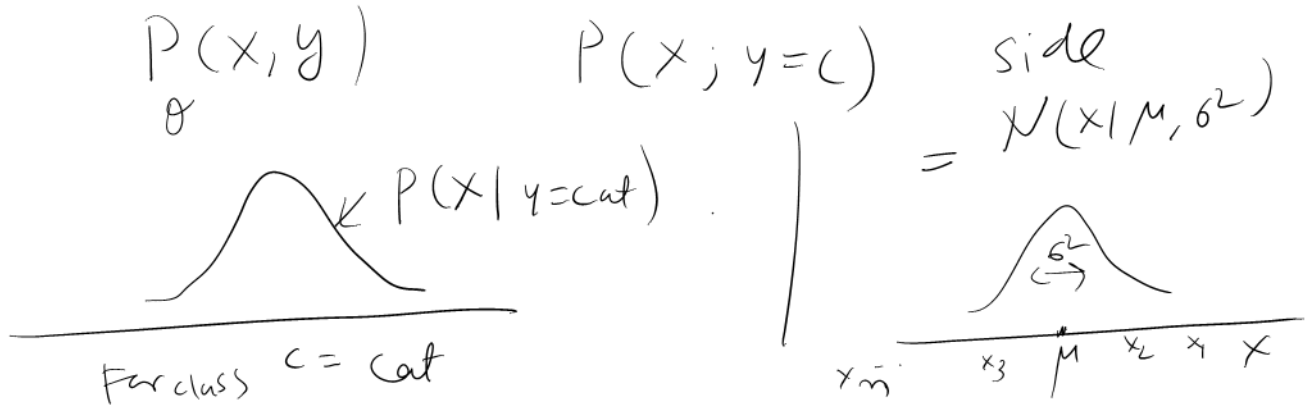


$$P(y=c|x;\theta) \propto \frac{P_\theta(X|y;\theta) P_\theta(y)}{N}$$

← generative classifier



$$P(Y=c | X; \theta)$$

In discriminative classifier we model this as a function $x; \theta$

2 We will write a direct function of X giving us the probability

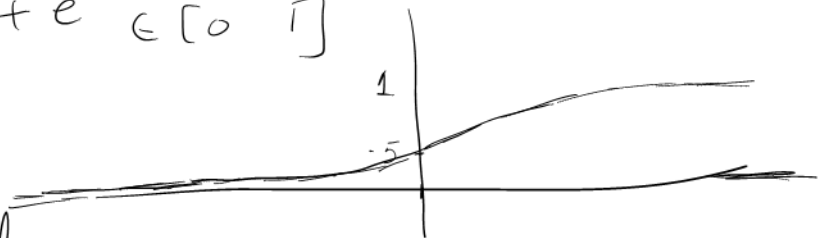
For Binary case logistic Regression

$$P(y=1|x;w) = \text{Ber}(y | \underbrace{g(w^T x)}_{= g(w^T x)}) \quad \left[\begin{array}{l} \text{classification} \\ \text{but via} \\ \text{regression} \end{array} \right]$$

side what is G (sigmoid function) $= \frac{1}{1 + e^{-w \cdot x}}$

~~(6.1)~~ $G(x) = \frac{1}{1 + e^{-x}}$ Here $x \in \mathbb{R}$
 \uparrow $\in [0, 1]$

Sigmoid takes a real number and map to $[0, 1)$ interval.



$$G'(x) = G(1 - G(x))$$

Popular ~~will~~ see that it is easy to extend
For multi-class

- using kernel trick can model non-linear $\phi(x)$

what is decision surface

$$4 \quad p(y=1 | x, w) = p(y=0 | x, w)$$

$$\frac{1}{1 + e^{-w^T x}} = 1 - \frac{1}{1 + e^{w^T x}}$$

$$\frac{e^{wTx}}{e^{wTx}} = \frac{e^{-wTx}}{1 + e^{wTx}} = \frac{1}{e^{wTx} + 1}$$

$$e^{wix} = 1$$

take \log_e

$$w^T x = 0 \quad \text{or} \quad w_1 x_1 + w_2 x_2 + w_d x_d = 0$$

hence decision boundary is
a line in high dim
or is hyperplane

Note $x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \\ 1 \end{bmatrix}$

Model Fitting (estimate w)

Conditional log likelihood

$$\log P(D) = \log \prod_{i=1}^N p(y_i | x_i; w)$$

$$= \log \prod_{i=1}^N \frac{1}{\Gamma(y_i)} \frac{\Gamma(y_i)}{\Gamma(y_i)} \frac{1}{(1 - \sigma(w^T x_i))^{1 - y_i}}$$

$$= \sum_{i=1}^N (y_i \log \sigma(w^T x_i) + (1 - y_i) \log (1 - \sigma(w^T x_i)))$$

$$NLL(D, w) = - \sum_{i=1}^N (y_i \log \sigma(w^T x_i) + (1 - y_i) \log (1 - \sigma(w^T x_i)))$$

$$p_{\text{ed}}(x_i) = \begin{bmatrix} \sigma(w^T x_i) \\ 1 - \sigma(w^T x_i) \end{bmatrix} \text{ also see } \text{oneHot}(y_i) = \begin{bmatrix} y_i \\ 1 - y_i \end{bmatrix}$$

$$= - \sum_{i=1}^N \text{cross} (p_{\text{ed}}(x_i), \text{oneHot}(y_i))$$

can't write MLE solution
in closed form?

solution

optimization algorithms (iterative)

$$\theta_{k+1} = \theta_k - \eta \left(\frac{\partial NLL(\theta)}{\partial \theta} \right)_k$$

effect of this is $NLL(\theta_{k+1}) \leq NLL(\theta_k)$

Learning rate η (step size)

derivative (direction of steep descent)

$$\text{Ber}(x, \theta) = \begin{cases} \theta & \text{if } x=1 \\ 1-\theta & \text{if } x=0 \end{cases}$$

$$\text{Ber}(x, \theta) = \theta^x (1-\theta)^{1-x}$$

$$a = \begin{bmatrix} a_1 \\ \vdots \\ a_K \end{bmatrix} \quad b = \begin{bmatrix} b_1 \\ \vdots \\ b_K \end{bmatrix}$$

$$\text{cross}(a, b) = - \sum_{i=1}^K a_i \log b_i$$

$$- \sum_{i=1}^K a_i \log a_i \quad \text{entropy}$$

$$\theta_{k+1} = \theta_k - \eta \left(\frac{\partial \mathcal{L}(\theta)}{\partial \theta} \right)_k + \mu (\theta_k - \theta_{k-1})$$

→ Momentum update

$$0 \leq \mu \leq 1$$

8.3.6

$$f(w) = \underbrace{\text{MLL}(w)}_{\text{+ve}} + \underbrace{\lambda \|w\|}_{\text{+ve}}$$

$$\begin{cases} D = \sum_{i=1}^N 2x_i y_i \\ T = \{x_i, y_i\}_{i=1}^{N/2} \quad V = \{x_i, y_i\}_{i=N/2+1}^N \end{cases}$$

can LASSO

hyperparameter (weight decay)



$$\text{or } f(w) = \text{MLL}(w) + \lambda \|w\|$$



