda_2 (\mathbb{E}[x] - \mu) + \lambda_3 (\text{Var}(x) - \sigma^2) + H[p].$$



By expanding all definitions to integrals, we get

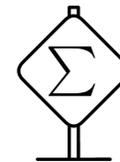
$$L(p; \mu, \sigma^2) = \int \left(\lambda_1 p(x) + \lambda_2 p(x)x + \lambda_3 p(x)(x - \mu)^2 - p(x) \log p(x) \right) dx - \lambda_1 - \mu\lambda_2 - \sigma^2\lambda_3.$$

The functional derivative of L is:

$$\frac{\partial}{\partial p(x)} L(p; \mu, \sigma^2) = \lambda_1 + \lambda_2 x + \lambda_3 (x - \mu)^2 - 1 - \log p(x) = 0.$$

Rearrangint the functional derivative of L :

$$\frac{\partial}{\partial p(x)} L(p; \mu, \sigma^2) = \lambda_1 + \lambda_2 x + \lambda_3 (x - \mu)^2 - 1 - \log p(x) = 0.$$



we obtain

$$p(x) = \exp \left(\lambda_1 + \lambda_2 x + \lambda_3 (x - \mu)^2 - 1 \right).$$

We can verify that setting $\lambda_1 = 1 - \log \sigma \sqrt{2\pi}$, $\lambda_2 = 0$ and $\lambda_3 = -1/(2\sigma^2)$ fulfils all the constraints, arriving at

$$p(x) = \mathcal{N}(x; \mu, \sigma^2).$$