

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

← generative classifier

$$\propto P_\theta(X, y; \theta)$$

$$P_\theta(X, y)$$

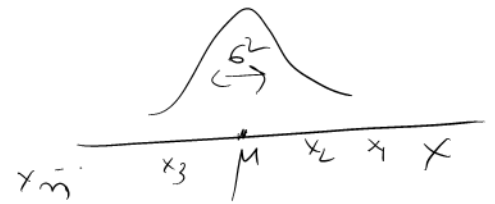
$$P(X; y=c)$$

$$\text{side } N(X|\mu, \sigma^2)$$



For class $c = \text{cat}$ feature of cat images
distribution of feature.

region of most probable feature for cat images



- Given probability distribution it is easy to evaluate the probability in an interval (continuous random variable) or at a point (discrete random variable) as
- But given probability distribution how do we generate sample x_1, x_2, \dots, x_n
(will cover sample generation later)

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

In discriminative classifier we model $P(Y=c | X; \theta)$ (this) as a function $X; \theta$.

We will write a direct function of X giving us the probability $P(Y=c | X; \theta)$
(Note: here we don't model joint distribution $P(X, Y)$)

For Binary case: logistic regression (LR)

Let $X \in \mathbb{R}^D, w \in \mathbb{R}^D$

$$P(Y=1 | X; w) = \text{Ber}(Y | G(w^T X)) \quad \left[\begin{array}{l} \text{classification} \\ \text{via} \\ \text{regression} \end{array} \right]$$

$\rightarrow G(w^T X)$

Side what is G (sigmoid function) $= \frac{1}{1 + e^{-w^T X}}$

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

Sigmoid takes a real number

and map it to $[0, 1]$ interval



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

LR is popular will see that it is easy to extend
(1) For multi-class

(2) using kernel trick can model non-linear decision boundaries

what is decision surface of LR?

$$P(Y=1 | X; w) = P(Y=0 | X; w)$$

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

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

take loge $e^{w^T X} = 1$
 $w^T X = 0$ or $w_1 x_1 + w_2 x_2 + \dots + w_D x_D = 0$

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

note $x = \begin{bmatrix} x_1 \\ \vdots \\ x_d \\ 1 \end{bmatrix}$
 modify add 1 as last feature

Model Fitting (estimate w)

let say we observe samples $D = \{ \{x_i, y_i\} \}_{i=1}^N$
 I.I.D samples

conditional log likelihood

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

$$= \log \prod_{i=1}^N \begin{matrix} \mathbb{I}(y_i=1) \\ \sigma(w^T x_i) \\ \mathbb{I}(y_i=0) \\ (1 - \sigma(w^T x_i)) \end{matrix}$$

side $\text{Ber}(x, \theta) = \begin{cases} \theta^x (1-\theta)^{1-x} & \text{for } x \in \{0, 1\} \end{cases}$

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

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

prediction via LR \rightarrow $\text{ped}(x_i) = \begin{bmatrix} \sigma(w^T x_i) \\ 1 - \sigma(w^T x_i) \end{bmatrix}$ also see $\text{oneHot}(y_i) = \begin{bmatrix} y_i \\ 1 - y_i \end{bmatrix}$

ground truth $\sum = 1$ vector sums up to 1

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

can't write MLE solution

in closed form?

But $NLL(D, w)$ is still a convex function in w

solution

we use optimization algorithms (iterative)

$$a = \begin{bmatrix} a_1 \\ \vdots \\ a_k \end{bmatrix} \quad b = \begin{bmatrix} b_1 \\ \vdots \\ b_k \end{bmatrix}$$

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

$$= - \sum_{i=1}^k a_i \log a_i \quad \text{entropy if } b = a$$

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

effect of this is

$NLL(\theta_{k+1}) \leq NLL(\theta_k)$ Called Learning rate (direction of steep descent!)

Hence we are guaranteed to not increase $NLL(\theta)$ in deep learning

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

→ Momentum update $0 \leq \mu \leq 1$

8.3.6

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

quadratic penalty term

side data split into

$$D = \{x_i, y_i\}_{i=1}^N$$

$$T = \{x_i, y_i\}_{i=1}^{N/2} \quad V = \{x_i, y_i\}_{i=N/2+1}^N$$

↑ training set ↑ validation set

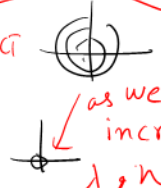
model selection. Will cover in detail

Can do LASSO too

hyper parameter (weight decay) regularization
called in deep learning

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

EFFECT of λ



as we increase λ, w
(decays towards 0)

or can mix penalty

$$f(w) = \text{MLL}(w) + \lambda_1 \|w\|^2 + \lambda_2 \|w\|$$

