

8. Supervised Learning: Softmax Regression

Softmax Regression: Model

Setting

We consider a K-class classification problem where the output value $y \in \{1, 2, \dots, K\}$ and the input variable $\mathbf{x} \in \mathbb{R}^{m+1}$.

We also consider that we have access to labelled data $(\mathbf{x}^{(i)}, y^{(i)})$, $i = 1, \dots, n$, with $y^{(i)} \in \{1, 2, \dots, K\}$ and $\mathbf{x}^{(i)} = [1 \quad x_1^{(i)} \quad x_2^{(i)} \quad \dots \quad x_m^{(i)}]$

We want to estimate $Pr(y|\mathbf{x})$ for $y \in \{1, 2, \dots, K\}$ given $\mathbf{x} = [1 \quad x_1 \quad x_2 \quad \dots \quad x_m] \in \mathbb{R}^{m+1}$

Model

The hypothesis is such that:

$$h_{\theta}(\mathbf{x}) = \begin{bmatrix} Pr(y = 1|\mathbf{x}; \boldsymbol{\theta}_1) \\ \dots \\ Pr(y = K|\mathbf{x}; \boldsymbol{\theta}_K) \end{bmatrix} = \frac{1}{\sum_{k=1}^K e^{\boldsymbol{\theta}_k^t \cdot \mathbf{x}}} \cdot \begin{bmatrix} e^{\boldsymbol{\theta}_1^t \cdot \mathbf{x}} \\ \dots \\ e^{\boldsymbol{\theta}_K^t \cdot \mathbf{x}} \end{bmatrix}$$

where $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_K \in \mathbb{R}^{m+1}$ are the parameters of the model

Softmax Regression: Cost Function

How does the learning algorithm select the linear hypothesis / model? We need a cost function...

Cost Function

The softmax cost function is given by:

$$J(\theta_1, \dots, \theta_K) = -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K 1\{y^{(i)} = k\} \cdot \log_2 \frac{e^{\theta_k^t \cdot x^{(i)}}}{\sum_{j=1}^K e^{\theta_j^t \cdot x^{(i)}}}$$

where $1\{\cdot\}$ is the indicator function given by: $1\{\text{true statement}\} = 1$ and $1\{\text{false statement}\} = 0$.

Why this cost function?

This cost function generalizes the logistic regression cost function.

Softmax Regression: Learning Algorithm



Learning Algorithm

We should select the logistic regression model parameters that minimize the cost function as follows:

$$\theta_1^*, \dots, \theta_K^* = \underset{\theta_1, \dots, \theta_K}{\operatorname{argmin}} J(\theta_1, \dots, \theta_K)$$

New Predictions

The new predictions are then given by:

$$\begin{bmatrix} \Pr(y = 1 | \mathbf{x}; \theta_1^*) \\ \dots \\ \Pr(y = K | \mathbf{x}; \theta_K^*) \end{bmatrix} = \frac{1}{\sum_{k=1}^K e^{\theta_k^{*t} \cdot \mathbf{x}}} \cdot \begin{bmatrix} e^{\theta_1^{*t} \cdot \mathbf{x}} \\ \dots \\ e^{\theta_K^{*t} \cdot \mathbf{x}} \end{bmatrix}$$

There is also no closed-form solution for the optimal parameters



Gradient descent algorithm

Softmax Regression: Learning Algorithm



Considerations

- The softmax regression has a redundant set of parameters because the hypothesis is not affected by subtracting some fixed vector ψ to each of our parameter vectors $\theta_1, \dots, \theta_K$ i.e.

$$Pr(y = k|\mathbf{x}) = \frac{e^{(\theta_k^t - \psi) \cdot \mathbf{x}}}{\sum_{j=1}^K e^{(\theta_j^t - \psi) \cdot \mathbf{x}}} = \frac{e^{\psi \cdot \mathbf{x}} \cdot e^{\theta_k^t \cdot \mathbf{x}}}{e^{\psi \cdot \mathbf{x}} \cdot \sum_{j=1}^K e^{\theta_j^t \cdot \mathbf{x}}} = \frac{e^{\theta_k^t \cdot \mathbf{x}}}{\sum_{j=1}^K e^{\theta_j^t \cdot \mathbf{x}}}$$

- Therefore, the softmax regression model is overparameterized in the sense that there are multiple parameter settings that give rise to exactly the same hypothesis function.
- The cost function is convex so gradient descent will not run into local optima problem
- The cost function Hessian is non-invertible so Newton's methods will run into numerical problems