



Logistic Regression

T-662-ARTI

Stefán Ólafsson
Spring 2023



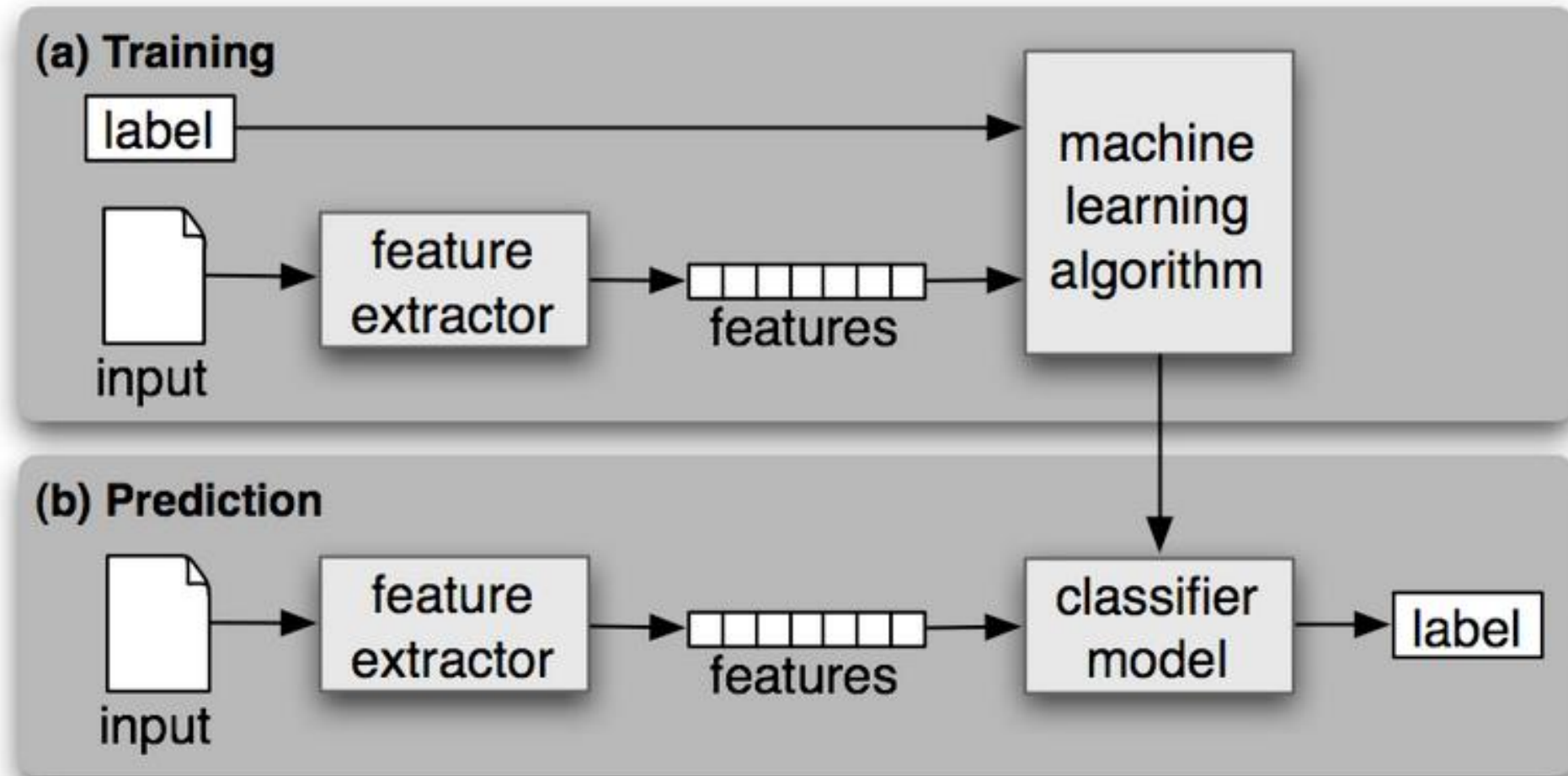


Classification

- The task of choosing the correct *class label* for a given input
- **Examples:**
 - Assigning subject categories, topics, or genres
 - Spam detection
 - Authorship identification
 - Language Identification
 - Sentiment analysis
 - ...



Supervised Classification





Generative vs. Discriminative Classifier

- **Naïve Bayes is a Generative Classifier**

$$c_{NB} = \underset{c \in \mathcal{C}}{\operatorname{argmax}} P(d|c)P(c)$$

- **$P(c|d)$ *not* computed directly**
- **$P(d|c)$ expresses how to ‘generate’ the features of d , given class c**



Generative vs. Discriminative Classifier

- **Logistic Regression is a discriminative classifier**
- **Attempts to directly compute $P(c|d)$**
- **Learns to assign high weight to features that improve its ability to “discriminate” between possible classes**
- **Cannot generate an example of one of the classes**



Logistic Regression

- **A training corpus of M observations:**
 - $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}) \dots (x^{(m)}, y^{(m)})$
- **Four components:**
 - For each input observation $x^{(i)}$:
 - Vector of features $[x_1, x_2, \dots, x_n]$
 - A classification function that computes the estimated class $\hat{y} \Rightarrow$ **sigmoid**
 - An objective function for learning, involving minimizing error on training examples \Rightarrow **cross-entropy loss**
 - An algorithm for optimizing the objective function \Rightarrow **gradient descent**



Logistic Regression

- **Estimate $P(y = 1 | x)$**
- **Learns a vector of weights and a bias term during training**
- **Each weight w_i is a real number associated with feature x_i**
 - Represents how important that feature is to the classification decision
 - If positive, then the feature is associated with the class
 - If negative, then the feature is not associated with the class



Classifying a test instance

- $z = \sum_{i=1}^n w_i x_i + b$
- $z = w * x + b$
- z expresses the weighted sum of the evidence for the class
- To create a probability, z is passed through the **sigmoid** function, $\sigma(z)$:
 - $y = \sigma(z) = \frac{1}{1+e^{-z}}$
 - $\hat{y} = \begin{cases} 1 & \text{if } p(y = 1|x) > 0.5 \\ 0 & \text{otherwise} \end{cases}$



The Sigmoid Function

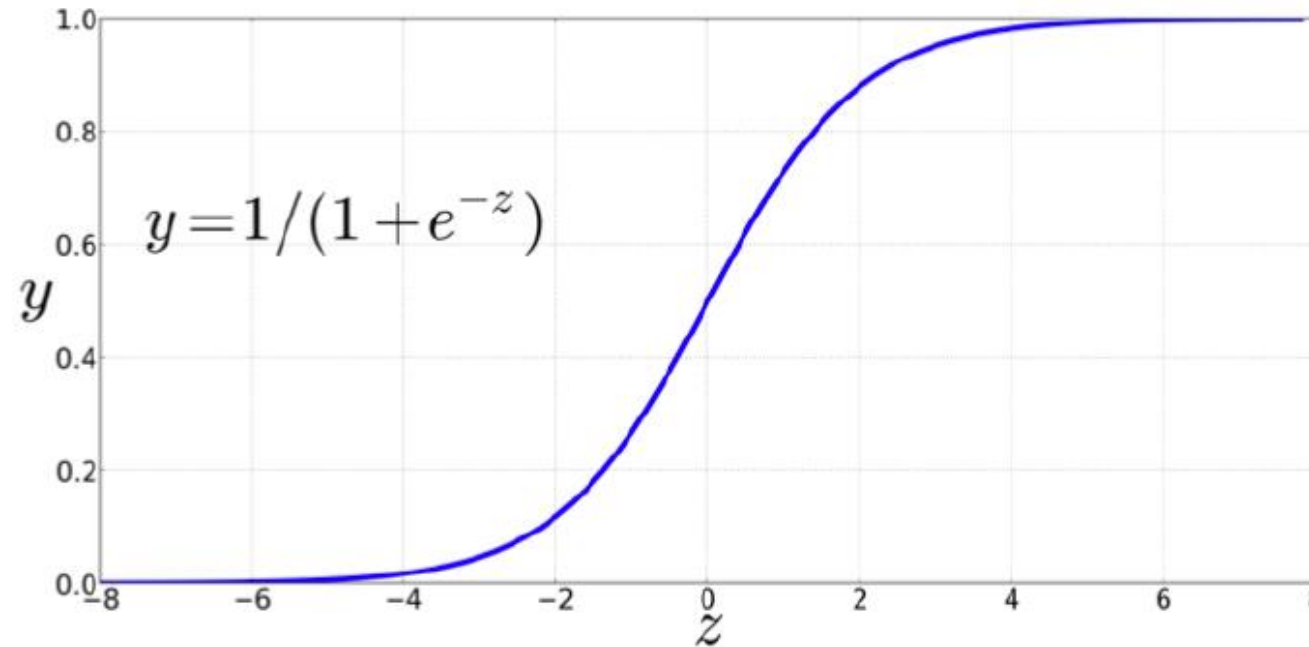


Figure 5.1 The sigmoid function $y = \frac{1}{1+e^{-z}}$ takes a real value and maps it to the range $[0, 1]$. It is nearly linear around 0 but outlier values get squashed toward 0 or 1.



Example: sentiment classification

It's **hokey**. There are virtually no surprises, and the writing is **second-rate**. So why was it so **enjoyable**? For one thing, the cast is **great**. Another **nice** touch is the music, **I** was overcome with the urge to get off the couch and start dancing. It sucked **me** in, and it'll do the same to **you**.

- **positive** = 3
- **negative** = 2
- “no” = 1
- 1st and 2nd person pronouns = 3
- “!” = 0
- $\log(\text{word count}) = \log(64) = 4.15$



Example: sentiment classification

Var	Definition	Value
x_1	count(positive lexicon \in doc)	3
x_2	count(negative lexicon \in doc)	2
x_3	$x_3 = \begin{cases} 1 & \text{if „no“} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$	1
x_4	Count(1 st and 2 nd person pronouns \in doc)	3
x_5	$x_5 = \begin{cases} 1 & \text{if „!“} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$	0
x_6	log(word count of doc)	4.19



Example: sentiment classification

- Let's assume we have already learned the weight vector **w**, and the bias **b**
- $w = [2.5, -5.0, -1.2, 0.5, 2.0, 0.7]$, $b = 0.1$
- $p(+|x) = p(Y=1|x) = \sigma(w * x + b)$
 $= \sigma([2.5, -5.0, -1.2, 0.5, 2.0, 0.7] * [3, 2, 1, 3, 0, 4.19] + 0.1)$
 $= \sigma(0.833) = \mathbf{0.70}$
- $p(-|x) = p(Y=0|x) = 1 - \sigma(w * x + b) = \mathbf{0.30}$



Learning in logistic regression

- Learn weights \mathbf{w} and bias \mathbf{b} that make \hat{y} for each training observation as close as possible to the true y
1. **loss (cost) function**
 2. **Iteratively updating the weights \Rightarrow gradient descent**



Cross-entropy loss function

- $L(\hat{y}, y)$ = How much \hat{y} differs from the true y
- We want to learn weights that maximize the probability of the correct label $p(y|x)$
- Two discrete outcomes (1 or 0)
- $p(y|x) = \hat{y}^y (1 - \hat{y})^{1-y}$
 - When $y=1$: $p(y|x) = \hat{y}$
 - When $y=0$: $p(y|x) = 1 - \hat{y}$



Cross-entropy loss function

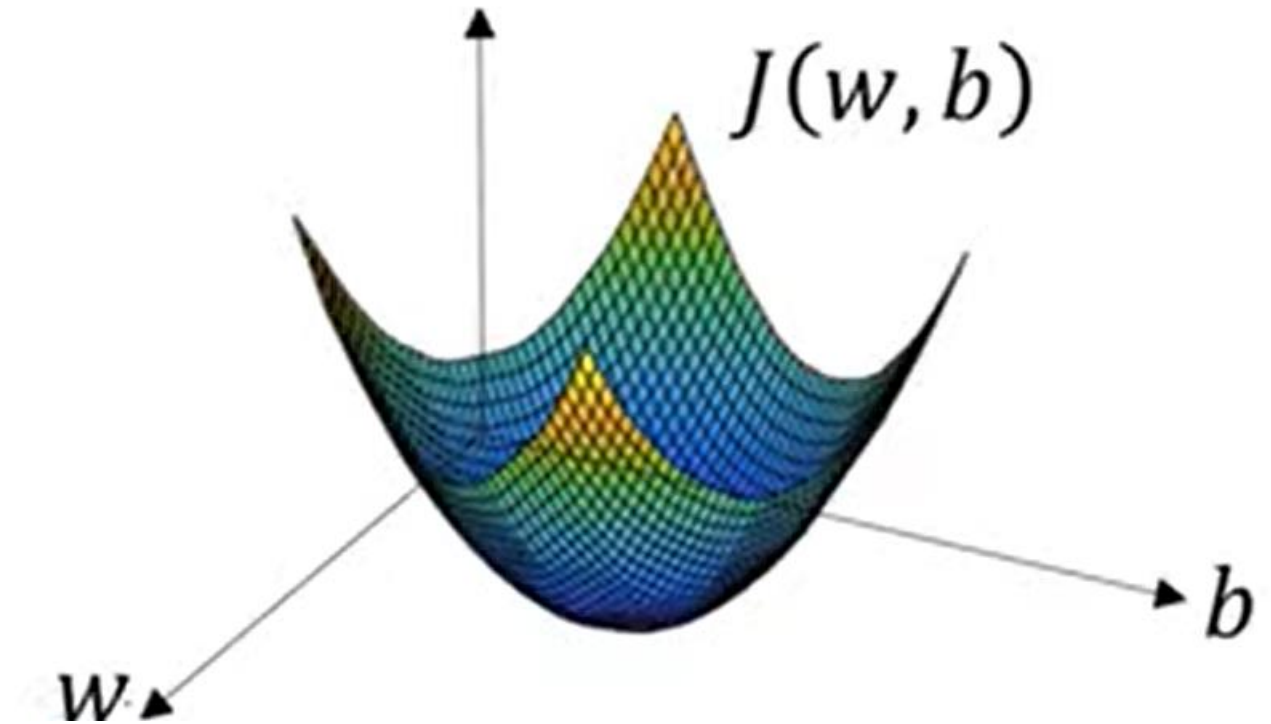
- $\log p(y|x) = y * \log \hat{y} + (1 - y) * \log(1 - \hat{y})$
- **To minimize:**
- $L_{CE}(\hat{y}, y) = -[y * \log \hat{y} + (1 - \hat{y}) * \log(1 - \hat{y})]$
 - When $y = 1$: this is $-\log \hat{y}$
 - When $y = 0$: this is $-\log(1 - \hat{y})$



Gradient Descent

- **Used to find the optimal weights**
- **Minimize the loss function**
 - Cost function:

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$



Andrew Ng - <https://www.youtube.com/watch?v=uJryes5Vk1o>

HÁSKÓLINN Í REYKJAVÍK | REYKJAVÍK UNIVERSITY

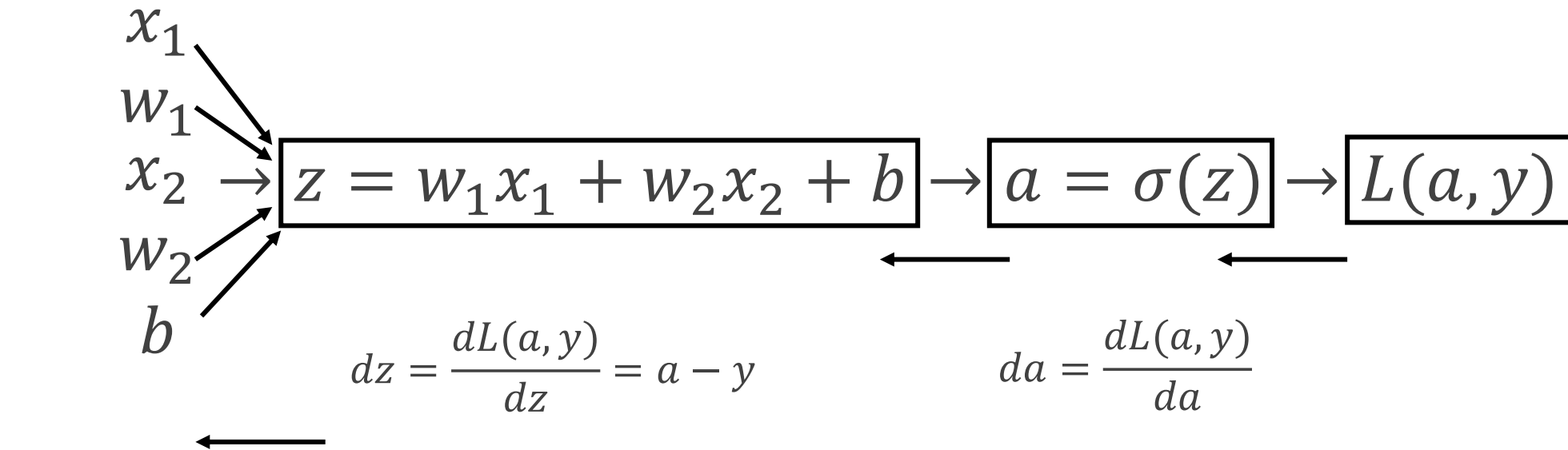


Gradient Descent





Logistic regression derivatives



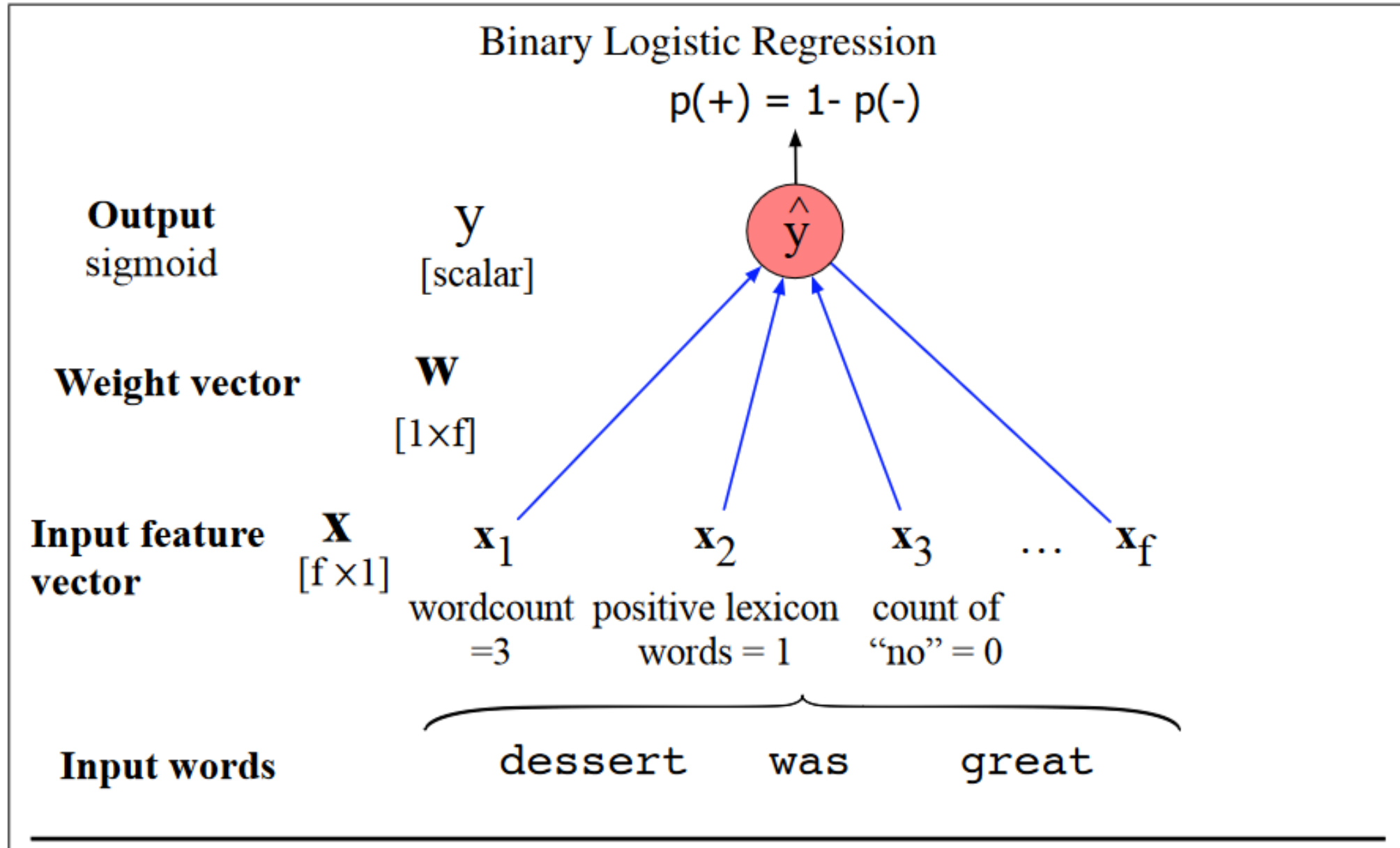
$$dw_1 = \frac{dL}{dw_1} = x_1 * dz$$

Parameter updates:

$$w_1 \leftarrow w_1 - \alpha * dw_1$$

$$w_2 \leftarrow w_2 - \alpha * dw_2$$

$$b \leftarrow b - \alpha * db$$





Multinomial Logistic Regression

