

CSC 561: Neural Networks and Deep Learning

Gradient Descent

Marco Alvarez

Department of Computer Science and Statistics
University of Rhode Island

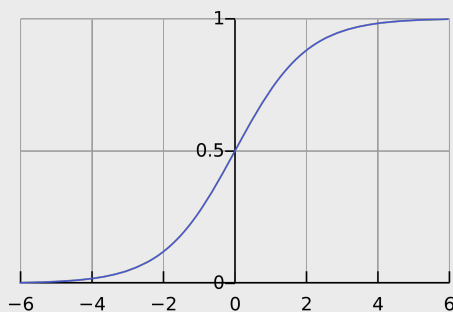
Spring 2024



Logistic regression

Logistic function

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



mapping \mathbb{R} to $[0,1]$
continuous and
differentiable

3

Logistic regression

• Binary classifier

- ✓ models $P(y | \mathbf{x})$, $\mathbf{x} \in \mathbb{R}^d$, $y \in \{+1, -1\}$
 - a threshold (e.g., $\theta = 0.5$) may be used for a final binary classification

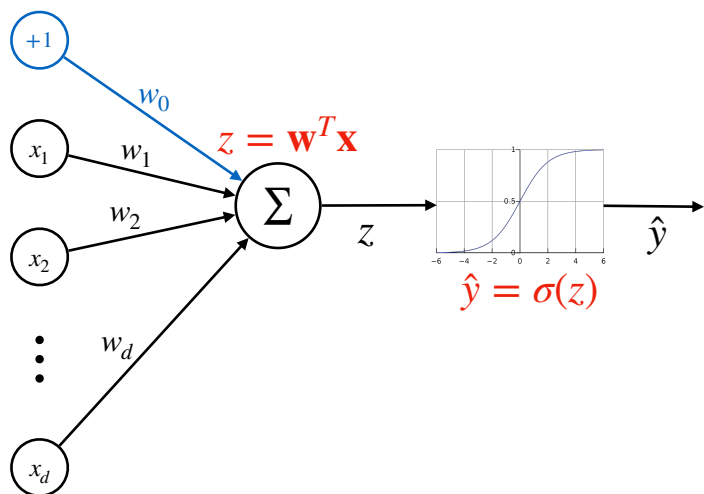
- ✓ uses the logistic function

• Considered a **linear classifier**

- ✓ although the “activation function” is nonlinear

4

NN-like



5

Probabilistic interpretation

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

(probability of class +1) $P(y = +1 | \mathbf{x}; \mathbf{w}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} = \sigma(\mathbf{w}^T \mathbf{x})$

(probability of class -1) $P(y = -1 | \mathbf{x}; \mathbf{w}) = \frac{1}{1 + e^{\mathbf{w}^T \mathbf{x}}} = \sigma(-\mathbf{w}^T \mathbf{x})$

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

note that $P(y|\mathbf{x}) > 0.5$ when $y\mathbf{w}^T \mathbf{x} > 0$ (correct classifications)

6

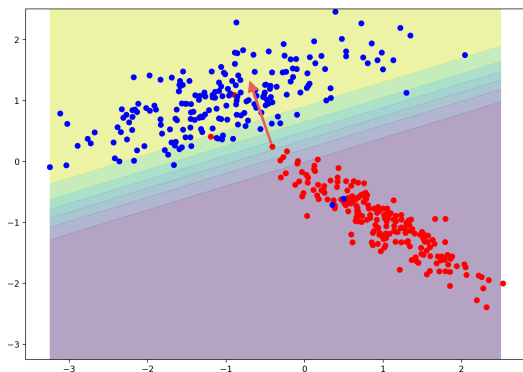
Linear decision boundary

$$P(y | \mathbf{x}; \mathbf{w}) = \frac{1}{1 + e^{-y\mathbf{w}^T \mathbf{x}}} = \frac{1}{2}$$

$$1 + e^{-y\mathbf{w}^T \mathbf{x}} = 2$$

$$e^{-y\mathbf{w}^T \mathbf{x}} = 1$$

$$\mathbf{w}^T \mathbf{x} = 0$$



7

Learning the parameters

MLE

Maximum likelihood estimation

- choose parameters \mathbf{w} that **maximize** the **conditional data likelihood** $P(\mathbf{y} | X; \mathbf{w})$, i.e., the probability of the observed values conditioned on the feature values

i.i.d. assumption $P(\mathbf{y} | X; \mathbf{w}) = \prod_{i=1}^n P(y^{(i)} | \mathbf{x}^{(i)}; \mathbf{w})$

- objective function:

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \prod_{i=1}^n P(y^{(i)} | \mathbf{x}^{(i)}; \mathbf{w})$$

9

MLE

$$\begin{aligned} \mathbf{w}^* &= \arg \max_{\mathbf{w}} \prod_{i=1}^n P(y^{(i)} | \mathbf{x}^{(i)}; \mathbf{w}) \\ &= \arg \max_{\mathbf{w}} \log \left(\prod_{i=1}^n P(y^{(i)} | \mathbf{x}^{(i)}; \mathbf{w}) \right) \\ &= \arg \max_{\mathbf{w}} \frac{1}{n} \log \left(\prod_{i=1}^n P(y^{(i)} | \mathbf{x}^{(i)}; \mathbf{w}) \right) \frac{1}{1 + e^{-y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)}}} \\ &= \arg \max_{\mathbf{w}} - \frac{1}{n} \sum_{i=1}^n \log \left(1 + e^{-y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)}} \right) \\ \text{negative log likelihood} &= \arg \min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \log \left(1 + e^{-y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)}} \right) \text{ per instance loss} \end{aligned}$$

10

Applying gradient descent

cross-entropy loss: no closed-form solution, but loss is convex

$$J(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \log \left(1 + e^{-y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)}} \right)$$

gradient

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = \left[\frac{\partial J(\mathbf{w})}{\partial w_0}, \dots, \frac{\partial J(\mathbf{w})}{\partial w_d} \right]$$

$$\begin{aligned} f(x) &= \log(1 + e^x) \\ f'(x) &= \frac{e^x}{1 + e^x} \\ &= \sigma(x) \end{aligned}$$

$$\frac{\partial J(\mathbf{w})}{\partial w_j} = - \frac{1}{n} \sum_{i=1}^n \sigma \left(-y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)} \right) y^{(i)} x_j^{(i)}$$

11

Show me the code

```
# ...
```

12

Extension to multiple classes

From binary to C>2 classes

- Binary logistic regression

$$P(y = +1 | \mathbf{x}; \mathbf{w}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} = \frac{e^{\mathbf{w}^T \mathbf{x}}}{e^{\mathbf{w}^T \mathbf{x}} + 1}$$

- C classes

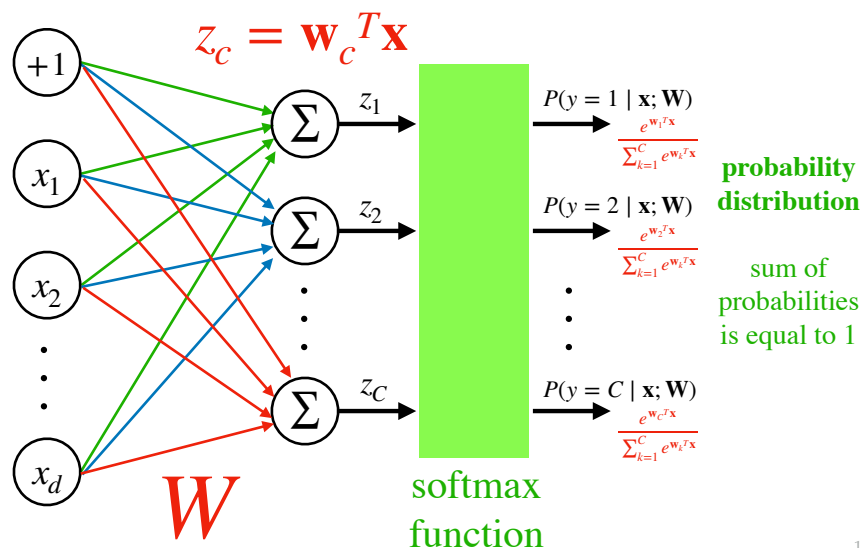
softmax function

$$P(y = c | \mathbf{x}; \mathbf{W}) = \frac{e^{\mathbf{w}_c^T \mathbf{x}}}{\sum_{k=1}^C e^{\mathbf{w}_k^T \mathbf{x}}}$$

$\mathbf{W}_{C \times (d+1)}$ is a matrix where rows are “per-class” weight vectors

14

NN-like



15

Multinomial logistic regression

- Replace the **logistic** by the **softmax** function
 - predict the class with the highest probability score

$$\hat{y} = \arg \max_c P(y = c | \mathbf{x}; \mathbf{W})$$

- How to learn the parameters?
 - use MLE to derive a **loss function** ... then apply **gradient descent**

16

MLE

$$\mathbf{W}^* = \arg \max_{\mathbf{W}} \frac{1}{n} \prod_{i=1}^n P(y^{(i)} \mid \mathbf{x}^{(i)}; \mathbf{W})$$

$$= \arg \max_{\mathbf{W}} \frac{1}{n} \prod_{i=1}^n \prod_{c=1}^C P(y^{(i)} = c \mid \mathbf{x}^{(i)}; \mathbf{W})^{t_{i,c}}$$

Consider a matrix $T_{n \times C}$ where every row is a **one-hot encoding** of the target variable

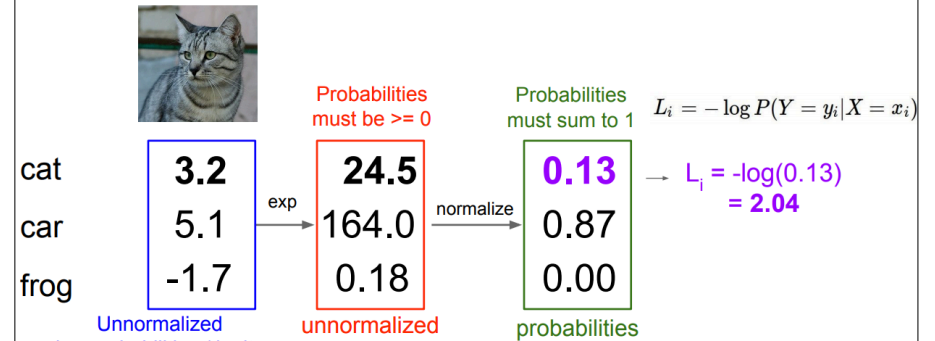
$$= \arg \max_{\mathbf{W}} \frac{1}{n} \sum_{i=1}^n \sum_{c=1}^C t_{i,c} \log (P(y^{(i)} = c \mid \mathbf{x}^{(i)}; \mathbf{W}))$$

$$= \arg \min_{\mathbf{W}} -\frac{1}{n} \sum_{i=1}^n \sum_{c=1}^C t_{i,c} \log (P(y^{(i)} = c \mid \mathbf{x}^{(i)}; \mathbf{W}))$$

$$= \arg \min_{\mathbf{W}} \frac{1}{n} \sum_{i=1}^n \sum_{c=1}^C -t_{i,c} \log \left(\frac{e^{\mathbf{w}_c \mathbf{x}^{(i)}}}{\sum_{k=1}^C e^{\mathbf{w}_k \mathbf{x}^{(i)}}} \right) \quad \begin{array}{l} \text{per instance} \\ \text{cross-entropy} \\ \text{loss} \end{array}$$

17

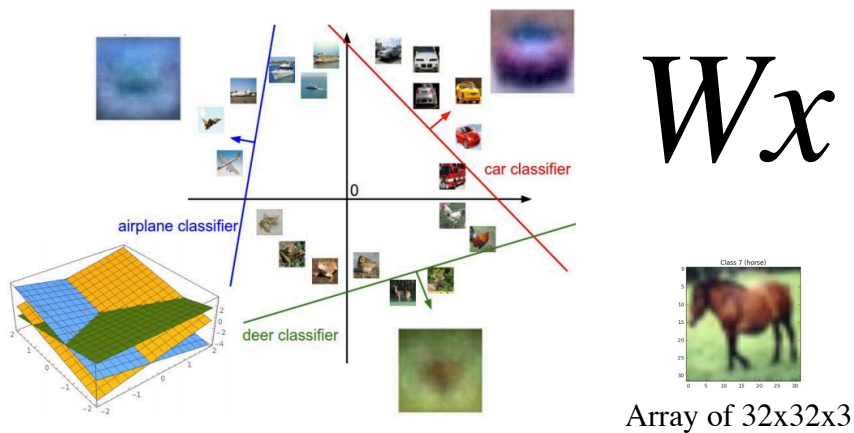
Softmax and the cross-entropy loss



Stanford University CS231n: Deep Learning for Computer Vision

18

Geometric interpretation



Stanford University CS231n: Deep Learning for Computer Vision

19