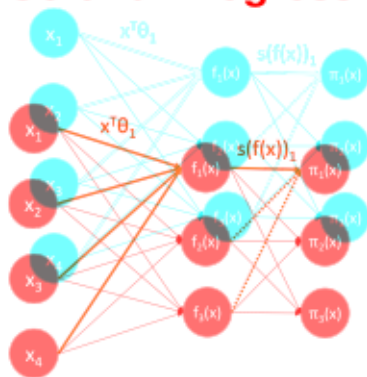


Introduction to Machine Learning



Multiclass Classification

Softmax Regression



Learning goals

- Know softmax regression

Learning goals

- Understand that softmax regression is a generalization of logistic regression
- Understand that softmax regression is a generalization of logistic regression

... TO SOFTMAX REGRESSION / 2

- The softmax function is a generalization of the logistic function. For $g = 2$, the logistic function and the softmax function are equivalent.
- Instead of the **Bernoulli** loss, we use the multiclass **logarithmic loss**

$$L(y, \pi(\mathbf{x})) = - \sum_{k=1}^g \mathbb{1}_{\{y=k\}} \log(\pi_k(\mathbf{x})).$$

- Note that the softmax function is a “smooth” approximation of the arg max operation, so $s((1, 1000, 2)^T) \approx (0, 1, 0)^T$ (picks out 2nd element!).
- Furthermore, it is invariant to constant offsets in the input:

$$s(f(\mathbf{x}) + \mathbf{c}) = \frac{\exp(\boldsymbol{\theta}_k^T \mathbf{x} + c)}{\sum_{j=1}^g \exp(\boldsymbol{\theta}_j^T \mathbf{x} + c)} = \frac{\exp(\boldsymbol{\theta}_k^T \mathbf{x}) \cdot \exp(c)}{\sum_{j=1}^g \exp(\boldsymbol{\theta}_j^T \mathbf{x}) \cdot \exp(c)} = s(f(\mathbf{x}))$$



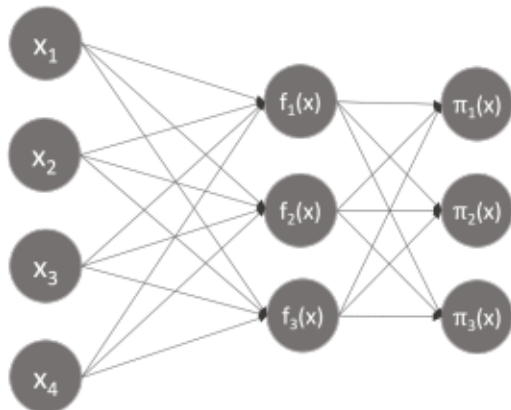
LOGISTIC VS. SOFTMAX REGRESSION



| | Logistic Regression | Softmax Regression |
|-------------------------|--|--|
| y | $\{0, 1\}$ | $\{1, 2, \dots, g\}$ |
| \hat{y} | $\{0, 1\}$ | $\{1, 2, \dots, g\}$ |
| Discriminant fun. | $f(\mathbf{x}) = \theta^T \mathbf{x}$ | $f_k(\mathbf{x}) = \theta_k^T \mathbf{x}, k = 1, 2, \dots, g$ |
| Discriminant fun. | $f(\mathbf{x}) = \theta^T \mathbf{x}$ | $f_k(\mathbf{x}) = \theta_k^T \mathbf{x}, k = 1, 2, \dots, g$ |
| Probabilities | $\pi(\mathbf{x}) = \frac{1}{1 + \exp(-\theta^T \mathbf{x})}$ | $\pi_k(\mathbf{x}) = \frac{\exp(\theta_k^T \mathbf{x})}{\sum_{j=1}^g \exp(\theta_j^T \mathbf{x})}$ |
| Probabilities | $\pi(\mathbf{x}) = \frac{1}{1 + \exp(-\theta^T \mathbf{x})}$ | $\pi_k(\mathbf{x}) = \frac{\exp(\theta_k^T \mathbf{x})}{\sum_{j=1}^g \exp(\theta_j^T \mathbf{x})}$ |
| $L(y, \pi(\mathbf{x}))$ | Bernoulli / logarithmic loss | Multiclass logarithmic loss |
| $L(y, \pi(\mathbf{x}))$ | $-y \log(\pi(\mathbf{x})) - (1-y) \log(1-\pi(\mathbf{x}))$ | $-\sum_{k=1}^g [y = k] \log(\pi_k(\mathbf{x}))$ |
| $L(y, \pi(\mathbf{x}))$ | $-y \log(\pi(\mathbf{x})) - (1-y) \log(1-\pi(\mathbf{x}))$ | $-\sum_{k=1}^g [y = k] \log(\pi_k(\mathbf{x}))$ |

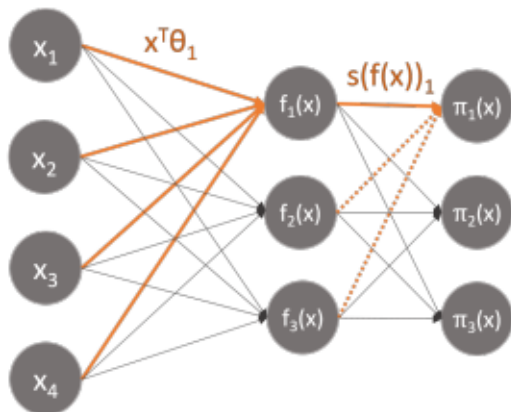
LOGISTIC VS. SOFTMAX REGRESSION

We can schematically depict softmax regression as follows:



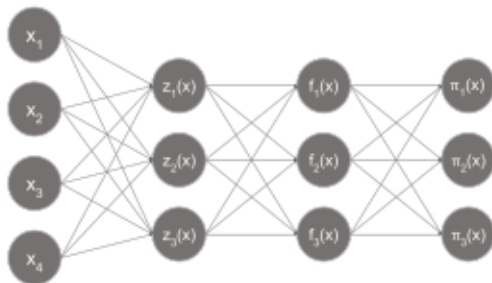
LOGISTIC VS. SOFTMAX REGRESSION

We can schematically depict softmax regression as follows:



GENERALIZING SOFTMAX REGRESSION / 2

For example for a **neural network** (note that softmax regression is also a neural network with no hidden layers):



Remark: For more details about neural networks please refer to the lecture **Deep Learning**.