

# CSC 561: Neural Networks and Deep Learning

## Logistic Regression

Marco Alvarez

Department of Computer Science and Statistics  
University of Rhode Island

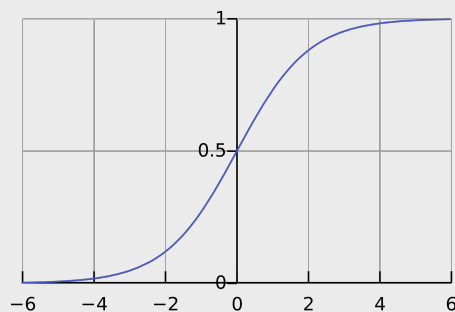
Spring 2024



# 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

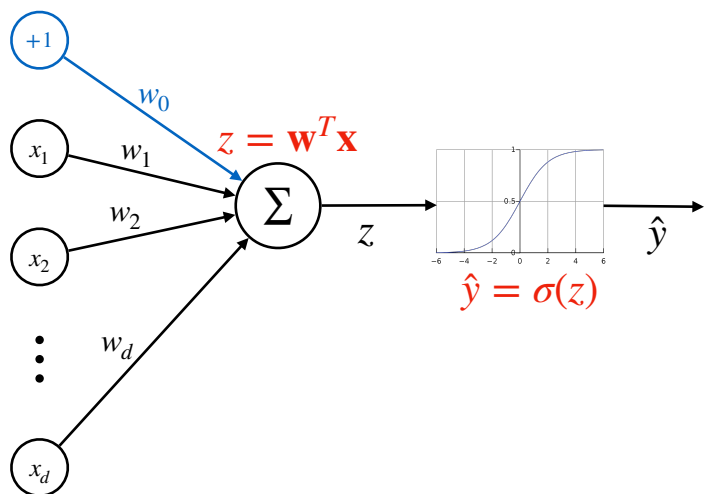
- ✓ uses the logistic function

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

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

compact expression

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

6

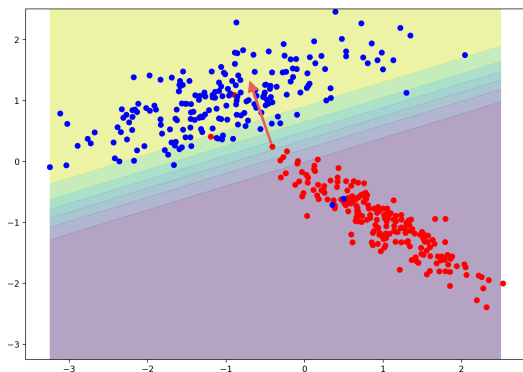
## Linear decision boundary

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

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

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

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



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

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

maximize  $P(y^{(i)} = +1 | \mathbf{x}^{(i)})$  for any  $\mathbf{x}^{(i)}$  with a positive label, and maximize  $P(y^{(i)} = -1 | \mathbf{x}^{(i)})$  for any  $\mathbf{x}^{(i)}$  with a negative label

$$\begin{aligned} &= \arg \max_{\mathbf{w}} \frac{1}{n} \log \left( \prod_{i=1}^n P(y^{(i)} | \mathbf{x}^{(i)}; \mathbf{w}) \right) \\ &= \arg \max_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \log (P(y^{(i)} | \mathbf{x}^{(i)}; \mathbf{w})) \quad \frac{1}{1 + e^{-y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)}}} \\ &= \arg \max_{\mathbf{w}} - \frac{1}{n} \sum_{i=1}^n \log (1 + e^{-y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)}}) \end{aligned}$$

**minimize the negative log likelihood**  $= \arg \min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \log (1 + e^{-y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)}})$  **per instance loss**

10

## Applying gradient descent

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

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

**gradient**

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

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

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

**partial derivative**

11

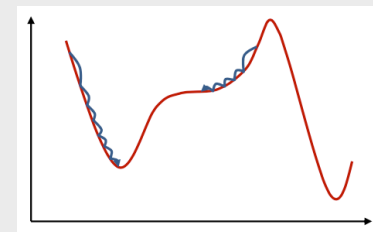
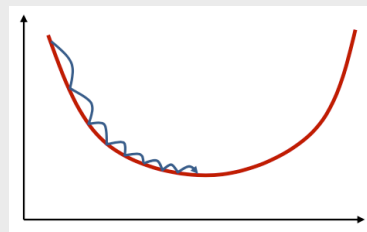
## Gradient descent and convex functions

### Convex functions

- for appropriate learning rates, GD will always find the minimum

### Non-convex functions

- GD may find a local minimum (or inflection point)



12

## Show me the code

```
def __forward(self, X):
    z = X @ self.weights
    y_pred = self.__sigmoid(z)
    return y_pred

def __loss(self, X, y):
    inv = -y * (X @ self.weights)
    loss = np.log(1 + np.exp(inv)).mean()
    dw = self.__sigmoid(inv) * -y * X
    dw = np.mean(dw, axis=0, keepdims=True).T
    return loss, dw
```

13

## Understanding matrix form

$$\mathbf{y} = \mathbf{X}\mathbf{w}$$

input vector  $\mathbf{x}^{(1)}$

$$\begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix} \approx \begin{bmatrix} (\mathbf{x}^{(1)})^T \\ \vdots \\ (\mathbf{x}^{(n)})^T \end{bmatrix} \cdot \begin{bmatrix} w_0 \\ \vdots \\ w_d \end{bmatrix}$$

14

## Extension to multiple classes

### From binary to $C > 2$ classes

#### • Multinomial logistic regression

- ✓ models  $P(y = c | \mathbf{x}; W)$ ,  $\mathbf{x} \in \mathbb{R}^d$ ,  $y \in \{0, \dots, c - 1\}$
- ✓ uses the **softmax** function

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

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

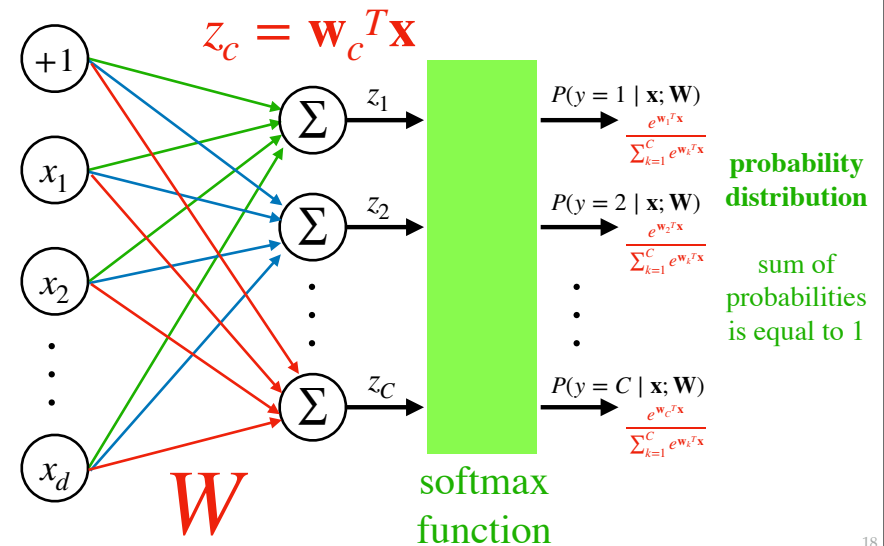
16

## One-hot encoding

	y	red	green	blue	yellow
red	0	1	0	0	0
green	1	0	1	0	0
blue	2	0	0	1	0
blue	2	0	0	1	0
green	1	0	1	0	0
blue	2	0	0	1	0
yellow	3	0	0	0	1
blue	2	0	0	1	0
red	0	1	0	0	0
red	0	1	0	0	0

17

## NN-like



18

## Multinomial logistic regression

- Predict the class with the highest probability

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

- How to learn the parameters?

- use MLE to derive a **loss function** ... then apply **gradient descent**

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

19

## MLE

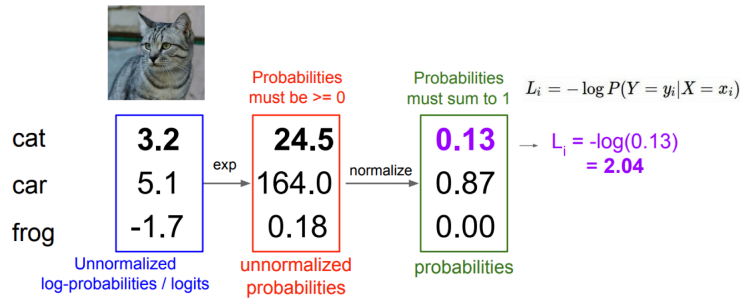
$$\begin{aligned} \mathbf{W}^* &= \arg \max_{\mathbf{W}} \frac{1}{n} \prod_{i=1}^n P(y^{(i)} | \mathbf{x}^{(i)}; \mathbf{W}) \\ &= \arg \max_{\mathbf{W}} \frac{1}{n} \prod_{i=1}^n \prod_{c=1}^C P(y^{(i)} = c | \mathbf{x}^{(i)}; \mathbf{W})^{t_{i,c}} \\ &= \arg \max_{\mathbf{W}} \frac{1}{n} \sum_{i=1}^n \sum_{c=1}^C t_{i,c} \log (P(y^{(i)} = c | \mathbf{x}^{(i)}; \mathbf{W})) \\ &= \arg \min_{\mathbf{W}} -\frac{1}{n} \sum_{i=1}^n \sum_{c=1}^C t_{i,c} \log (P(y^{(i)} = c | \mathbf{x}^{(i)}; \mathbf{W})) \\ &= \arg \min_{\mathbf{W}} \frac{1}{n} \sum_{i=1}^n \left[ \sum_{c=1}^C -t_{i,c} \log \left( \frac{e^{\mathbf{w}_c^T \mathbf{x}^{(i)}}}{\sum_{k=1}^C e^{\mathbf{w}_k^T \mathbf{x}^{(i)}}} \right) \right] \end{aligned}$$

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

per instance cross-entropy loss

20

## Softmax and the cross-entropy loss



0.13	compare	1.00
0.87	Kullback-Leibler divergence	0.00
0.00	$D_{KL}(P  Q)$	0.00
$Q$		$P$

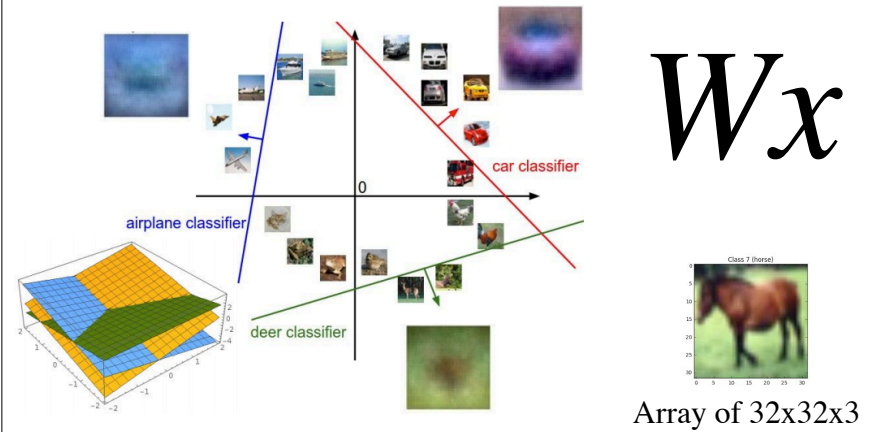
$$D_{KL}(P||Q) = \sum P(y) \log \frac{P(y)}{Q(y)}$$

$$H(P, Q) = D_{KL}(P||Q) - H(P)$$

cross-entropy

entropy

## Geometric interpretation



Stanford University CS231n: Deep Learning for Computer Vision

22

## Regularization

### Regularization

- Adding a **penalty** to the weights to control the complexity of the model
  - usually penalizing higher weights (**except intercept**)
  - results in **simpler** or **more sparse** solutions
- Impact of regularization can be controlled by a hyperparameter (*lambda*)

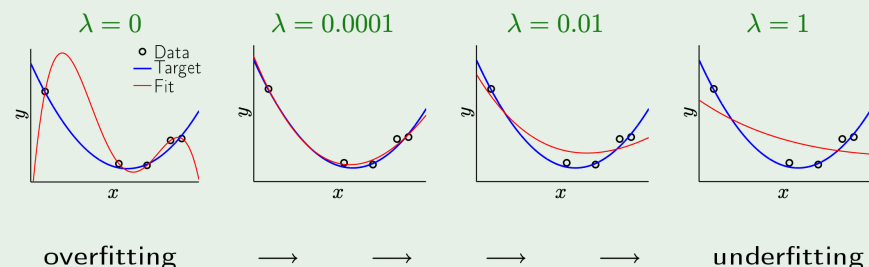
$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n l(h_{\mathbf{w}}, x_i, y_i) + \lambda R(\mathbf{w})$$

better predictions

control overfitting

24

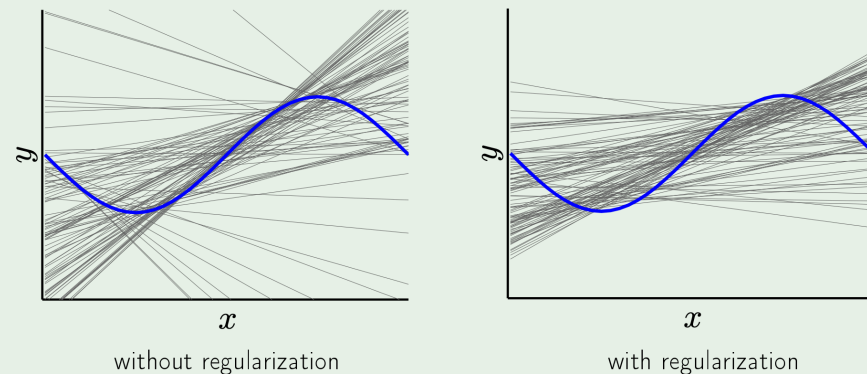
## Regularization prefers simpler models



<https://work.caltech.edu/slides/slides12.pdf>

25

## Restricting the hypothesis space



<https://work.caltech.edu/slides/slides12.pdf>

26

## Regularization

### • L1

- tends to drive weights to zero (**sparse solution**) effectively performing feature selection, which can be beneficial for interpretability and reducing model complexity
- less sensitive to outliers compared to L2 regularization
- absolute value function used in L1 regularization is not smooth, making the optimization process slightly more complex compared to L2

$$R(W) = \sum |W_{ij}|$$

### • L2

- encourages smaller coefficients by penalizing large ones more heavily
- by shrinking weights, L2 regularization can improve the stability of the model, making it less sensitive to small changes in the data
- can be more sensitive to outliers compared to L1 regularization
- most common form of regularization (superior performance)

$$R(W) = \sum W_{ij}^2$$

### • Other forms of regularization in neural networks

- elastic net, dropout, batch normalization, early stopping, data augmentation

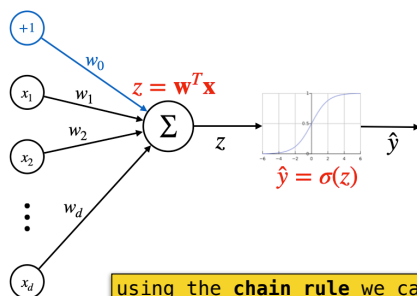
27

The importance of using differentiable functions

# Differentiable activation functions

## Continuous activation functions

- ✓ sigmoid, tanh, RELU, etc.
- ✓ differentiable (almost) everywhere



$$\frac{d\hat{y}}{dz} = \sigma'(z)$$

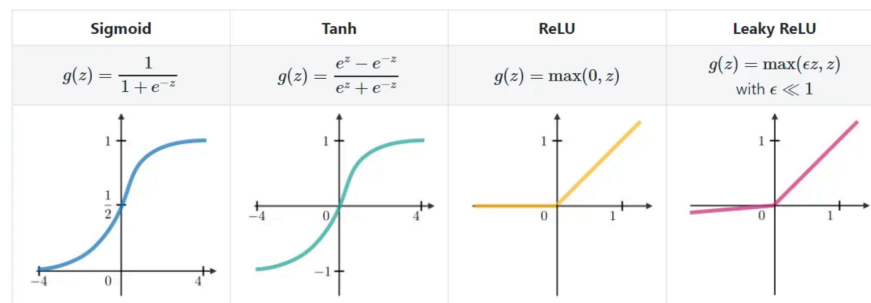
$$\frac{d\hat{y}}{dw_i} = \frac{d\hat{y}}{dz} \frac{dz}{dw_i} = \sigma'(z)x_i$$

$$\frac{d\hat{y}}{dx_i} = \frac{d\hat{y}}{dz} \frac{dz}{dx_i} = \sigma'(z)w_i$$

using the **chain rule** we can compute the change in the output for small changes to the input/weights

29

# A few examples



ReLU is a good default choice for most networks

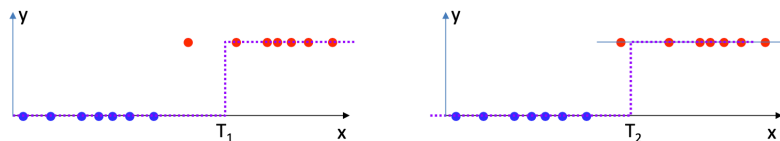
<https://medium.com/analytics-vidhya/understanding-activation-functions-data-science-for-the-rest-of-us-b652048a064f>

30

# Differentiable loss functions

## Threshold activation

- ✓ shifting threshold does not change classification error



## Continuous loss and activation

- ✓ can quantify a loss between continuous output and the desired target
- the loss changes as the weights change, even if the classification error remains the same

