

Week 04 - Part 02

So far: We talked about Logistic Regression and linear classification for binary classes.

Today: Multi-class Logistic Regression

Multi-class Logistic Regression

■ Label: ...

■ Hypothesis Set: Let $\Omega = \{ \underline{w}_{(1)}, \underline{w}_{(2)}, \dots, \underline{w}_{(c)} \}$ be the weight vectors for c classes.

□ Hypothesize that

$$P[y_n = i | \underline{x}_n] = \underline{\hspace{2cm}}$$

■ Error Criterion:

■ To minimize $E_{in}(\Omega)$, we need to use GD. Hence, we must find the gradient.

$$\nabla_{\Omega} e_n(\Omega) =$$

$=$

SGD update:

In each iteration t :

Computing $\nabla_{\underline{w}(i)} \ell_n(\Omega_t)$

$$\begin{aligned} \blacksquare \ell_n(\Omega_t) &= -\log \hat{P}_{\Omega_t}(y_n | \underline{x}_n) = -\log \frac{e^{\underline{w}^T(y_n) \underline{x}_n}}{\sum_{j=1}^c e^{\underline{w}^T(j) \underline{x}_n}} \\ &= -\underline{w}^T(y_n) \underline{x}_n + \log \left(\sum_{j=1}^c e^{\underline{w}^T(j) \underline{x}_n} \right) \end{aligned}$$

□ For $i = y_n$,

□ For $i \neq y_n$,

$$\nabla_{\underline{w}(i)} \ell_n(\Omega_t) = \nabla_{\underline{w}(i)} \left(-\underline{w}(y_n)^T \underline{x}_n + \log \left(\sum_{j=1}^S e^{\underline{w}(j)^T \underline{x}_n} \right) \right)$$

=

Softmax Logistic Regression for Binary Classification (C=2)

■ What is the relation between Softmax logistic regression for C=2 and binary logistic regression that we studied last week.

$$\square \hat{P}_{\Omega}(1 | \underline{x}_n) = \frac{e^{\underline{w}^{(1)T} \underline{x}_n}}{e^{\underline{w}^{(1)T} \underline{x}_n} + e^{\underline{w}^{(2)T} \underline{x}_n}} = \frac{e^{(\underline{w}^{(1)} - \underline{w}^{(2)})^T \underline{x}_n}}{1 + e^{(\underline{w}^{(1)} - \underline{w}^{(2)})^T \underline{x}_n}}$$

$$\square \hat{P}_{\Omega}(2 | \underline{x}_n) = \frac{e^{\underline{w}^{(2)T} \underline{x}_n}}{e^{\underline{w}^{(1)T} \underline{x}_n} + e^{\underline{w}^{(2)T} \underline{x}_n}} = 1 - \hat{P}_{\Omega}(1 | \underline{x}_n)$$

□ It is logistic Regression with $\underline{w} = \underline{w}^{(1)} - \underline{w}^{(2)}$

Can we use GD/SGD for Linear Regression? Yes, you can

■ $E_{in}(\underline{w}) = \frac{1}{N} \sum_{n=1}^N e_n(\underline{w}) = \frac{1}{N} \sum_{n=1}^N (\underline{w}^T \underline{x}_n - y_n)^2.$

■ \underline{w}_{ls} , the optimal solution.

□ $\underline{w}_{ls} =$

□ Complexity:

■ $E_{in}(\underline{w})$ is convex. So, with GD/SGD, as $t \rightarrow \infty$, \underline{w}_t converges to the optimal solution, i.e., \underline{w}_{ls} .

□ $\nabla_{\underline{w}} e_n(\underline{w}) =$

□ Time Complexity of SGD, Full-GD, Min-batch GD:

● SGD:

● Full GD:

● Mini GD:

■ Benefits of GD over Least-squares method:

1-

2-

