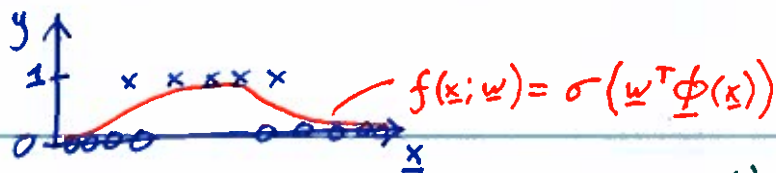


Logistic Regression

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

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

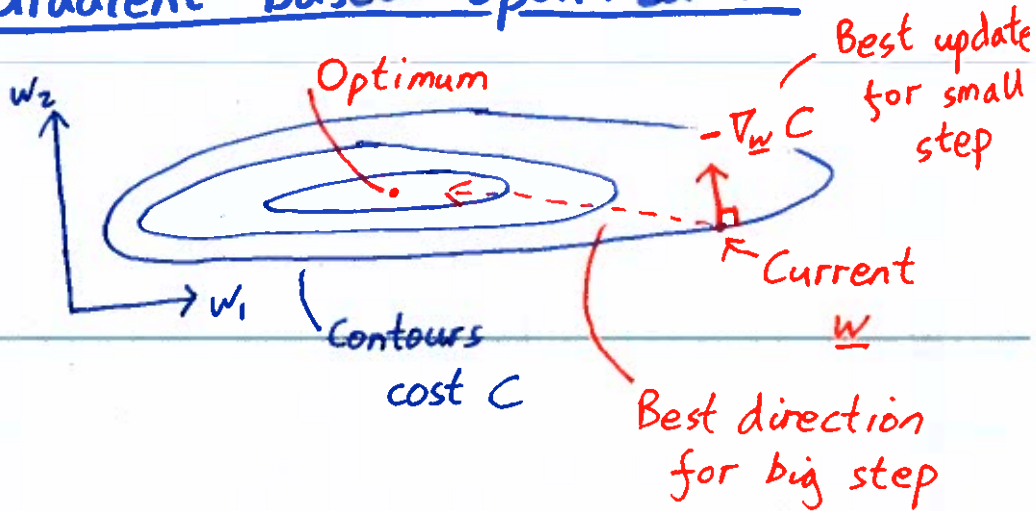


Maximum Likelihood (minimize -ve log likelihood) $2y^{(n)} - 1$

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

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

Gradient-based optimization



Finding better directions

- Non-linear conjugate gradients
 - L-BFGS
 - Newton's method ...
- Details non-examinable
- (might cover later)

Often "batch" methods:

C and $\nabla_w C$ use whole data set

Stochastic / Online or "Minibatch" use data subsets

Stochastic Gradient Descent

Average Gradient over examples:

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

(Cost for n th example)

Monte Carlo Approximation:

Sample mini-batch of B examples:

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

$n_b \sim$ sampled from $\{1 \dots N\}$

SGD

Initialize $\underline{w} \leftarrow \underline{w}^{(0)}$ (0 ok for logistic regression)

for $t = 1 \dots T$:

$$\underline{w} \leftarrow \underline{w} - \underbrace{\eta}_{\text{step-size}} \hat{\underline{g}}^{(t)}$$

gradient est. at time t using current \underline{w}

Softmax Regression

Multi-class classification

$$\text{Target } y^{(n)} = [0 \ 0 \ 0 \ \dots \ 0 \ 1 \ 0 \ \dots \ 0]^T$$

↑ c^{th} location
Example n has label " $y=c$ "

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

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

Positive score for each class:

$$s_k = e^{\underline{w}^{(k)T} \underline{x}} \quad k=1 \dots K$$

↓ # classes

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

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

$$\underline{f}_{K \times 1} = \text{softmax} \left(\underset{K \times D}{W} \underset{D \times 1}{\underline{x}} \right)$$

↑ # features

Model has parameters W
 $K \times D$

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

Maximize likelihood of W

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

- log prob. of this example given W

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

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

f_k

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

Logistic Regression

$k=2$, two classes

$$P(y=1 | \underline{x}, W) = \frac{e^{\underline{w}^{(1)} \cdot \underline{x}}}{e^{\underline{w}^{(1)} \cdot \underline{x}} + e^{\underline{w}^{(0)} \cdot \underline{x}}}$$

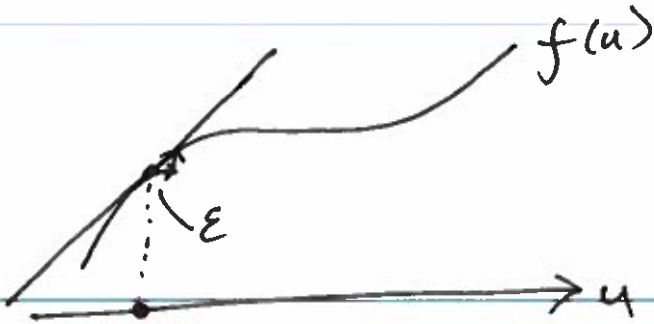
$$= \frac{1}{1 + e^{\underline{w}^{(0)} \cdot \underline{x} - \underline{w}^{(1)} \cdot \underline{x}}}$$

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

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

Parameters are redundant.

Check your derivatives



$$\frac{df}{du} = f'(u) \approx \frac{f(u+\epsilon) - f(u)}{\epsilon}$$

correct $\lim \epsilon \rightarrow 0$

With "doubling" floating point:
 $\epsilon = 10^{-5}$, error $O(\epsilon)$

Notes:
Central difference

$$= \frac{f(u + \epsilon/2) - f(u - \epsilon/2)}{\epsilon}$$

Error $O(\epsilon^2)$

(Non-examinable)

Probability estimation question

Iain is juggling.



Assume he has a 1% chance of dropping each catch.

(not really independent,
but let's pretend)

Estimate without a computer/calculator

the probability he manages

100 catches in a row
on first attempt.