

# Logistic Regression

given  $\{(x_n, y_n)\}_{n=1}^N$ ,  $x_n \in \mathbb{R}^{D_x}$ ,  $y_n \in \mathbb{Z}^{D_y}$

$$f(x) = \sigma(\underbrace{xw + b}_{\text{linear model}}) = \hat{y}$$

activation function for discrete real values

## Logistic sigmoid activation func

$$\sigma = \frac{e^{xw+b}}{e^{xw+b} + 1} = \frac{1}{1 + e^{-(xw+b)}}$$

Sigmoid squashes data between 0 to 1

Bernouli prob density funct = 0 or 1

Likelihood based on Bernouli

$$P(Y|X, w) = \prod_{n=1}^N \underbrace{P(Y|x_n)}_{\text{pred}}^{\underbrace{y_n}_{\text{actual}}} (1 - P(Y|x_n))^{1-y_n}$$
$$= \prod_{n=1}^N f(x_n)^{y_n} (1 - f(x_n))^{1-y_n}$$



Optimizing involves max likelihood of  $P(Y|X, \theta)$

which is the same as min the neg log likelihood (min is always easier)

$$\text{recall: } \frac{\partial}{\partial \omega} (-\log(P(Y|X, \omega))) = 0$$

Sub in the ~~ER formula~~ Bernoulli likelihood formula Distribution & apply logarithm rules

the result is a function of a function ... which also contains the  $\omega$  we want to diff

Solution is to open up the Derivation Function

$$\frac{\partial L}{\partial \omega} = \frac{\partial L}{\partial f(\theta)} \cdot \frac{\partial f(\theta)}{\partial \theta} \cdot \frac{\partial \theta}{\partial \omega}$$



(3)

"Applying rule of derivative of a function of a function"

This is the chain rule

- Expand & simplify
- Cancel common term to num & den
- Divide by -1

Final state

$$\sum x_n (f(x_n) - y_n) = 0$$

- note we cannot make the form = 0
- no closed solution
- cannot be solved analytically to find optimal  $w$

The Alternative solution is a ~~inter~~ iterative solution using gradient descent algorithm

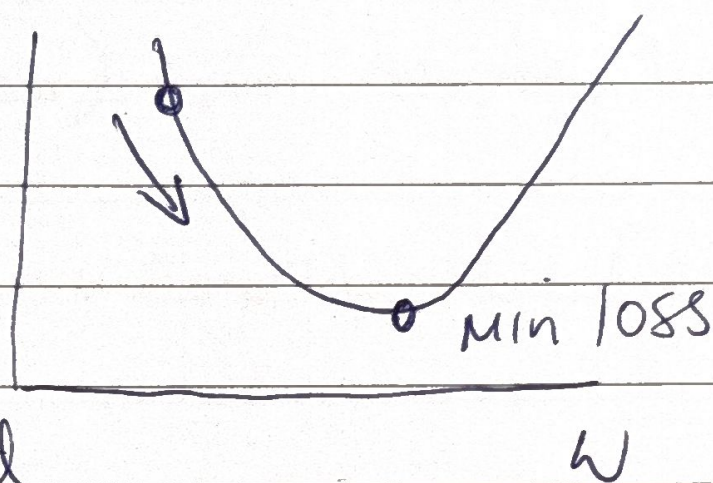
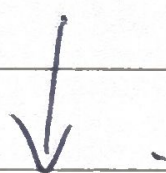


# Gradient Descent Algorithm

(4)

Optimal Parameters are at:

$$\frac{\partial L}{\partial w} = 0$$



Analytical Method  
to optimize solution

What we do instead is go in the dir  
of 0

Descend iteratively

iters

(0) init  $w$  get loss

(1) update  $w$   $w_{\text{new}} = w_{\text{old}} - \lambda \frac{dL(w)}{dw}$

$\lambda$  = learning rate

(2) iterate

(3) Stop when  $w_{\text{new}} \approx w_{\text{old}}$



(5)

log reg Address Overfit  
w/ regularization

$$L(w) = \sum -y_n \log f(x_n) - (1 - y_n) \log (1 - f(x))$$

cross entropy loss

$+ \frac{a \|w\|^2}{2}$  regularization term (2

Add reg parameter

## Summary

1. log reg is classifier
2. loss func = cross entropy loss
3. Optimized by using gradient Descent which has a learning rate hyperparam
4. Overfit Dealt w/ L2 reg