

Training One Neuron

Deep Learning (ENGG*6600*07)

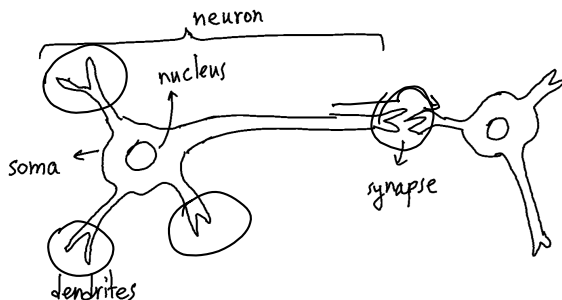
School of Engineering,
University of Guelph, ON, Canada

Course Instructor: Benyamin Ghogh
Fall 2023

Neuron and The McCulloch-Pitts Model

Biological Neuron

- **Nucleus**: The nucleus in the neuron cell
- **Soma**: Soma is the cell body of neuron containing the nucleus
- **Dendrite**: the branched protoplasmic extensions of a nerve cell that propagate the electrochemical stimulation received from other neural cells to the cell body, or soma, of the neuron
- **Synapse**: the small gap between the dendrites of connecting neurons.



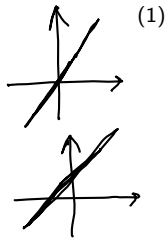
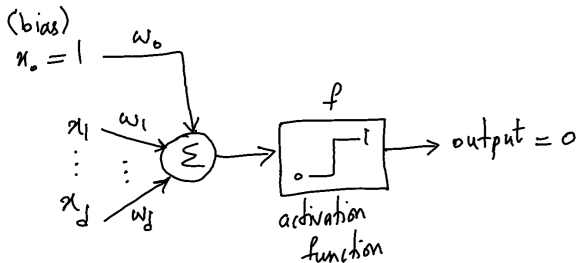
McCulloch-Pitts Model

$$affine \rightarrow a^T x + b, a \cdot b$$

- in 1943; First model of neuron was invented by McCulloch (physiologist) and Pitts (logician) [1]. It was later named the McCulloch-Pitts model.
- The model had two inputs and a single output.
- A neuron would not activate if only one of the inputs was active.
- Weights for each input were equal and fixed and the output was binary.
- Until inputs summed up to a certain threshold level, output would remain zero.
- This model is known nowadays as a logical gate, such as AND, OR, NOT, etc.
- Formula:

$$o = f\left(\sum_{j=1}^d w_j x_j + w_0\right) = f\left(\sum_{j=0}^d w_j x_j\right),$$

where $x_0 = 1$ is for bias (intercept).

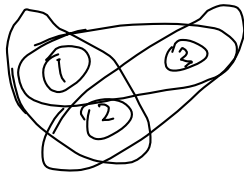


Perceptron

Hebbian learning

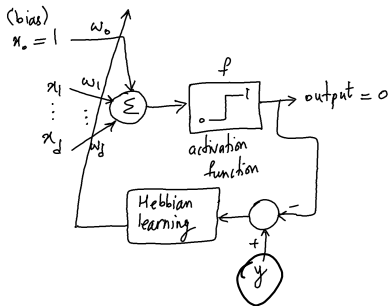
- **Hebbian theory** is a neuropsychological theory claiming that an increase in synaptic efficacy arises from a presynaptic cell's repeated and persistent stimulation of a postsynaptic cell.
- In simple words: the more two adjacent neurons are activated (fired) together, the stronger the connection becomes between them.
- **Hebbian learning** - proposed by Donald Hebb in his 1949 book "The Organization of Behavior" [2]:

Perceptron



- **Perceptron** tried to learn by neuron, as done in brain.
- Perceptron is the building block of nowadays' neural networks.
- Perceptron was implemented by Rosenblatt (physiologist) in 1958, at Cornell Aeronautical Laboratory [3].
- It was for **binary classification**. It was useful because:
 - ▶ Binary logic could do computer operations.
 - ▶ Even when we have multiple classes, we can consider pairs of classes.
- Rosenblatt randomly connected Perceptrons and changed the weights in order to achieve "learning". In later attempts, Hebbian learning [2] was used for learning in Perceptron.
- In a 1958 press conference organized by the US Navy, Rosenblatt gave a speech about the perceptron. Afterwards, New York Times reported the Perceptron to be "The embryo of an electronic computer that [the Navy] expects will be able to walk, talk, see, write, reproduce itself, and be conscious of its existence." [4]
- They thought they have solved AI!

Perceptron



- Hebbian learning [2]:

$$\begin{aligned} w_j &:= w_j + \eta(y_i - o_i)x_{ij}, \quad \forall j \in \{1, \dots, d\}, \forall i \in \{1, \dots, n\}, \\ w_0 &:= w_0 + \eta(y_i - o_i), \quad \forall i \in \{1, \dots, n\}, \end{aligned} \quad (2)$$

where $\eta > 0$ is the learning rate, $\mathbf{x}_i = [x_{i1}, \dots, x_{id}]^T \in \mathbb{R}^d$ is the i -th input data, y_i is its target label and o_i is the corresponding output of Perceptron for that input data.

- We can merge these two into one formula:

$$w_j := w_j + \eta(y_i - o_i)x_{ij}, \quad \forall j \in \{0, \dots, d\}, \forall i \in \{1, \dots, n\}, \quad (3)$$

where $x_{i0} := 1$.

Perceptron

- In batch learning of Hebbian learning:

$$\begin{aligned}w_j &:= w_j + \eta \underbrace{\sum_{i=1}^b (y_i - o_i)}_{\text{batch size}} x_{ij}, \quad \forall j \in \{1, \dots, d\}, \\w_0 &:= w_0 + \eta \underbrace{\sum_{i=1}^b (y_i - o_i)}_{\text{batch size}},\end{aligned}\tag{4}$$

where b is the batch size.

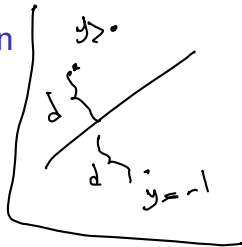
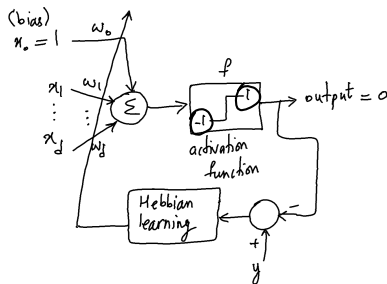
- Interpretation:

$$\star \quad \underbrace{w_j := w_j + \overbrace{\eta(y_i - o_i)}^{\text{error}} x_{ij}}_{\text{update}}, \quad \forall j \in \{1, \dots, d\}, \forall i \in \{1, \dots, n\}.$$

When the output o_i is the same as the target y_i , then we should not have any update:

$$o_i = y_i \implies \eta(y_i - o_i)x_{ij} = 0.$$

Perceptron with Signum Activation Function



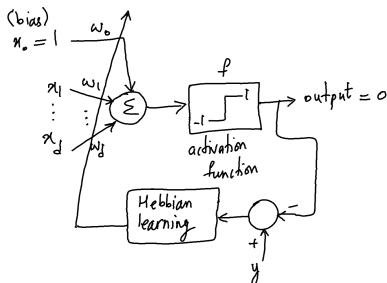
- If the activation function is the sign (signum) function, then we have:

$$\star \left[w_j := \begin{cases} w_j + 2\eta y_i x_{ij} & \text{if } o_i \neq y_i \\ w_j & \text{if } o_i = y_i, \end{cases} \right. \quad (5)$$

$$\underline{\forall j \in \{1, \dots, d\}, \forall i \in \{1, \dots, n\}.}$$

$$\left[w_0 := \begin{cases} w_0 + 2\eta y_i & \text{if } o_i \neq y_i \\ w_0 & \text{if } o_i = y_i, \end{cases} \right. \quad (6)$$

Perceptron with Signum Activation Function



- If the activation function is the sign (signum) function, in batch learning of Hebbian learning:

$$\forall j \in \{1, \dots, d\}.$$

$$\left\{ \begin{array}{l} w_j := w_j + 2\eta \sum_{i=1}^b y_i x_{ij} \mathbb{I}(\overbrace{o_i \neq y_i}), \end{array} \right. \quad (7)$$

$$\left\{ \begin{array}{l} w_0 := w_0 + 2\eta \sum_{i=1}^b y_i \mathbb{I}(o_i \neq y_i), \end{array} \right. \quad (8)$$

Perceptron with Signum Activation Function

- Proof 1:

$$w_j := w_j + \overbrace{\eta(y_i - o_i)x_{ij}}, \quad \forall j \in \{1, \dots, d\}, \forall i \in \{1, \dots, n\}.$$

When the output o_i and the target label y_i have levels ± 1 :

$$\begin{aligned} & \underline{y_i, o_i \in \{-1, 1\}}, \\ & \left. \begin{aligned} y_i = \underline{1}, o_i = \underline{1} &\Rightarrow \eta(y_i - o_i)x_{ij} = \eta(0) = \underline{0}, \\ y_i = \underline{1}, o_i = \underline{-1} &\Rightarrow \eta(y_i - o_i)x_{ij} = \eta(1 - (-1))x_{ij} = \underline{2\eta x_{ij}} = \underline{2\eta y_i x_{ij}}, \\ y_i = \underline{-1}, o_i = \underline{1} &\Rightarrow \eta(y_i - o_i)x_{ij} = \eta(-1 - 1)x_{ij} = \underline{-2\eta x_{ij}} = \underline{2\eta y_i x_{ij}}, \\ y_i = \underline{-1}, o_i = \underline{-1} &\Rightarrow \eta(y_i - o_i)x_{ij} = \eta(-1 - (-1))x_{ij} = \underline{0}, \end{aligned} \right\} \end{aligned}$$

Perceptron with Sigmoid Activation Function

Proof 2 (using gradient descent):

- We want to find a linear decision boundary where one class falls in one side and the other class falls on the other side.
- The equation of a line for linear decision boundary:

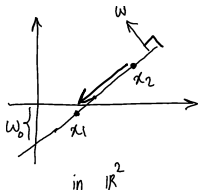
$$\boxed{\mathbf{w}^T \mathbf{x} + w_0 = 0}, \quad \star \quad (9)$$

where $\mathbf{x} \in \mathbb{R}^d$ is the data point, $\mathbf{w} \in \mathbb{R}^d$ is the normal vector of the linear line, and w_0 is the bias (intercept) of the line.

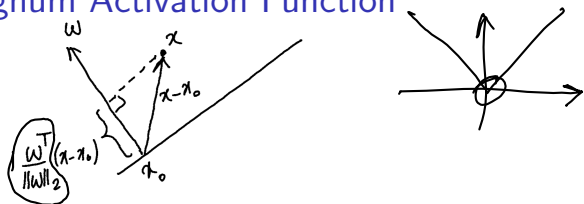
- Consider any two points \mathbf{x}_1 and \mathbf{x}_2 on the decision boundary. As the line passes through each of them, they both satisfy Eq. (9):

$$\boxed{\mathbf{w}^T \mathbf{x}_1 + w_0 = 0}, \quad \boxed{\mathbf{w}^T \mathbf{x}_2 + w_0 = 0} \Rightarrow \boxed{\mathbf{w}^T \mathbf{x}_1 + w_0 = \mathbf{w}^T \mathbf{x}_2 + w_0},$$
$$\Rightarrow \boxed{\mathbf{w}^T (\mathbf{x}_1 - \mathbf{x}_2) = 0} \Rightarrow \boxed{\mathbf{w} \perp (\mathbf{x}_1 - \mathbf{x}_2)},$$

which verifies that \mathbf{w} is the normal vector of the decision boundary.



Perceptron with Signum Activation Function



- Consider a point x_0 on the decision boundary and a point x on one of the sides of the decision boundary. Therefore:

$$\mathbf{w}^T \mathbf{x}_0 + w_0 = 0 \Rightarrow w_0 = -\mathbf{w}^T \mathbf{x}_0. \quad (10)$$

- Assume the normal vector \mathbf{w} is normalized, i.e., it has unit length. The distance of point x from the decision boundary is:

$$d = \mathbf{w}^T (\mathbf{x} - \mathbf{x}_0) = \mathbf{w}^T \mathbf{x} - \mathbf{w}^T \mathbf{x}_0 \stackrel{(10)}{=} \mathbf{w}^T \mathbf{x} + w_0. \quad (11)$$

- The distance should be non-negative so we have:

$$d = |\mathbf{w}^T \mathbf{x} + w_0|$$

However, the absolute value is non-smooth and non-differentiable. Therefore, we can multiply the distance with the target label to make it always non-negative:

$$\left[\begin{aligned} y = +1 &\Rightarrow \mathbf{w}^T \mathbf{x} + w_0 > 0, \\ y = -1 &\Rightarrow \mathbf{w}^T \mathbf{x} + w_0 < 0, \end{aligned} \right. \quad (12)$$

$$\Rightarrow d := y(\mathbf{w}^T \mathbf{x} + w_0). \quad (13)$$

Perceptron with Signum Activation Function

- A possible cost function: number of misclassified data points to be minimized. But it is discrete and hard to optimize. So, let's optimize the distance as the cost function.
- Approach 1: We want to maximize the distance of points from the decision boundary, so we use gradient ascent for maximizing the distances of points from the decision boundary¹:

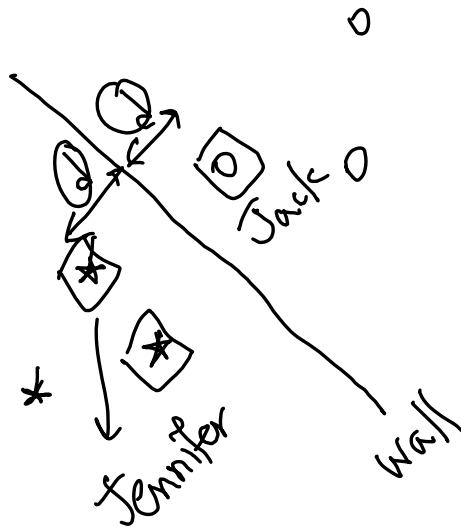
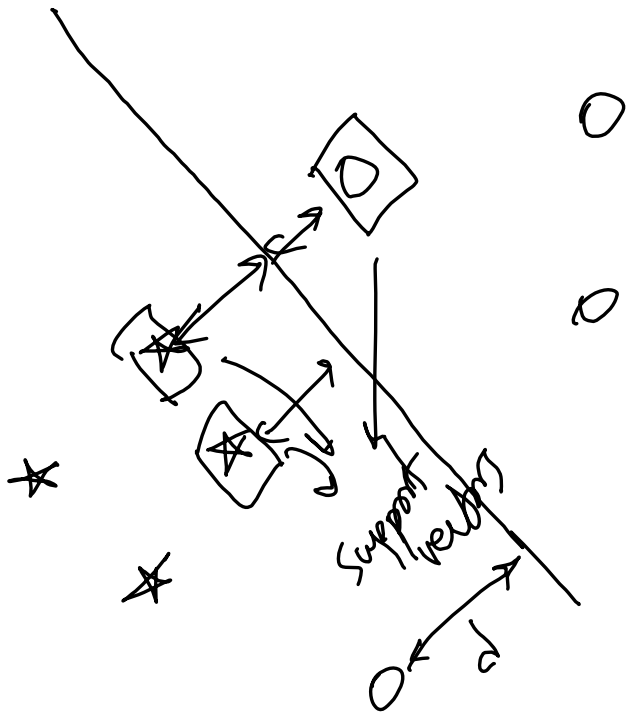
$$\begin{aligned}\text{distance} &= \sum_{i=1}^b y_i (\mathbf{w}^\top \mathbf{x}_i + w_0), \\ \rightarrow \frac{\partial \text{distance}}{\partial \mathbf{w}} &= \sum_{i=1}^b y_i \mathbf{x}_i, \quad \frac{\partial \text{distance}}{\partial w_0} = \sum_{i=1}^b y_i\end{aligned}$$

We should use gradient ascent to maximize the distances:

$$\left\{ \begin{aligned} \mathbf{w} &:= \mathbf{w} + \eta \sum_{i=1}^b y_i \mathbf{x}_{ij}, \quad \forall j \in \{1, \dots, d\}, \\ w_0 &:= w_0 + \eta \sum_{i=1}^b y_i, \end{aligned} \right. \quad (14)$$

where we can update only for $x_i \neq y_i$ cases to have Eqs. (7) and (8).

¹ Note that it is not like support vector machine which maximizes the distance of only support vectors, and not all points, from the decision boundary.



Perceptron with Signum Activation Function

- Approach 2: We take the distances of misclassified data points as the cost function. According to Eq. (12), the distance in Eq. (13) is for correctly classified data points. So, we should multiply distance by -1 to have the distances of misclassified points:

$$\begin{aligned}\underline{\text{error}} &= \sum_{i=1}^b \ominus y_i (\mathbf{w}^\top \mathbf{x}_i + w_0) = \ominus \sum_{i=1}^b y_i (\mathbf{w}^\top \mathbf{x}_i + w_0), \\ \frac{\partial \text{error}}{\partial \mathbf{w}} &= \ominus \underbrace{\sum_{i=1}^b y_i \mathbf{x}_i}, \quad \frac{\partial \text{error}}{\partial w_0} = \ominus \underbrace{\sum_{i=1}^b y_i}.\end{aligned}$$

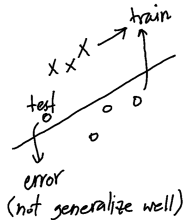
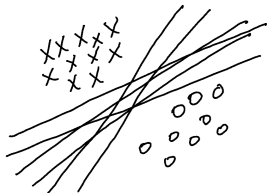
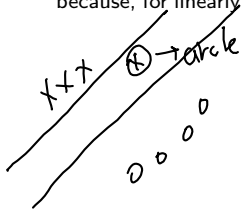
We should use gradient descent to minimize the error:

$$\left\{ \begin{array}{l} \mathbf{w} := \mathbf{w} + \eta \sum_{i=1}^b y_i \mathbf{x}_{ij}, \quad \forall j \in \{1, \dots, d\}, \\ w_0 = w_0 + \eta \sum_{i=1}^b y_i, \end{array} \right. \quad (15)$$

where we can update only for $o_i \neq y_i$ cases to have Eqs. (7) and (8).

Problems with Perceptron

- Perceptron is for binary classification.
- In 1969, Minsky and Papert published a book titled "Perceptrons" [5] and showed that Perceptron can only solve linearly separable problems. For example, they showed that it cannot classify XOR classes, which is a nonlinear classification problem.
- Therefore, researchers lost interest in Perceptron and artificial neural networks!
- Researchers guessed that they should have multilayer Perceptrons but they did not know how to train multilayer Perceptrons.
- Also, Perceptron cannot generalize well enough because, for linearly separable classes, it finds one of the many possible decision boundaries.



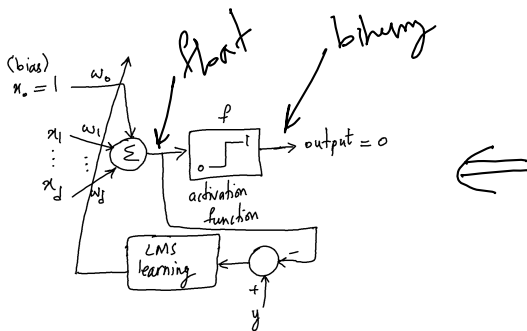
- Because of this problem, two things were developed:
 - ▶ ADALINE which generalizes better.
 - ▶ Support Vector Machines (SVM) which found the best decision boundary by optimization of the distances of the support vectors from the decision boundary.

ADALINE

ADALINE

- In 1960: Widrow and his student Hoff, at Stanford University, proposed a method, named ADALINE, for adjusting weights [6, 7].
- In 1960: articles claimed that robots can think!
- It has more generalization compared to Perceptron.
- ADALINE minimizes least mean squares (LMS) error using gradient descent. This training method is referred to as LMS algorithm or Widrow-Hoff learning rule.

ADALINE



- In ADALINE, the target label is compared with the output before activation function. This is while Perceptron compares the target label with the output after activation function.
- LMS error:

$$e = \left(\frac{1}{2} \sum_{i=1}^b \left(y_i - \left(\sum_{j=1}^d w_j x_{ij} + w_0 \right) \right)^2 \right), \quad (16)$$

where b is the batch size, $\mathbf{x}_i = [x_{i1}, \dots, x_{id}]^T \in \mathbb{R}^d$ is the i -th input data, and y_i is its target label.

ADALINE

- We had the LMS error:

$$e = \frac{1}{2} \sum_{i=1}^b \left(y_i - \left(\sum_{j=1}^d w_j x_{ij} + w_0 \right) \right)^2 \quad (17)$$

Handwritten annotations:
 - "Scalars" written above the summation term.
 - A bracket under the term $\left(\sum_{j=1}^d w_j x_{ij} + w_0 \right)$ with an arrow pointing to the x_{ij} term in the equation above.
 - A star $*$ below the summation term.

- The gradients:

$$\star \quad \frac{\partial e}{\partial w_j} = - \sum_{i=1}^b \left(y_i - \left(\sum_{j=1}^d w_j x_{ij} + w_0 \right) \right) x_{ij}, \quad \forall j \in \{1, \dots, d\},$$

$$\frac{\partial e}{\partial w_0} = - \sum_{i=1}^b \left(y_i - \left(\sum_{j=1}^d w_j x_{ij} + w_0 \right) \right).$$

Handwritten annotations:
 - A star $*$ to the left of the first equation.
 - A bracket under the term $\forall j \in \{1, \dots, d\}$ in the first equation.

- The gradient descent updates:

$$\begin{cases} w_j := w_j - \eta \frac{\partial e}{\partial w_j}, \\ w_0 := w_0 - \eta \frac{\partial e}{\partial w_0}, \end{cases} \quad \forall j \in \{1, \dots, d\}, \quad (19)$$

where $\eta > 0$ is the learning rate.

$$C = \sum_{j=1}^d w_j x_{ij} = \underbrace{w_1 x_{i1} + w_2 x_{i2} + \dots + w_d x_{id}}$$

$$\frac{\partial C}{\partial w_j} \rightarrow \frac{\partial C}{\partial w_1} = x_{i1} \quad \frac{\partial C}{\partial w_2} = x_{i2}$$

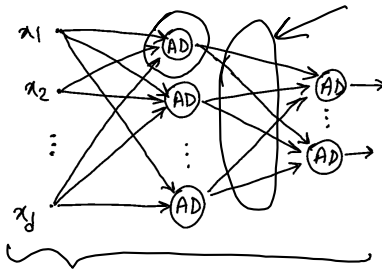
$$\frac{\partial C}{\partial w_j} = x_{ij}$$

MADALINE

MADALINE

- At Stanford university, Widrow and his students stacked several ADALINE neurons to be able to have nonlinear classification. They proposed **MADALINE** (Many ADALINEs).
- MADALINE Rule 1 (MRI)**: proposed in 1962 [8] and could not adapt the weights of the hidden-output layer.
- MADALINE Rule 2 (MRII)**: proposed in 1988 [9] and improved MRI to be able to also train the weights of the hidden-output layer.
- MADALINE Rule 3 (MRIII)**: proposed in 1990 [10] and changed signum activation function to sigmoid function for having float outputs rather than merely binary outputs.

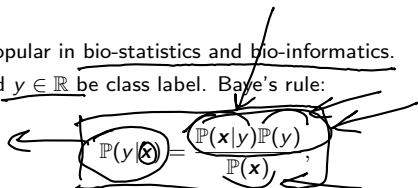
inspired
by
fuzzy logic



Logistic Regression

Logistic Regression

- Logistic regression is popular in bio-statistics and bio-informatics.
- Let $\mathbf{x} \in \mathbb{R}^d$ be data and $y \in \mathbb{R}$ be class label. Baye's rule:



A hand-drawn diagram of the equation $\mathbb{P}(y|\mathbf{x}) = \frac{\mathbb{P}(\mathbf{x}|y)\mathbb{P}(y)}{\mathbb{P}(\mathbf{x})}$. The equation is enclosed in a hand-drawn rectangle. Arrows point from the text above to specific parts of the equation: one arrow points to $\mathbb{P}(\mathbf{x}|y)$, another to $\mathbb{P}(y)$, and a third to $\mathbb{P}(\mathbf{x})$. The term $\mathbb{P}(y|\mathbf{x})$ is circled, and the entire equation is also circled.

$$\mathbb{P}(y|\mathbf{x}) = \frac{\mathbb{P}(\mathbf{x}|y)\mathbb{P}(y)}{\mathbb{P}(\mathbf{x})}, \quad (20)$$

where $\mathbb{P}(y|\mathbf{x})$ and $\mathbb{P}(\mathbf{x}|y)$ are the posterior and likelihood, respectively, and $\mathbb{P}(\mathbf{x})$ and $\mathbb{P}(y)$ are the priors.

- In contrast to Linear Discriminant Analysis (LDA), logistic regression works on the posterior $\mathbb{P}(y|\mathbf{x})$ directly rather than working on likelihood $\mathbb{P}(\mathbf{x}|y)$ and prior $\mathbb{P}(y)$.

Logistic Regression



- Logistic regression is a binary classifier where it assigns probability between zero and one for belonging to one of the classes.
- The logistic function, used in logistic regression, was initially proposed in 1845 for modeling the population growth [11]. It was further improved in the 20th century [12]. See [13] for the history of logistic regression.
- It considers the classification problem as a regression problem where it regresses (predicts) the probability of belonging to a class. It first considers a linear regression $\beta^T \mathbf{x} + \beta_0$. However, in order to not have the bias, it assumes that \mathbf{x} is $d + 1$ dimensional with an additional element of 1 for bias, i.e., $\mathbf{x} = [x_1, \dots, x_d, 1]^T$. The $\beta \in \mathbb{R}^{d+1}$ is the learnable parameter of the logistic regression model. As a result, the linear regression becomes $\beta^T \mathbf{x}$.
- However, there is no bound on this regression while logistic regression desires the output to be in the range $[0, 1]$ to behave like a probability. Therefore, Logistic regression models the posterior using a logistic function, also called the sigmoid function, to make this regression between zero and one.

$$\beta^T \mathbf{x} + \beta_0 \rightarrow \beta^T \mathbf{x} \rightarrow x_0 = 1$$

\downarrow \downarrow
 $d+1$ $d+1$

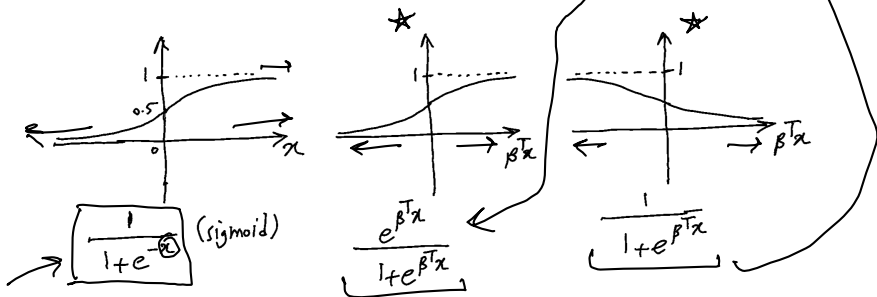
Logistic Regression

- Assume we have two classes $y \in \{0, 1\}$.
- Logistic regression models the posterior using a logistic function, also called the **sigmoid function**:

$$\star \left[\mathbb{P}(y = 1 | X = x) = \frac{e^{\beta^\top x}}{1 + e^{\beta^\top x}} \right] \quad \leftarrow \text{linear regression} \quad (21)$$

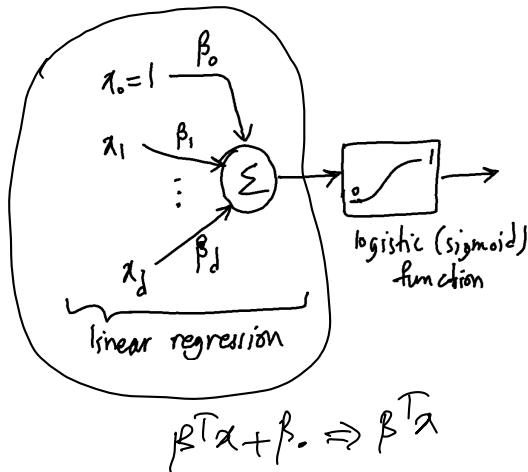
$$\star \left[\mathbb{P}(y = 0 | X = x) = 1 - \mathbb{P}(y = 1 | X = x) = \frac{1}{1 + e^{\beta^\top x}} \right] \quad (22)$$

where $\beta \in \mathbb{R}^d$ is the learnable parameter of the logistic regression model.



Logistic Regression as a Neural Network

- Logistic regression can be seen as a neural network with one neuron where the activation function is the nonlinear sigmoid (logistic) function.



Logistic Regression

- Consider n data points $\{(x_i, y_i)\}_{i=1}^n$ in the dataset. Assuming that they are independent and identically distributed (i.i.d), the posterior over all data points is:

$$\mathbb{P}(y|X) = \prod_{i=1}^n \left(\underbrace{\mathbb{P}(y_i = 1|X = x_i)}_{\text{prob of } y_i=1} \underbrace{\mathbb{I}(y_i = 1)}_{\text{indicator}} + \underbrace{\mathbb{P}(y_i = 0|X = x_i)}_{\text{prob of } y_i=0} \underbrace{\mathbb{I}(y_i = 0)}_{\text{indicator}} \right), \quad (23)$$

where $\mathbb{I}(\cdot)$ is the indicator function which is one if its condition is satisfied and is zero otherwise.

- As the labels are either zero or one, i.e., $y_i \in \{0, 1\}$, this equation can be restated as:

$$\mathbb{P}(y|X) = \prod_{i=1}^n \left(\underbrace{\mathbb{P}(y_i = 1|X = x_i)}_{\text{prob of } y_i=1} \right)^{y_i} \underbrace{\left(\mathbb{P}(y_i = 0|X = x_i) \right)^{1-y_i}}_{\text{prob of } y_i=0} \quad (24)$$

- Substituting Eqs. (21) and (22) in this equation gives:

$$\mathbb{P}(y|X) = \prod_{i=1}^n \left(\underbrace{\frac{e^{\beta^\top x_i}}{1 + e^{\beta^\top x_i}}}_{\text{prob of } y_i=1} \right)^{y_i} \underbrace{\left(\frac{1}{1 + e^{\beta^\top x_i}} \right)^{1-y_i}}_{\text{prob of } y_i=0}. \quad (25)$$

Logistic Regression

$$L(\beta) \quad \ell(\beta)$$

$$\begin{aligned} \log\left(\frac{a}{b}\right) &= \log a - \log b \\ \log(ab) &= \log a + \log b \\ \log \frac{1}{a} &= -\log a \end{aligned}$$

- The log posterior is:

$$\begin{aligned} \ell(\beta) &:= \mathbb{P}(y|X=x) = \log \prod_{i=1}^n \left(\frac{e^{\beta^\top x_i}}{1 + e^{\beta^\top x_i}} \right)^{y_i} \left(\frac{1}{1 + e^{\beta^\top x_i}} \right)^{1-y_i} \\ &= \sum_{i=1}^n \left(\log \left(\frac{e^{\beta^\top x_i}}{1 + e^{\beta^\top x_i}} \right)^{y_i} + \log \left(\frac{1}{1 + e^{\beta^\top x_i}} \right)^{1-y_i} \right) \\ &= \sum_{i=1}^n \left(y_i \log(e^{\beta^\top x_i}) - y_i \log(1 + e^{\beta^\top x_i}) - (1 - y_i) \log(1 + e^{\beta^\top x_i}) \right) \\ &= \sum_{i=1}^n \left(y_i \beta^\top x_i - y_i \log(1 + e^{\beta^\top x_i}) - \log(1 + e^{\beta^\top x_i}) + y_i \log(1 + e^{\beta^\top x_i}) \right) \\ &= \sum_{i=1}^n \left(y_i \beta^\top x_i - \log(1 + e^{\beta^\top x_i}) \right). \end{aligned}$$

Logistic Regression

- The log posterior is:

$$\star \ell(\beta) = \sum_{i=1}^n (y_i \beta^\top x_i - \log(1 + e^{\beta^\top x_i})).$$

- Newton's method can be used to find the optimum β . The first derivative, or the gradient, is:

$$\frac{\partial \ell(\beta)}{\partial \beta} = \sum_{i=1}^n (y_i x_i - \frac{1}{1 + e^{\beta^\top x_i}} e^{\beta^\top x_i} x_i) = \sum_{i=1}^n (y_i - \frac{e^{\beta^\top x_i}}{1 + e^{\beta^\top x_i}}) x_i. \quad (26)$$

Its transpose is:

$$\frac{\partial \ell(\beta)}{\partial \beta^\top} = \sum_{i=1}^n (y_i - \frac{e^{\beta^\top x_i}}{1 + e^{\beta^\top x_i}}) x_i^\top.$$

Logistic Regression

- The second derivative is:

$$\begin{aligned} \frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta^\top} &= \frac{\partial}{\partial \beta} \left(\frac{\partial \ell(\beta)}{\partial \beta^\top} \right) = \frac{\partial}{\partial \beta} \left(\sum_{i=1}^n \left(y_i - \frac{e^{\beta^\top x_i}}{1 + e^{\beta^\top x_i}} \right) x_i^\top \right) \\ &= \sum_{i=1}^n \left(\ominus \frac{\partial}{\partial \beta} \left(\frac{e^{\beta^\top x_i}}{1 + e^{\beta^\top x_i}} \right) \right) x_i^\top. \end{aligned}$$

- We define:

$$\mathbb{P}(x_i | \beta) := \frac{e^{\beta^\top x_i}}{1 + e^{\beta^\top x_i}}. \quad (27)$$

Therefore:

$$\frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta^\top} = \ominus \sum_{i=1}^n \left(\frac{\partial}{\partial \beta} (\mathbb{P}(x_i | \beta)) \right) x_i^\top. \quad (28)$$

Logistic Regression

- We have:

$$\begin{aligned}
 \frac{\partial}{\partial \beta} (\mathbb{P}(\mathbf{x}_i | \beta)) &= \frac{\partial}{\partial \beta} \left(\frac{e^{\beta^\top \mathbf{x}_i}}{1 + e^{\beta^\top \mathbf{x}_i}} \right) \\
 &= \frac{1}{(1 + e^{\beta^\top \mathbf{x}_i})^2} \left(e^{\beta^\top \mathbf{x}_i} \mathbf{x}_i (1 + e^{\beta^\top \mathbf{x}_i}) - e^{\beta^\top \mathbf{x}_i} (e^{\beta^\top \mathbf{x}_i} \mathbf{x}_i) \right) \\
 &= \frac{e^{\beta^\top \mathbf{x}_i}}{(1 + e^{\beta^\top \mathbf{x}_i})^2} (1 + e^{\beta^\top \mathbf{x}_i} - e^{\beta^\top \mathbf{x}_i}) \mathbf{x}_i = \frac{e^{\beta^\top \mathbf{x}_i}}{(1 + e^{\beta^\top \mathbf{x}_i})^2} \mathbf{x}_i \\
 &= \underbrace{\frac{e^{\beta^\top \mathbf{x}_i}}{(1 + e^{\beta^\top \mathbf{x}_i})}}_{\mathbb{P}(\mathbf{x}_i | \beta)} \underbrace{\frac{1}{(1 + e^{\beta^\top \mathbf{x}_i})} \mathbf{x}_i}_{(1 - \mathbb{P}(\mathbf{x}_i | \beta)) \mathbf{x}_i} \stackrel{(27)}{=} \mathbb{P}(\mathbf{x}_i | \beta) (1 - \mathbb{P}(\mathbf{x}_i | \beta)) \mathbf{x}_i
 \end{aligned}$$

- Substituting it in Eq. (28) gives the second derivative, i.e., the Hessian matrix:

$$\frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta^\top} = - \sum_{i=1}^n \left(\mathbb{P}(\mathbf{x}_i | \beta) (1 - \mathbb{P}(\mathbf{x}_i | \beta)) \mathbf{x}_i \mathbf{x}_i^\top \right) \quad (29)$$

Logistic Regression

- It is possible to write the Newton's method in matrix form. We define:

$$\begin{aligned}\mathbb{R}^{(d+1) \times n} \ni \mathbf{X} &:= \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \end{bmatrix}, \\ \mathbb{R}^{n \times n} \ni \mathbf{W} &:= \text{diag} \left(\mathbb{P}(x_i|\beta)(1 - \mathbb{P}(x_i|\beta)) \right), \\ \mathbb{R}^n \ni \mathbf{y} &:= [y_1, \dots, y_n]^\top, \\ \mathbb{R}^n \ni \mathbf{p} &:= \left[\frac{e^{\beta^\top x_1}}{1 + e^{\beta^\top x_1}}, \dots, \frac{e^{\beta^\top x_n}}{1 + e^{\beta^\top x_n}} \right]^\top.\end{aligned}$$

- The Eqs. (26) and (29) can be restated as:

$$\mathbb{R}^{(d+1)} \ni \frac{\partial \ell(\beta)}{\partial \beta} = \mathbf{X}(\mathbf{y} - \mathbf{p}), \quad (30)$$

$$\mathbb{R}^{(d+1) \times (d+1)} \ni \frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta^\top} = -\mathbf{X} \mathbf{W} \mathbf{X}^\top. \quad (31)$$

- Using Newton's method for maximization of the log posterior is:

$$\begin{aligned}\beta^{(\tau+1)} &:= \beta^{(\tau)} + \left(\frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta^\top} \right)^{-1} \frac{\partial \ell(\beta)}{\partial \beta} \Rightarrow \\ \beta^{(\tau+1)} &:= \beta^{(\tau)} - (\mathbf{X} \mathbf{W} \mathbf{X}^\top)^{-1} \mathbf{X}(\mathbf{y} - \mathbf{p}),\end{aligned} \quad (32)$$

where τ is the iteration index. It is repeated until convergence of β .

Logistic Regression

- In the test phase, the class of a point \mathbf{x} is determined as:

$$\star y = \begin{cases} 1 & \text{if } \frac{e^{\beta^\top \mathbf{x}}}{1 + e^{\beta^\top \mathbf{x}}} \geq 0.5, \\ 0 & \text{Otherwise.} \end{cases} \quad (33)$$

- Comparison to LDA:

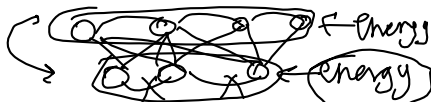
- ▶ Logistic regression estimates $(d + 1)$ parameters in β , but LDA estimates many more parameters:
 - ★ prior of each class: 1. We have two classes: $2 \times 1 = 2$.
 - ★ mean of each class: d . We have two classes: $2 \times d = 2d$.
 - ★ covariance matrix of each class: $d(d + 1)/2$. We have two classes: $2 \times (d(d + 1)/2) = d(d + 1)$.
 - ★ so, in total: $2 + 2d + d(d + 1) = \underline{d^2 + 2d + 2}$.
- ▶ LDA assumes the distribution of each class is Gaussian which may not be true. However, logistic regression does not assume anything about the distribution of data.

Other History

Other History

- in 1969: Arthur E. Bryson and Yu-Chi Ho described backpropagation as a multi-stage dynamic system optimization method [14, 15].
- Starting 1969, people started inventing and re-inventing backpropagation algorithm for training multilayer Perceptron.
- in 1972: Stephen Grossberg proposed networks capable of learning XOR function.
- in 1986: the main and succesful backpropagation was proposed by David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams [16].
- in 1980's: successful era of neural networks.
- Kernel support vector machines [17] resulted in the winter of neural networks in the last years of previous century until around 2006.

Other History



- Hinton et. al. had proposed Boltzmann Machine (BM) and Restricted Boltzmann Machine (RBM) in 1983 and 1985 [18, 19].
- During the winter of neural networks, Hinton tried to save neural networks from being forgotten in the history of machine learning. So, he returned to his previously proposed RBM and proposed a learning method for RBM with the help of some other researchers including Max Welling [20, 21].
- They proposed training the weights of BM and RBM using maximum likelihood estimation. BM and RBM can be seen as generative models where new values for neurons can be generated using Gibbs sampling [22].
- Hinton noticed RBM because he knew that the set of weights between every two layers of a neural network is an RBM. It was in the year 2006 [23, 24] that he thought it is possible to train a network in a greedy way² [25] where the weights of every layer of network is trained using RBM training.
- This stack of RBM models with a greedy algorithm for training was named Deep Belief Network (DBN) [24, 26]. DBN allowed the networks to become deep by preparing a good initialization of weights (using RBM training) for backpropagation. This good starting point for backpropagation optimization did not face the problem of vanishing gradients anymore.

²A greedy algorithm makes every decision based on the most benefit at the current step and does not care about the final outcome at the final step. This greedy approach hopes that the final step will obtain a good result by small best steps based on their current benefits.

Other History

- Since the breakthrough in 2006 [23], the winter of neural networks started to end gradually because the networks could get deep to become more nonlinear and handle more nonlinear data.
- DBN was used in different applications including speech recognition [27, 28, 29] and action recognition [30].
- Hinton was very excited about the success of RBM and was thinking that the future of neural networks belongs to DBN.
- However, two important techniques were proposed, which were the ReLU activation function (2011) [31] and the dropout technique (2014) [32]. These two regularization methods prevented overfitting [33] and resolved vanishing gradients even without RBM pre-training.
- Hence, backpropagation could be used alone if the new regularization methods were utilized. The success of neural networks was found out more [34] by its various applications, for example in image recognition [35].

Acknowledgment

- For more information on early history of artificial neural networks, see the book [36]:
“Fundamentals of neural networks: architectures, algorithms and applications”
- Some slides of this slide deck were inspired by teachings of Prof. Ali Ghodsi (at University of Waterloo, Department of Statistics), Prof. Fakhri Karray (at University of Waterloo, Department of Electrical and Computer Engineering), and Prof. Saeed Bagheri Shouraki (at Sharif University of Technology, Department of Electrical Engineering).

References

- [1] W. S. McCulloch and W. Pitts, "A logical calculus of the ideas immanent in nervous activity," *The bulletin of mathematical biophysics*, vol. 5, pp. 115–133, 1943.
- [2] D. O. Hebb, *The Organization of Behavior*.
New York: Wiley & Sons, 1949.
- [3] F. Rosenblatt, "The Perceptron – a perceiving and recognizing automaton project para," tech. rep., Report 85-460-1, Cornell Aeronautical Laboratory., 1957.
- [4] M. Olazaran, "A sociological study of the official history of the perceptrons controversy," *Social Studies of Science*, vol. 26, no. 3, pp. 611–659, 1996.
- [5] M. Minsky and S. A. Papert, *Perceptrons*.
MIT press, 1969.
- [6] B. Widrow and M. E. Hoff, "Adaptive switching circuits," tech. rep., Stanford University, California, Stanford Electronics Labs, 1960.
- [7] B. Widrow, "An adaptive "adaline" neuron using chemical "memistors"," tech. rep., No. 1553-2, Stanford Electronics Laboratories, 1960.
- [8] B. Widrow, "Generalization and information storage in networks of adaline neurons," *Self-organizing systems*, pp. 435–461, 1962.
- [9] C. R. Winter and B. Widrow, "Madaline rule ii: A training algorithm for neural networks," in *Second Annual International Conference on Neural Networks*, pp. 1–401, 1988.

References (cont.)

- [10] B. Widrow and M. A. Lehr, "30 years of adaptive neural networks: perceptron, madaline, and backpropagation," *Proceedings of the IEEE*, vol. 78, no. 9, pp. 1415–1442, 1990.
- [11] P. F. Verhulst, "Resherches mathematiques sur la loi d'accroissement de la population," *Nouveaux memoires de l'academie royale des sciences*, vol. 18, pp. 1–41, 1845.
- [12] S. H. Walker and D. B. Duncan, "Estimation of the probability of an event as a function of several independent variables," *Biometrika*, vol. 54, no. 1-2, pp. 167–179, 1967.
- [13] J. S. Cramer, "The origins of logistic regression," 2002.
- [14] A. E. Bryson Jr, W. F. Denham, and S. E. Dreyfus, "Optimal programming problems with inequality constraints," *AIAA journal*, vol. 1, no. 11, pp. 2544–2550, 1963.
- [15] A. E. Bryson and H. Yu-Chi, *Applied optimal control: optimization, estimation and control*.
CRC Press, 1969.
- [16] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning representations by back-propagating errors," *nature*, vol. 323, no. 6088, pp. 533–536, 1986.
- [17] B. E. Boser, I. M. Guyon, and V. N. Vapnik, "A training algorithm for optimal margin classifiers," in *Proceedings of the fifth annual workshop on Computational learning theory*, pp. 144–152, 1992.

References (cont.)

- [18] G. E. Hinton and T. J. Sejnowski, "Optimal perceptual inference," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, vol. 448, IEEE, 1983.
- [19] D. H. Ackley, G. E. Hinton, and T. J. Sejnowski, "A learning algorithm for Boltzmann machines," *Cognitive science*, vol. 9, no. 1, pp. 147–169, 1985.
- [20] G. E. Hinton, "Training products of experts by minimizing contrastive divergence," *Neural computation*, vol. 14, no. 8, pp. 1771–1800, 2002.
- [21] M. Welling, M. Rosen-Zvi, and G. E. Hinton, "Exponential family harmoniums with an application to information retrieval.," in *Advances in neural information processing systems*, vol. 4, pp. 1481–1488, 2004.
- [22] S. Geman and D. Geman, "Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images," *IEEE Transactions on pattern analysis and machine intelligence*, vol. PAMI-6, no. 6, pp. 721–741, 1984.
- [23] G. E. Hinton and R. R. Salakhutdinov, "Reducing the dimensionality of data with neural networks," *Science*, vol. 313, no. 5786, pp. 504–507, 2006.
- [24] G. E. Hinton, S. Osindero, and Y.-W. Teh, "A fast learning algorithm for deep belief nets," *Neural computation*, vol. 18, no. 7, pp. 1527–1554, 2006.
- [25] Y. Bengio, P. Lamblin, D. Popovici, and H. Larochelle, "Greedy layer-wise training of deep networks," in *Advances in neural information processing systems*, pp. 153–160, 2007.

References (cont.)

- [26] G. E. Hinton, “Deep belief networks,” *Scholarpedia*, vol. 4, no. 5, p. 5947, 2009.
- [27] A.-r. Mohamed, G. Dahl, G. Hinton, et al., “Deep belief networks for phone recognition,” in *Nips workshop on deep learning for speech recognition and related applications*, vol. 1, p. 39, Vancouver, Canada, 2009.
- [28] A.-r. Mohamed and G. Hinton, “Phone recognition using restricted Boltzmann machines,” in *2010 IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 4354–4357, IEEE, 2010.
- [29] A.-r. Mohamed, G. E. Dahl, and G. Hinton, “Acoustic modeling using deep belief networks,” *IEEE transactions on audio, speech, and language processing*, vol. 20, no. 1, pp. 14–22, 2011.
- [30] G. W. Taylor, G. E. Hinton, and S. T. Roweis, “Modeling human motion using binary latent variables,” in *Advances in neural information processing systems*, pp. 1345–1352, 2007.
- [31] X. Glorot, A. Bordes, and Y. Bengio, “Deep sparse rectifier neural networks,” in *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pp. 315–323, JMLR Workshop and Conference Proceedings, 2011.
- [32] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, “Dropout: a simple way to prevent neural networks from overfitting,” *The journal of machine learning research*, vol. 15, no. 1, pp. 1929–1958, 2014.

References (cont.)

- [33] B. Ghoggh and M. Crowley, “The theory behind overfitting, cross validation, regularization, bagging, and boosting: tutorial,” *arXiv preprint arXiv:1905.12787*, 2019.
- [34] Y. LeCun, Y. Bengio, and G. Hinton, “Deep learning,” *nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [35] A. Karpathy and L. Fei-Fei, “Deep visual-semantic alignments for generating image descriptions,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3128–3137, 2015.
- [36] L. V. Fausett, *Fundamentals of neural networks: architectures, algorithms and applications*. Pearson Education India, 2006.