

Logistic Regression

$$P(y=1 | \underline{x}, \underline{w}) = \sigma(\underline{w}^T \underline{x})$$

Match $y^{(n)}$ to $\sigma(\underline{w}^T \underline{x}^{(n)})$

Negative log likelihood:

$$NLL = - \sum_{n=1}^N \log \sigma(\underline{z}^{(n)} \underline{w}^T \underline{x}^{(n)})$$

↑

Label $\in \{-1, +1\}$

$$\underline{z} = 2y - 1$$

$$\nabla_{\underline{w}}^{NLL} = - \sum_{n=1}^N (1 - \sigma_n) \underline{z}^{(n)} \underline{x}^{(n)}$$

Stochastic Gradient Descent (SGD)

Average gradient

$$\frac{1}{N} \nabla_{\underline{w}} NLL = \frac{1}{N} \sum_{n=1}^N \underline{g}_n$$

eg. $-(1-\sigma_n) z^n \underline{x}^{(n)}$

Approximation

Batch of B examples at random

$$\approx \frac{1}{B} \sum_{b=1}^B \underline{g}_b$$

Stochastic / Online g.d.

uses a batch with $B=1$

"Minibatch" uses $B \approx 64, 128$



contour of
cost f^n

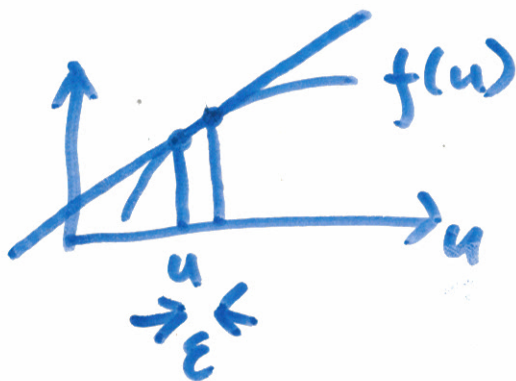
I quoted §10.6

"Numerical Recipes in Fortran"

(2nd Ed.)

Check your Derivatives

$$f'(u) \approx \frac{f(u+\varepsilon) - f(u)}{\varepsilon}$$



Error $O(\varepsilon)$

$$\varepsilon \approx 10^{-5}$$

$$\approx \frac{f(u+\varepsilon/2) - f(u-\varepsilon/2)}{\varepsilon}$$

Error $O(\varepsilon^2)$

Softmax Regression

Fit one-hot-encoded labels

Target $y = [0 \ 0 \ 0 \ \dots \ 0 \ 1 \ 0 \ \dots \ 0]$

↑
cth location
indicating class c

Fit f . Our model:

$$P(\text{class } c) = f_c$$

Each f_k to be positive:

$$f_k \propto e^{\underline{w}^{(k)T} \underline{x}}$$

We want f to be normalized

$$\sum_k f_k = 1:$$

$$f_k = \frac{e^{\underline{w}^{(k)T} \underline{x}}}{\sum_{k'} e^{\underline{w}^{(k')T} \underline{x}}}$$

])
Softmax

Model has params

$$W = \{ \underline{w}^{(k)} \}$$

Max. the likelihood of W

In one example at \underline{x} and ^{we} see that it has class c :

$$\log f_c = \underline{w}^{(c)T} \underline{x} - \log \sum_{k'} e^{\underline{w}^{(k')}T \underline{x}}$$

$$\nabla_{\underline{w}^{(k)}} \log f_c = \delta_{ck} \underline{x} - \frac{1}{\sum_{k'} \dots} e^{\underline{w}^{(k)T} \underline{x}} \underline{x}$$

↑
Kronecker delta

$$= (y_k - f_k) \underline{x}$$

Logistic Regression?

Two classes

$$\begin{aligned} P(y=1 | \underline{x}, \underline{w}) &= \frac{e^{\underline{w}^{(1)T} \underline{x}}}{e^{\underline{w}^{(1)T} \underline{x}} + e^{\underline{w}^{(2)T} \underline{x}}} \\ &= \frac{1}{1 + e^{(\underline{w}^{(2)T} - \underline{w}^{(1)T}) \underline{x}}} \\ &= \sigma(\underbrace{(\underline{w}^{(1)} - \underline{w}^{(2)})^T \underline{x}}_{\text{"w"}}) \end{aligned}$$