# Multinomial Logistic Regression

## CSC 461: Machine Learning

Fall 2020

Prof. Marco Alvarez
University of Rhode Island

---

# Logistic regression

‣ Binary classification

✓ uses a **logistic function** (type of sigmoid function)

✓ models **probability** of output in terms of input

‣ What if we want k > 2 classes?

✓ can try one-vs-all

✓ learn a binary classifier per class; relabel training data with samples of that class as positive and all other as negatives; predict using the highest score from all classifiers

✓ can try one-vs-one

✓ learn k (k-1) / 2 binary classifiers; each learns to distinguish between two classes; predict using a voting scheme

---

# Issues with OvA or OvO

‣ Class imbalance

‣ Scale of scores may differ from classifier to classifier

‣ Computational cost (both train and predict)

---

# MNIST

‣ The MNIST database is a large database of handwritten digits

✓ contains 60,000 training images and 10,000 testing images

✓ convolutional neural networks, manages to get an error rate of 0.23%

✓ original paper reports an error rate of 0.8% with SVMs

**http://yann.lecun.com/exdb/mnist/**

# Basics of multiclass classification

‣ Data instance

  ✓ in general, $x \in \mathcal{X}, \mathcal{X} = \mathbb{R}^d$

  ✓ $y \in \mathcal{Y}, \mathcal{Y} = \{1, 2, \ldots, C\}$

‣ Hypothesis

  ✓ each hypothesis **g** is a classifier

$$g : \mathcal{X} \mapsto \mathcal{Y}, g \in \mathcal{H}$$

# From binary to k classes

‣ Binary logistic regression:

$$P(y = +1 \mid \mathbf{x}; \mathbf{w}) = \frac{1}{1 + e^{-\mathbf{w}^T\mathbf{x}}} = \frac{e^{\mathbf{w}^T\mathbf{x}}}{e^{\mathbf{w}^T\mathbf{x}} + 1}$$
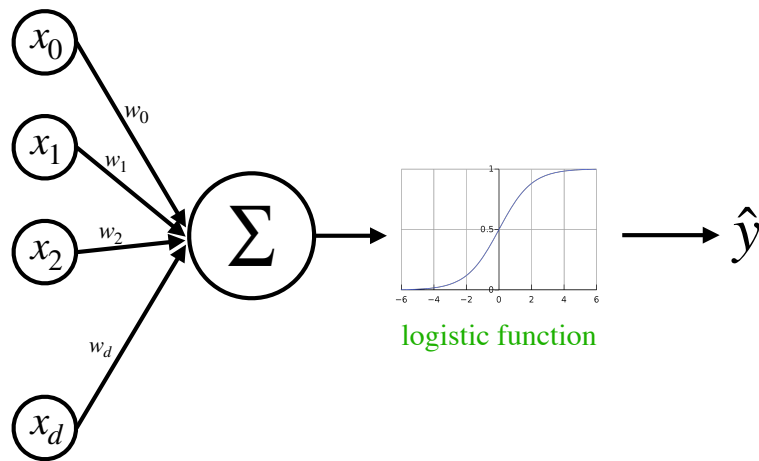
‣ Can be extended to:

$$P(y = c \mid \mathbf{x}; \mathbf{W}) = \frac{e^{\mathbf{w}_c\mathbf{x}}}{\sum_{k=1}^{C} e^{\mathbf{w}_k\mathbf{x}}}$$
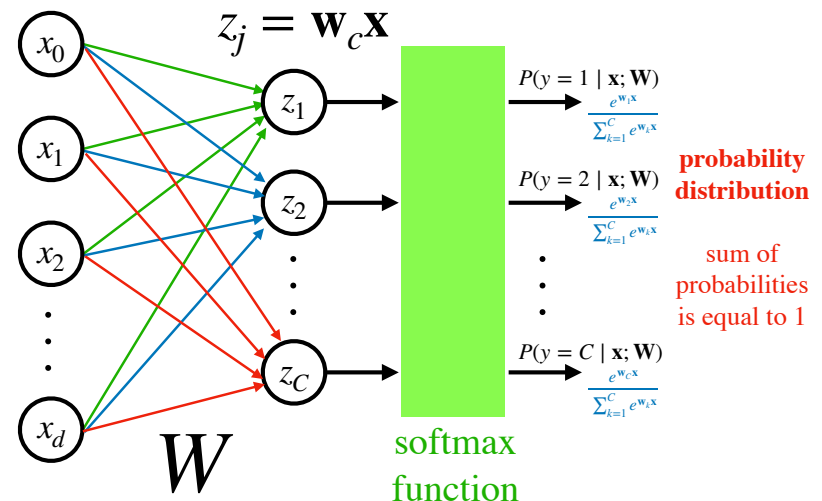
$\mathbf{w}_c$ is the row vector c from **W**

$W_{C \times d+1}$ is a matrix where every row is a "class" weight vector

# Logistic regression (binary classification)



logistic function

# Multinomial logistic regression



$$z_j = \mathbf{W}_c\mathbf{x}$$

$P(y = 1 \mid \mathbf{x}; \mathbf{W})$   $\frac{e^{\mathbf{w}_1\mathbf{x}}}{\sum_{k=1}^{C} e^{\mathbf{w}_k\mathbf{x}}}$

$P(y = 2 \mid \mathbf{x}; \mathbf{W})$   $\frac{e^{\mathbf{w}_2\mathbf{x}}}{\sum_{k=1}^{C} e^{\mathbf{w}_k\mathbf{x}}}$

$P(y = C \mid \mathbf{x}; \mathbf{W})$   $\frac{e^{\mathbf{w}_C\mathbf{x}}}{\sum_{k=1}^{C} e^{\mathbf{w}_k\mathbf{x}}}$

**probability distribution**

sum of probabilities is equal to 1

$W$

softmax function

# Example

‣ What is the value of softmax($\mathbf{z}$), given that $\mathbf{z}^T = [-10, 10, 5, 4.3, 7]$?

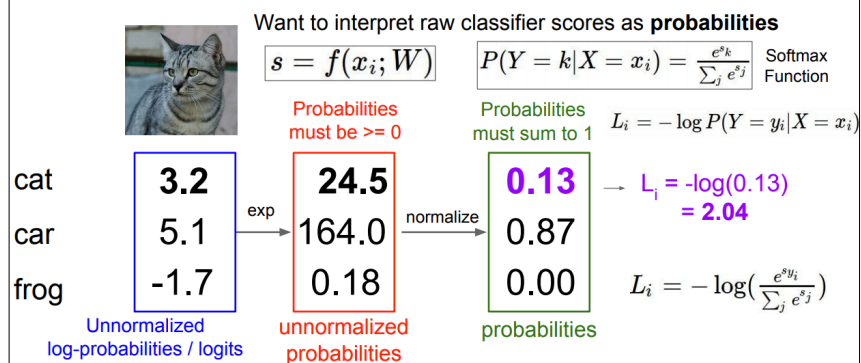# Multinomial logistic regression

‣ Use the **softmax function** for activation

‣ Predict the label with the highest probability score (forward pass)

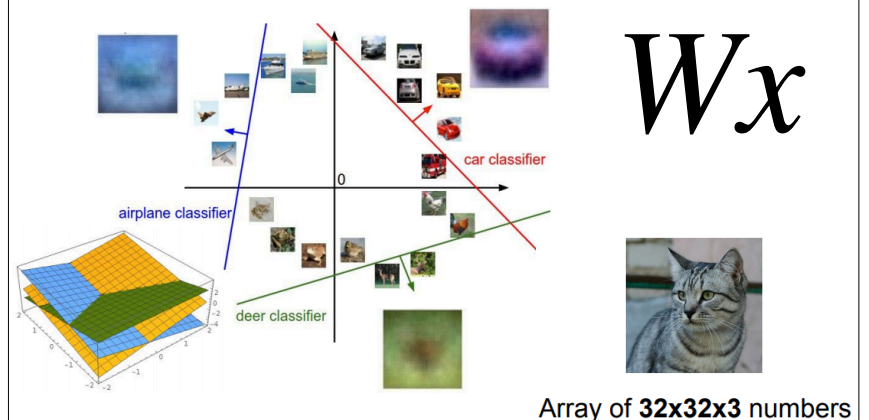$$\hat{y} = \arg\max_c P(y = c \mid \mathbf{x}; \mathbf{W})$$

‣ How to learn the weights?

✓ need to define a **Loss Function** … then apply **gradient descent**

✓ loss function can be derived using **MLE** (similar to binary logistic regression)
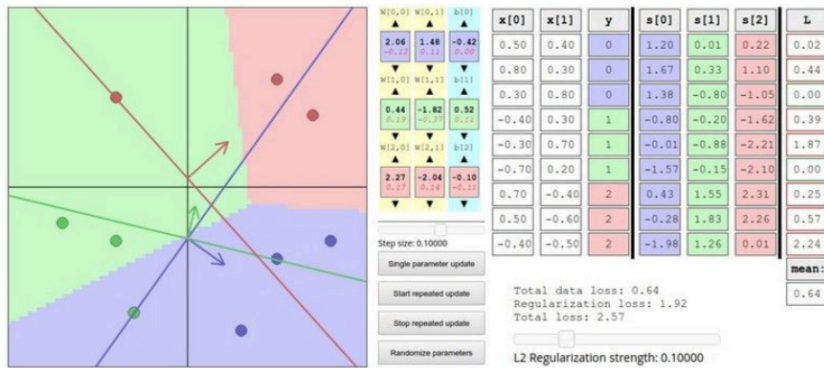
# Softmax and the loss function

Want to interpret raw classifier scores as **probabilities**

$$s = f(x_i; W) \qquad P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}} \quad \text{Softmax Function}$$

| | Unnormalized log-probabilities / logits | | unnormalized probabilities | | probabilities | |
|---|---|---|---|---|---|---|
| cat | **3.2** | | **24.5** | | **0.13** | |
| car | 5.1 | exp | 164.0 | normalize | 0.87 | |
| frog | -1.7 | | 0.18 | | 0.00 | |

Probabilities must be >= 0

Probabilities must sum to 1

$$L_i = -\log P(Y = y_i | X = x_i)$$

$\rightarrow$ L$_i$ = -log(0.13) = 2.04

$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

http://vision.stanford.edu/teaching/cs231n/slides/2019/cs231n_2019_lecture03.pdf

# Geometric interpretation

$$Wx$$

car classifier

airplane classifier

deer classifier

Array of **32x32x3** numbers

http://vision.stanford.edu/teaching/cs231n/slides/2019/cs231n_2019_lecture02.pdf

# Interactive web demo



http://vision.stanford.edu/teaching/cs231n-demos/linear-classify/