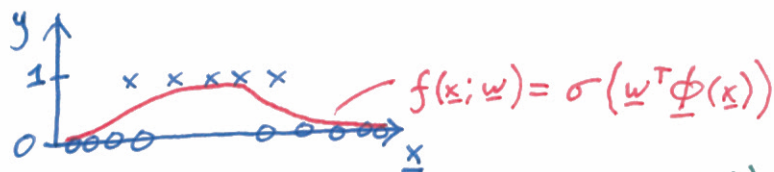


Logistic Regression

for $y \in \{0, 1\}$

Model the outputs: $P(y=1 | \underline{x}, \underline{w}) = f(\underline{x}; \underline{w}) = \sigma(\underline{w}^T \underline{x})$

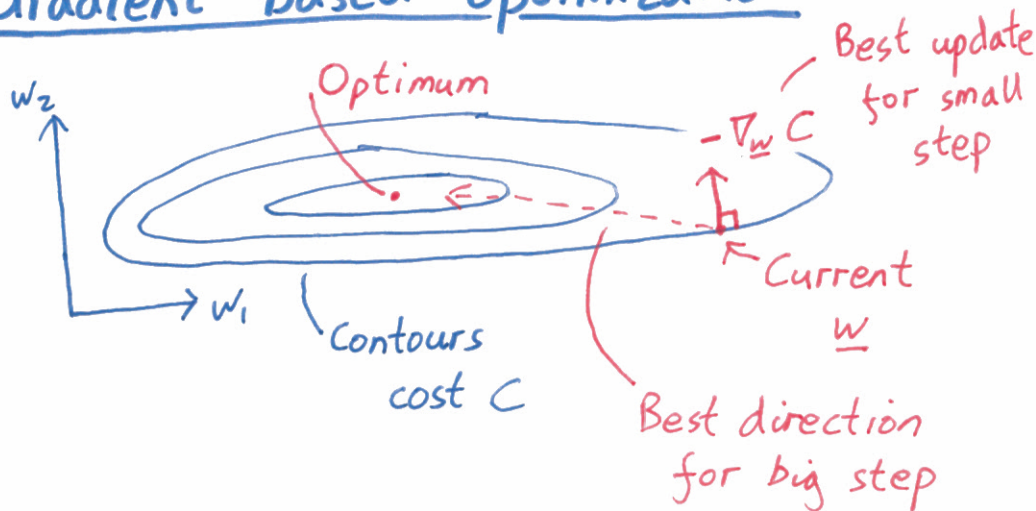


Maximum Likelihood (minimize -ve log likelihood)

$$NLL = -\sum_n \log P(y^{(n)} | \underline{x}^{(n)}, \underline{w}) = -\sum_n \log \underbrace{\sigma(z^{(n)} \underline{w}^T \underline{x}^{(n)})}_{\sigma_n}^{2y^{(n)} - 1}$$

$$\begin{aligned} \nabla_{\underline{w}} NLL &= -\sum_n \nabla_{\underline{w}} \log \sigma_n && \left. \begin{array}{l} \frac{d \log a}{da} = \frac{1}{a} + \text{chain rule} \\ \frac{d \sigma(a)}{da} = \sigma(a)(1 - \sigma(a)) \end{array} \right\} \\ &= -\sum_n \frac{1}{\sigma_n} \nabla_{\underline{w}} \sigma_n \\ &= -\sum_n \frac{1}{\sigma_n} \sigma_n (1 - \sigma_n) \nabla_{\underline{w}} z^{(n)} \underline{w}^T \underline{x}^{(n)} \\ &= -\sum_n \underbrace{(1 - \sigma_n)}_{\sigma(-z^{(n)} \underline{w}^T \underline{x}^{(n)})} z^{(n)} \underline{x}^{(n)} = -\sum_n (y^{(n)} - \underbrace{f^{(n)}}_{\sigma(\underline{w}^T \underline{x}^{(n)})}) \underline{x}^{(n)} \end{aligned}$$

Gradient-based optimization



Finding better directions

- Non-linear conjugate gradients
- L-BFGS
- Newton's method ...

Often "batch" methods:

C and $\nabla_{\underline{w}} C$ use whole data set

Stochastic / Online or "Minibatch" use data subsets

Stochastic Gradient Descent

Average Gradient over examples: $-(1-\sigma_n)z^{(n)}x^{(n)}$

$$\frac{1}{N} \nabla_{\underline{w}} C = \frac{1}{N} \sum_{n=1}^N \nabla_{\underline{w}} C_n$$

Cost for n^{th} example.

Monte Carlo Approximation:

Sample mini-batch of B examples:

$$\approx \frac{1}{B} \sum_{b=1}^B \nabla_{\underline{w}} C_b = \hat{\underline{g}}$$

SGD

Initialize $\underline{w} \leftarrow \underline{w}^{(0)}$ (maybe $\underline{0}$)

for $t = 1 \dots T$:

$$\underline{w} \leftarrow \underline{w} - \eta_t \hat{\underline{g}}^{(t)}$$

step-size

Softmax Regression

Multi-class classification

Target $\underline{y}^{(n)} = [0 \ 0 \ 0 \ \dots \ 0 \ 1 \ 0 \ \dots \ 0]^T$
 \uparrow c^{th} location

If example n has label " $y=c$ "

Fit \underline{f} , predictions from our model:

$$P(\text{class} = c \mid \underline{x}, W) = f_c(\underline{x}; W)$$

Positive score for each class:

$$s_k = e^{\underline{w}^{(k)T} \underline{x}}$$

Want \underline{f} to be normalized: $\sum_k f_k = 1$

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

Softmax($W\underline{x}$)

Model has parameters W
 $K \times D$ \rightarrow # features

$$W = \begin{pmatrix} \text{---} & \underline{w^{(1)}} & \text{---} \\ & \vdots & \\ \text{---} & \underline{w^{(k)}} & \text{---} \end{pmatrix}$$

Maximize likelihood of \underline{w}

For SGD we use one example at \underline{x}
with class label c

- log prob. of this example given W

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

$$- \nabla_{\underline{w}^{(k)}} \log f_c = - \underbrace{\delta_{kc}}_{\text{Kronecker Delta}} \underline{x} + \underbrace{\frac{1}{\sum_{k'} \dots}}_{f_k} e^{\underline{w}^{(k)T} \underline{x}} \underline{x}$$

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

Logistic Regression

Two classes

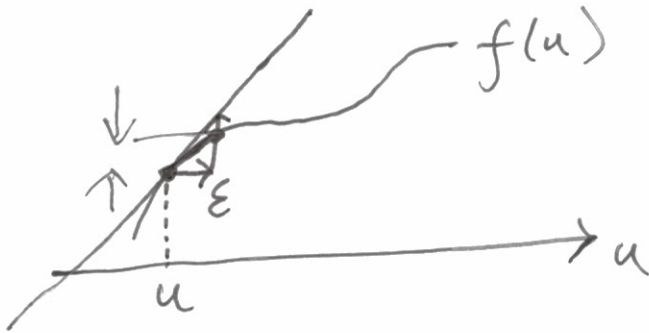
$$P(y=1 | \underline{x}, W) = \frac{e^{\underline{w}^{(1)T} \underline{x}}}{e^{\underline{w}^{(1)T} \underline{x}} + e^{\underline{w}^{(0)T} \underline{x}}}$$
$$= \frac{1}{1 + e^{\underbrace{\underline{w}^{(0)T} \underline{x} - \underline{w}^{(1)T} \underline{x}}_{\text{"-a"}}}}$$

$$\left[\sigma(a) = \frac{1}{1 + e^{-a}} \right]$$

$$= \sigma\left(\underbrace{(\underline{w}^{(1)} - \underline{w}^{(0)})^T \underline{x}}_{\text{"w"}}$$

Parameters are redundant.

Check your derivatives



$$\left. \frac{df}{du} \right|_u = f'(u) \approx \frac{f(u+\epsilon) - f(u)}{\epsilon}$$

correct $\lim \epsilon \rightarrow 0$

approx $\epsilon = 10^{-5}$

