

# Teaching Evaluations

<https://heliocampusac.wisc.edu/>

# Multiclass Logistic Regression and Neural Networks

- cross entropy
- logistic regression
- neural network basics

# Entropy, Cross-Entropy

- entropy is a fundamental quantity that measures the number of bits required to encode a sequences of i.i.d. random variables
- we can encode/compress  $n$  i.i.d. random variables with  $nH(X)$  bits
- to build a code, we needed to know  $p(x)$ .

$$X_i \stackrel{i.i.d.}{\sim} \text{bernoulli}(\theta) \quad H_2(\theta) = \theta \log \frac{1}{\theta} + (1 - \theta) \log \frac{1}{1-\theta}$$

$$H(X) = \sum_i p_i \log_2 \left( \frac{1}{p_i} \right)$$

- approach: only encode/compress vectors that had about  $\theta n$  ones.
- code is optimized for  $p(x)$
- field of study on constructing codes with desirable properties

# Entropy, Cross-Entropy

```
model.compile(optimizer='sgd', loss=tf.keras.losses.CategoricalCrossentropy())
```

- entropy:

$$H(\mathbf{p}) = \sum_i p_i \log \frac{1}{p_i}$$

- number of bits required to encode a sequences of i.i.d. random variables

- cross-entropy:

$$H(\mathbf{p}, \mathbf{q}) = \sum_i p_i \log \frac{1}{q_i}$$

- number of bits required to encode a sequences of i.i.d. random variables from  $\mathbf{p}$  using a code optimized for  $\mathbf{q}$ .
- $H(\mathbf{p}, \mathbf{q})$  is big if  $q_i \ll p_i$  for some  $i$

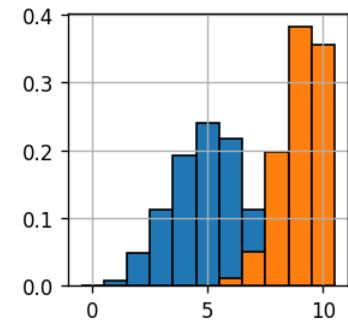
- Kullback-Leibler Divergence:

$$D(\mathbf{p} || \mathbf{q}) = \sum_i p_i \log \frac{p_i}{q_i}$$

- distance between  $\mathbf{p}$  and  $\mathbf{q}$
- important in binary hypothesis testing

- relationship:

$$H(\mathbf{p}, \mathbf{q}) = H(\mathbf{p}) + D(\mathbf{p} || \mathbf{q})$$

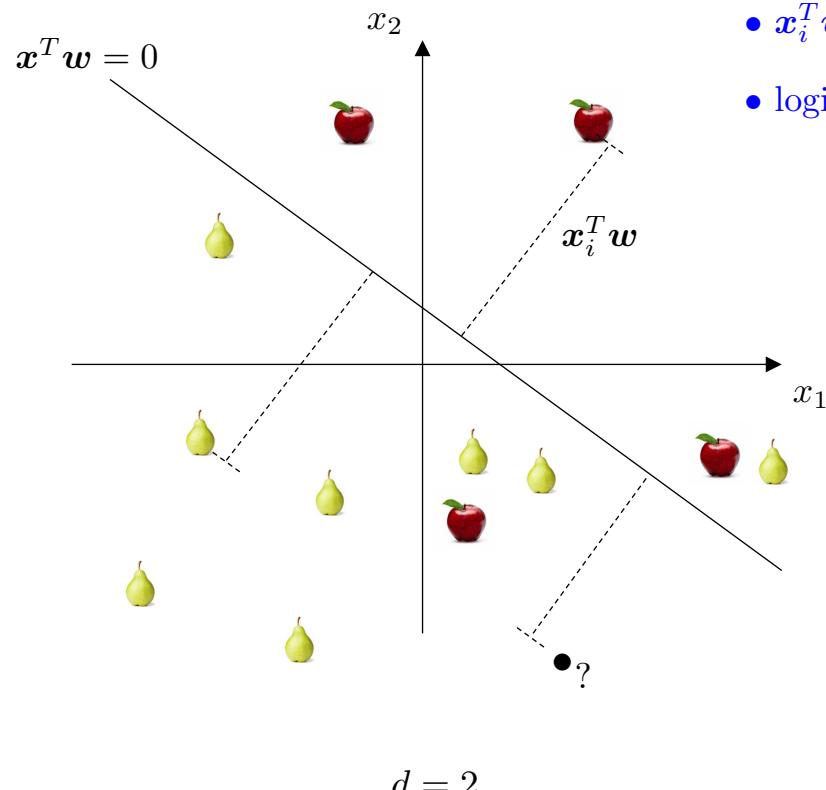


# Two Class Logistic Regression

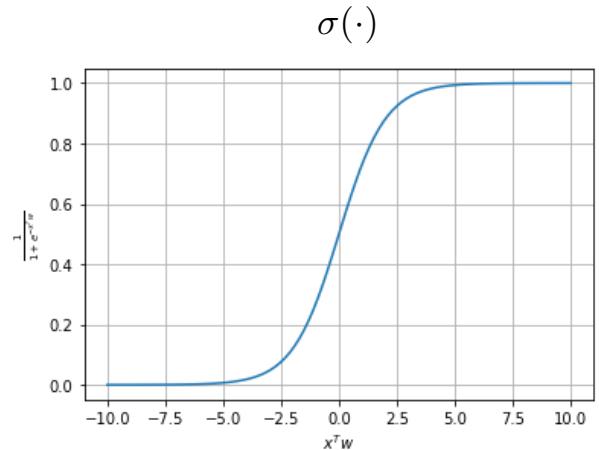
- logistic regression is a *discriminative* classifier

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

directly fit  $p(y|\mathbf{x})$  to data by treating  $\mathbf{x}$  as non-random



- decision boundary is  $\mathbf{x}^T \mathbf{w} = 0$
- $\mathbf{x}_i^T \mathbf{w}$  is the signed distance from the decision boundary
- logistic function takes  $\mathbf{x}_i^T \mathbf{w}$  and squishes it to make it a probability

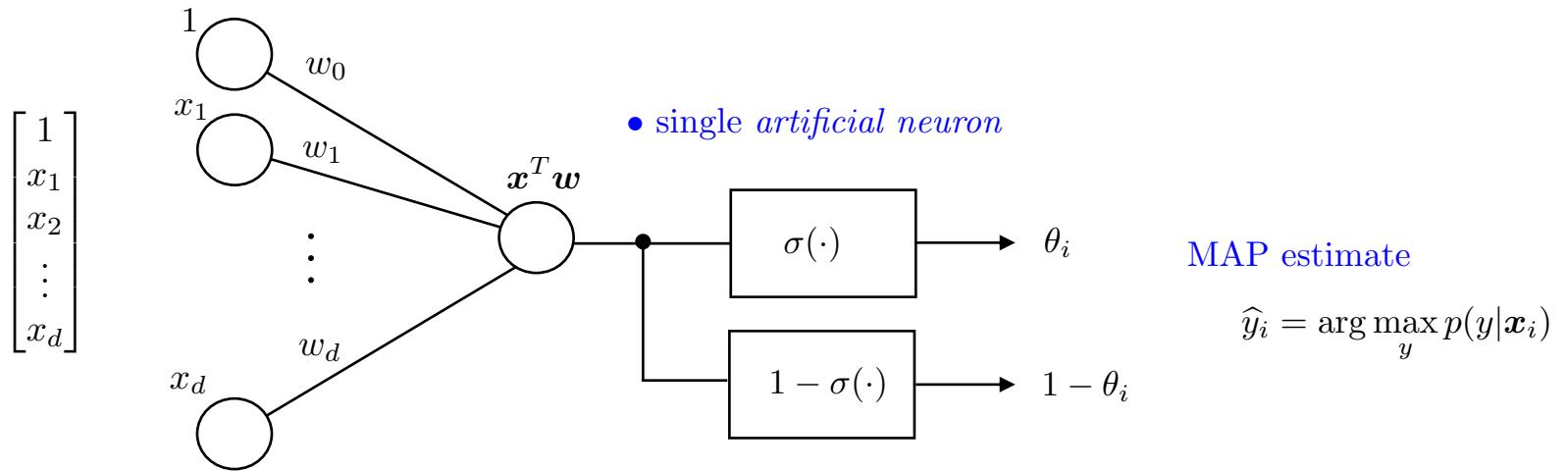


- probability that ? is a pear/apple depends distance from boundary

# Logistic Regression

- testing

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



- training

- notation

$$\mathcal{D} = \{(\mathbf{x}_i, y)\}_{i=1}^n$$

$$p_{\mathbf{w}}(y_1, y_2, \dots, y_n) = \prod_i \theta_i^{\mathbb{I}\{y_i=1\}} (1 - \theta_i)^{\mathbb{I}\{y_i=-1\}}$$

$$\theta_i = \frac{1}{1 + e^{-\mathbf{x}_i^T \mathbf{w}}}$$

$$y_i = 1$$

$$\mathbf{p}_i = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$y_i = -1$$

$$\mathbf{p}_i = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$\mathbf{q}_i = \begin{bmatrix} \theta_i \\ 1 - \theta_i \end{bmatrix}$$

# Training

$$\theta_i = \frac{1}{1 + e^{-\mathbf{x}_i^T \mathbf{w}}}$$

- maximum likelihood principle

- notation

$$p_{\mathbf{w}}(y_1, y_2, \dots, y_n) = \prod_i \theta_i^{\mathbb{I}\{y_i=1\}} (1 - \theta_i)^{\mathbb{I}\{y_i=-1\}}$$

$$y_i = 1 \quad \mathbf{p}_i = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$\log L(\mathbf{w}) = \sum_i \mathbb{I}_{\{y_i=1\}} \log \theta_i + \sum_i \mathbb{I}_{\{y_i=-1\}} \log(1 - \theta_i)$$

$$y_i = -1 \quad \mathbf{p}_i = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$\mathbf{q}_i = \begin{bmatrix} \theta_i \\ 1 - \theta_i \end{bmatrix}$$

$$= - \sum_i H(\mathbf{p}_i, \mathbf{q}_i)$$

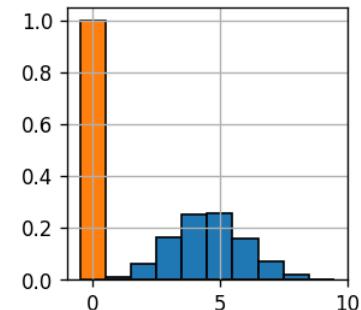
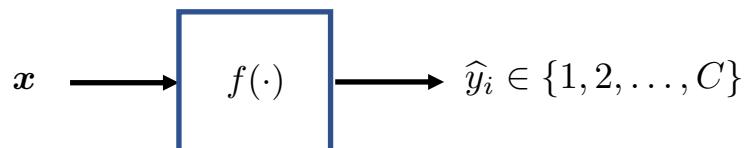
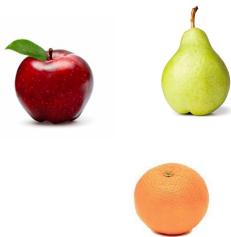
```
model.compile(optimizer='sgd', loss=tf.keras.losses.CategoricalCrossentropy())
```

$$\text{nll}(\mathbf{w}) = \sum_i H(\mathbf{p}_i, \mathbf{q}_i)$$

$$= \sum_i (H(\mathbf{p}_i) + D(\mathbf{p}_i || \mathbf{q}_i))$$

$$= \sum_i D(\mathbf{p}_i || \mathbf{q}_i)$$

# Multiclass Classification



$$\mathcal{D} = \{(\mathbf{x}_i, y)\}_{i=1}^n \longrightarrow \mathcal{D} = \{(\mathbf{x}_i, \mathbf{p}_i)\}_{i=1}^n \quad \mathbf{p}_i = \text{vector of zeros with a 1 in the } y_i \text{ location}$$

- example:

$$\mathbf{x}_i = \begin{matrix} \text{apple image} \end{matrix} \longrightarrow \boxed{f(\cdot)} \longrightarrow \mathbf{q}_i = \begin{bmatrix} 0.92 \\ 0.01 \\ \vdots \\ 0.03 \end{bmatrix}$$
$$\mathbf{p}_i = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$
A diagram showing a feature vector  $\mathbf{x}_i$  (represented by an apple image) entering a function block  $f(\cdot)$ , which outputs a probability vector  $\mathbf{q}_i$ . Below it, the corresponding one-hot encoding  $\mathbf{p}_i$  is shown as a column vector with a 1 at the first position and zeros elsewhere.

- loss of  $i$ th training example:

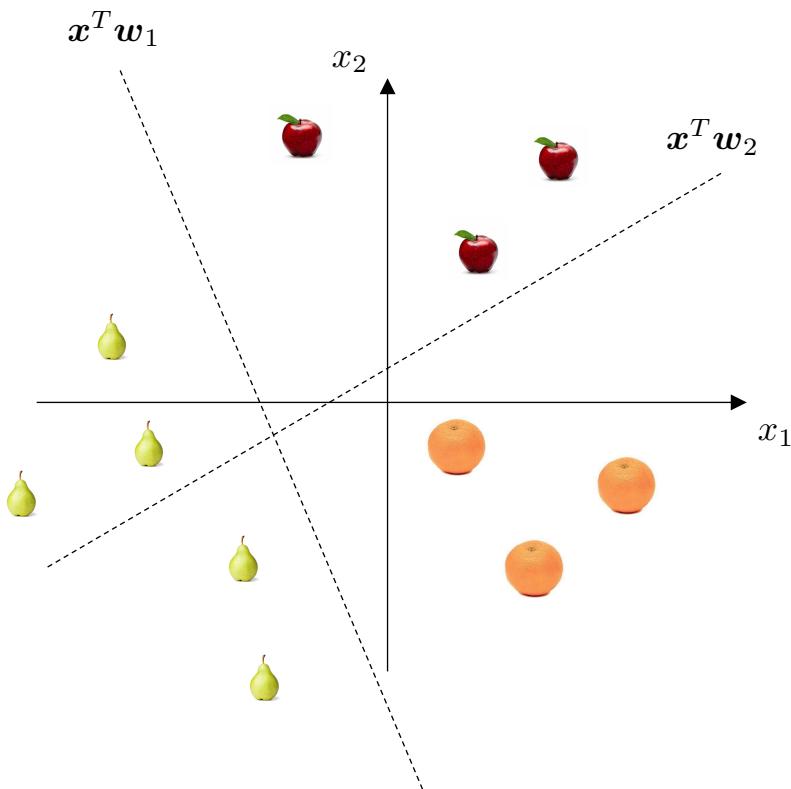
$$\ell(\hat{y}, y) = D(\mathbf{p}_i || \mathbf{q}_i)$$

# Multinomial Logistic Regression

- logistic regression is a *discriminative* classifier

directly fit  $p(y|\mathbf{x})$  to data by treating  $\mathbf{x}$  as non-random

$$p(y = c|\mathbf{x}_i) = \frac{e^{\mathbf{x}_i^T \mathbf{w}_c}}{\sum_{c'=1}^C e^{\mathbf{x}_i^T \mathbf{w}_{c'}}}$$



- $\mathbf{x}_i^T \mathbf{w}_c$  is a **signed and scaled distance** from  $\{\mathbf{x} : \mathbf{x}_c^T \mathbf{w} = 0\}$
- softmax function takes  $\mathbf{x}^T \mathbf{w}_c, c = 1, \dots, C$  and squishes distance to make a probability
- classification rule:
$$\hat{y} = \arg \max_y p(y|\mathbf{x})$$
- boundaries are linear
$$\frac{p(y = c|\mathbf{x}_i)}{p(y = c'|\mathbf{x}_i)} = \frac{e^{\mathbf{x}_i^T \mathbf{w}_c}}{e^{\mathbf{x}_i^T \mathbf{w}_{c'}}} \gtrless 1$$
$$\mathbf{x}_i^T (\mathbf{w}_c - \mathbf{w}_{c'}) \gtrless 0$$
- convention: set  $\mathbf{w}_C = \mathbf{0}$  to ensure identifiability

# Multinomial Logistic Regression

- logistic regression is a *discriminative* classifier

$$y_i = \begin{cases} 1 & \text{with probability } \theta_i \\ -1 & \text{with probability } 1 - \theta_i \end{cases} \quad \theta_i = \frac{1}{1 + e^{-\mathbf{x}_i^T \mathbf{w}}} = \frac{e^{\mathbf{x}_i^T \mathbf{w}}}{e^{\mathbf{x}_i^T \mathbf{w}} + 1}$$



- multiclass (multinomial) logistic regression:

$$y_i = \begin{cases} 1 & \text{with probability } \theta_{i,1} \\ 2 & \text{with probability } \theta_{i,2} \\ \vdots \\ c & \text{with probability } \theta_{i,C} \end{cases} \quad \theta_{i,c} = p(y = c | \mathbf{x}_i) = \frac{e^{\mathbf{x}_i^T \mathbf{w}_c}}{\sum_{c'=1}^C e^{\mathbf{x}_i^T \mathbf{w}_{c'}}$$

- softmax:

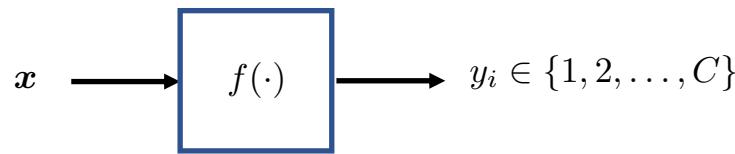
$$[\sigma(\mathbf{z})]_i = \frac{e^{z_i}}{\sum_j e^{z_j}}$$

$$\sigma(\mathbf{z}) = \frac{e^{\mathbf{z}}}{\sum_j e^{z_j}}$$

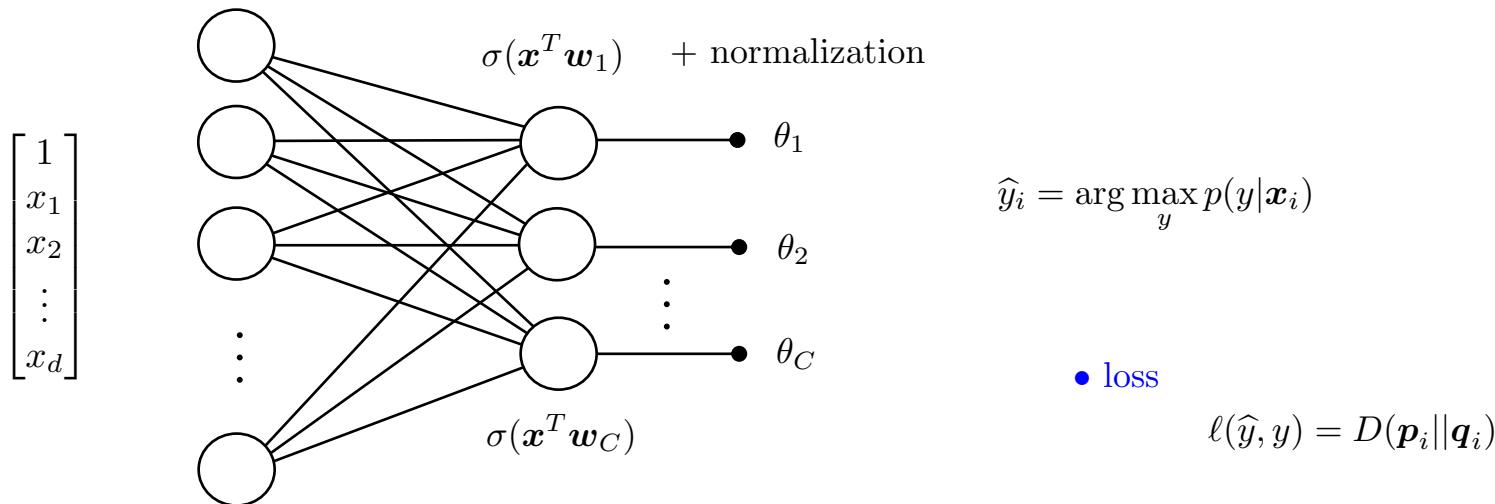
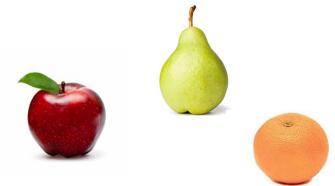
- softmax example:

$$\begin{bmatrix} [-18.77687567] \\ [39.26910416] \\ [50.56209096] \\ [47.72731802] \\ [34.95432764] \end{bmatrix} \xrightarrow{\hspace{1cm}} \begin{bmatrix} [0. & \dots] \\ [0. & \dots] \\ [0.9445] \\ [0.0555] \\ [0. & \dots] \end{bmatrix}$$

# Multinomial Logistic Regression



- testing
  - single *single layer of artificial neural network*

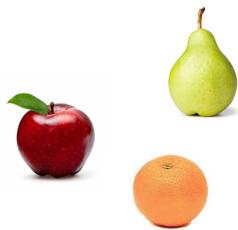


- training

$$\mathcal{D} = \{(\mathbf{x}_i, y)\}_{i=1}^n$$

$$p_{\mathbf{w}}(y_1, y_2, \dots, y_n) = \prod_i \theta_{i,1}^{\mathbb{I}\{y_i=1\}} \theta_{i,2}^{\mathbb{I}\{y_i=2\}} \dots \theta_{i,C}^{\mathbb{I}\{y_i=C\}}$$

# Training



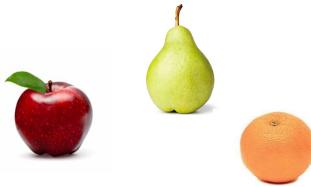
- training: how do we find a good  $\mathbf{w}$ ?
- maximum likelihood principle

- given data (and our model) what's the most likely  $\mathbf{w}$ ?

$$\begin{aligned}
 p_{\mathbf{W}}(y_1, y_2, \dots, y_n) &= \prod_{i=1}^n \prod_{c=1}^C \theta_{i,c}^{\mathbb{I}\{y_i=c\}} \\
 &= \prod_{i=1}^n \prod_{c=1}^C \left( \frac{e^{\mathbf{w}_c^T \mathbf{x}_i}}{\sum_{c'=1}^C e^{\mathbf{w}_{c'}^T \mathbf{x}_i}} \right)^{\mathbb{I}\{y_i=c\}} \\
 \log p_{\mathbf{W}} &= \sum_{i=1}^n \sum_{c=1}^C \mathbb{I}\{y_i=c\} \log \left( \frac{e^{\mathbf{w}_c^T \mathbf{x}_i}}{\sum_{c'=1}^C e^{\mathbf{w}_{c'}^T \mathbf{x}_i}} \right)
 \end{aligned}$$

$$= - \sum_i H(\mathbf{p}_i, \mathbf{q}_i) \quad \mathbf{p}_i = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \mathbf{q}_i = \begin{bmatrix} e^{\mathbf{w}_1^T \mathbf{x}_i} \\ e^{\mathbf{w}_2^T \mathbf{x}_i} \\ \vdots \\ e^{\mathbf{w}_C^T \mathbf{x}_i} \end{bmatrix} \times \frac{1}{\sum_{c'} e^{\mathbf{w}_{c'}^T \mathbf{x}_i}}$$

# Training



- training: how do we find a good  $\mathbf{w}$ ?
- maximum likelihood principle

$$p(y = c | \mathbf{x}_i) = \frac{e^{\mathbf{w}_c^T \mathbf{x}_i}}{\sum_{c'=1}^C e^{\mathbf{w}_{c'}^T \mathbf{x}_i}}$$

- given data (and our model) what's the most likely  $\mathbf{w}$ ?

$$p_{\mathbf{W}}(y_1, y_2, \dots, y_n) = \prod_{i=1}^n \prod_{c=1}^C \theta_{i,c}^{\mathbb{I}\{y_i=c\}}$$

$$= \prod_{i=1}^n \prod_{c=1}^C \left( \frac{e^{\mathbf{w}_c^T \mathbf{x}_i}}{\sum_{c'=1}^C e^{\mathbf{w}_{c'}^T \mathbf{x}_i}} \right)^{\mathbb{I}\{y_i=c\}}$$

$$\log p_{\mathbf{W}} = \sum_{i=1}^n \sum_{c=1}^C \left( \mathbb{I}\{y_i=c\} \mathbf{w}_c^T \mathbf{x}_i - \mathbb{I}\{y_i=c\} \log \left( \sum_{c'=1}^C e^{\mathbf{w}_{c'}^T \mathbf{x}_i} \right) \right)$$

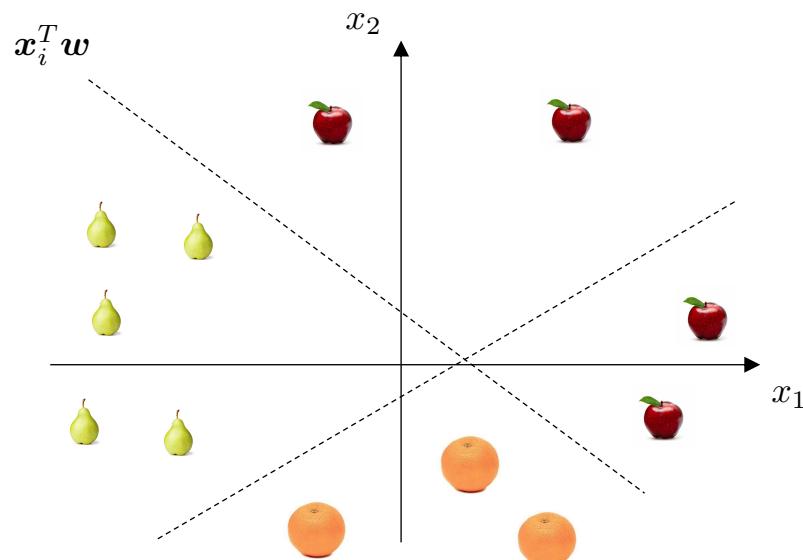
$$= \sum_{i=1}^n \left( \left( \sum_{c=1}^C \mathbb{I}\{y_i=c\} \mathbf{w}_c^T \mathbf{x}_i \right) - \log \left( \sum_{c'=1}^C e^{\mathbf{w}_{c'}^T \mathbf{x}_i} \right) \right)$$

$$= \sum_{i=1}^n \left( \mathbf{w}_{y_i}^T \mathbf{x}_i - \log \left( \sum_{c'=1}^C e^{\mathbf{w}_{c'}^T \mathbf{x}_i} \right) \right)$$

# Optimizing

- training reduces to solving the following optimization

$$\arg \max_{\mathbf{w}_1, \dots, \mathbf{w}_c} \sum_{i=1}^n \left( \mathbf{w}_{y_i}^T \mathbf{x}_i - \log \left( \sum_{c'=1}^C e^{\mathbf{w}_{c'}^T \mathbf{x}_i} \right) \right)$$



- changing  $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_c$  changes decision boundaries and steepness of logistic function
- no closed form solution
- iterative methods for finding  $\mathbf{w}$
- concave and differentiable

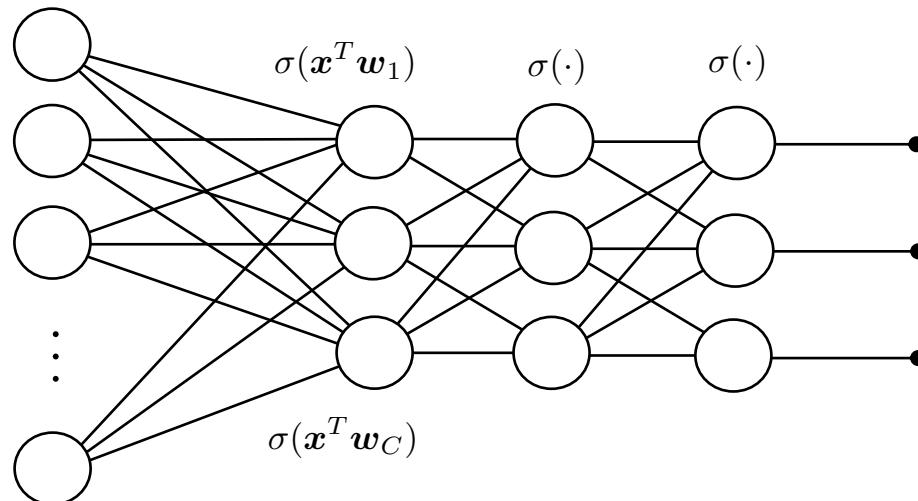
# Computing the Gradient

$$\arg \max_{\mathbf{w}_1, \dots, \mathbf{w}_c} \sum_{i=1}^n \left( \mathbf{w}_{y_i}^T \mathbf{x}_i - \log \left( \sum_{c'=1}^C e^{\mathbf{w}_{c'}^T \mathbf{x}_i} \right) \right)$$

$$\nabla_{\mathbf{w}_c} = \sum_i \left( \mathbb{I}_{\{y_i=c\}} - \frac{e^{\mathbf{w}_c^T \mathbf{x}_i}}{\sum_{c'} e^{\mathbf{w}_{c'}^T \mathbf{x}_i}} \right) \mathbf{x}_i$$

# Artificial Neural Networks and Logistic Regression

- binary logistic regression:
  - single *artificial neuron*
- multiclass logistic regression:
  - single layer *artificial neural network*
- neural network:
  - important notes:
    - with stacked layers, objective not convex
    - decision boundaries are no longer linear
    - can be used for regression



# Differential Entropy

- definition - entropy:

$$H(X) = E \left[ \log_2 \left( \frac{1}{p(x)} \right) \right]$$

- discrete  $x \in \mathcal{X}$ :

$$H(X) = \sum_x p(x) \log_2 \left( \frac{1}{p(x)} \right)$$

- amount of uncertainty in a random variable

- *differential entropy* is the equivalent definition of entropy, for continuous random variables

- continuous  $x \in \mathbb{R}$ :

$$H(X) = \int_x p(x) \log_2 \left( \frac{1}{p(x)} \right) dx$$

- example:  $X \sim \mathcal{N}(\mu, \sigma^2)$

$$H(X) = \frac{1}{2} \log_2(2\pi e \sigma^2) \text{ bits}$$