# Logistic Regression for Classification:

Can we build a classification scheme that models and optimizes $p(y|x)$ without
the need to model relationships among different features/attributes?

$\longrightarrow$ Discriminative Modeling $p(y|x)$ ✓

Intuition from linear regression $\quad f(x) = \omega^T \phi(x) \quad$ [optimizing for $\omega$]

Can we learn a function $g \quad X \rightarrow y$

So that if $\quad \omega^T \phi(x) > 0 \longrightarrow g(x) = 1$
$\qquad \omega^T \phi(x) < 0 \longrightarrow g(x) = 0$

$\hookrightarrow$ Separating hyperplane in
the $X$ space, dividing
into two spaces/classes

We want to come up with a probabilistic model for $p(y|x)$

In the above, if $\quad p(y=1|x) = \begin{cases} 1 & , \omega^T \phi(x) > 0 \\ 0 & , \omega^T \phi(x) < 0 \end{cases}$

$p(y_1 \cdots y_n | x_1 \cdots x_n) = \prod p(y_i | x_i) = 0$

A single mistake will make the dataset to have a zero probability

We need to have confidence score about $p(y=1|x)$ farther from hyperplane
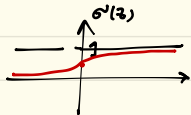then more confident in it being class 0 or 1 $\longrightarrow$ use exponential

miss confident

less confident

$$p(y=1 \mid x) \propto e^{\frac{1}{2} \omega^T \phi(x)} = \frac{1}{z} e^{\frac{1}{2} \omega^T \phi(x)}$$

$$p(y=0 \mid x) \propto e^{-\frac{1}{2} \omega^T \phi(x)} \quad \text{to have a symmetric classifier} = \frac{1}{z} e^{-\frac{1}{2} \omega^T \phi(x)}$$

$$\longrightarrow p(y=1 \mid x) + p(y=0 \mid x) = 1 \longrightarrow \frac{1}{z}\left[ e^{\frac{1}{2}\omega^T\phi(x)} + e^{-\frac{1}{2}\omega^T\phi(x)} \right] = 1$$

$$\longrightarrow z = e^{\frac{1}{2}\omega^T\phi(x)} + e^{-\frac{1}{2}\omega^T\phi(x)} \qquad \text{normalizing constant}$$

$$\longrightarrow p(y=1 \mid x) = \frac{e^{\frac{1}{2}\omega^T\phi(x)}}{e^{\frac{1}{2}\omega^T\phi(x)} + e^{-\frac{1}{2}\omega^T\phi(x)}} = \frac{1}{1+e^{-\omega^T\phi(x)}} = \sigma\left(\omega^T\phi(x)\right)$$

logistic function

$$\sigma(z) = \frac{1}{1+e^{-z}} \qquad \text{logistic function}$$



squashing $\omega^T\phi(x)$
(linear reg) to
$[0,1]$

* What is data ?

$$D = \left\{ (x^1,y^1), \quad , (x^n,y^n) \right\} \qquad \text{where } x^i = \begin{pmatrix} x^i_1 \\ x^i_2 \end{pmatrix}, \quad y^i \in \{0,1, ,L-1\}$$

categorical / continuous

categorical   $x^i = \{0,1,2\}$  ↗ sunny, cloudy, rainy

$$x^i = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \text{ or } \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \text{ or } \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} \qquad \text{"one hot" encoding} \qquad \text{all but 1 are zeros}$$

* Parameters to learn ?

weight vectors $w \in \mathbb{R}^{|\phi(x)|}$ which model $p(y=1 \mid x) = \dfrac{1}{1+e^{-\phi(x)^T \omega}}$

\* Cost function?    conditional log-likelihood

$$\ell(w) = \sum_{j=1}^{N} \log p(y^j | x^j, w)$$

$$p(y^j | x^j, w) = p(y^j=1 | x^j, w)^{1(y^j=1)} \times p(y^j=0 | x^j, w)^{1(y^j=0)}$$

$$\ell(w) = \sum_{j=1}^{N} 1(y^j=1) \log p(y^j=1 | x^j, w) + 1(y^j=0) \log p(y^j=0 | x^j, w)$$

$$= \sum_{j=1}^{N} 1(y^j=1) \log \frac{1}{1+e^{-\phi(x_j)^T w}} + \underbrace{1(y^j=0)}_{1-1(y^j=1)} \log \left(1 - \frac{1}{1+e^{-\phi(x_j)^T w}}\right)$$

$$= \sum_{j=1}^{N} \log \left(\frac{1}{1+e^{+w^T \phi(x_j)}}\right) + \sum_{j=1}^{N} 1(y^j=1) \log \frac{\frac{1}{1+e^{-\phi(x_j)^T w}}}{\frac{1}{1+e^{\phi(x_j)^T w}}}$$

$$= \sum_{j=1}^{N} \log \left(\frac{1}{1+e^{+w^T \phi(x_j)}}\right) + \sum_{j=1}^{N} 1(y^j=1) \log \frac{\frac{1}{1+e^{-\phi(x_j)^T w}}}{\frac{1}{1+e^{\phi(x_j)^T w}}}$$

$$= \sum_{j=1}^{N} \log \frac{1}{1+e^{w^T \phi(x_j)}} + \sum_{j=1}^{N} y^j \log \frac{1+e^{\phi(x_j)^T w}}{1+e^{-\phi(x_j)^T w}} \longrightarrow e^{\phi(x_j)^T w}\left(1+e^{-\phi(x_j)^T w}\right)$$

$$= -\sum_{j=1}^{N} \log \left(1+e^{w^T \phi(x_j)}\right) + \sum_{j=1}^{N} y^j \phi(x_j)^T w$$

$$\Rightarrow \ell(w) = \sum_{j=1}^{N} \left[ y^j \phi(x^j)^T w - \log \left(1+e^{w^T \phi(x_j)}\right) \right]$$

We want to maximize $\ell(w)$ or minimize $-\ell(w)$ !

$$J(w) = -\ell(w) = \sum_{j=1}^{N} \left[ -y^j \phi(x^j)^T w + \log \left(1+e^{w^T \phi(x_j)}\right) \right]$$

$$J(w) = -\ell(w) = \sum_{i=1}^{N} \left[ -y^i \phi(x^i)^T w + \log\left(1+e^{w^T \phi(x_i)}\right) \right]$$

$$\frac{\partial J}{\partial w} = \sum_{i=1}^{N} -y^i \phi(x^i) + \frac{\partial \log\left(1+e^{w^T \phi(x_i)}\right)}{\partial w} = -\frac{\frac{\partial}{\partial w}\sigma'(w^T\phi(x_i))}{\sigma'(w^T\phi(x_i))} \quad ))$$

$$-\phi(x_i)\,\sigma'(w^T\phi(x_i))(1-\sigma'(w^T\phi(x_i)))$$

$$\log\left(1+e^{w^T\phi(x_i)}\right) = -\log\left(\sigma'(-w^T\phi(x_i))\right)$$

$$\frac{\partial J}{\partial w} = \sum_i -\phi(x^i)\left[ y^i - P(y^i=1 \mid x^i, w) \right] \qquad \nearrow \frac{1}{1+e^{-\phi(x_i)^T w}}$$

No closed-form sol
for $\frac{\partial J}{\partial w} = 0$!

Gradient descent
+ Initialize $w^{(0)} = 0$
+ At iteration K, update $w^{(k)} = w^{(k-1)} - \alpha \frac{\partial J}{\partial w}$

## Overfitting

$$= w^{(k-1)} + \alpha \sum_i \phi(x^i)\left[y^i - P(y=1 \mid x^i w)\right]$$

what is the decision boundary for logistic reg? where $P(y=1 \mid x, w) = \frac{1}{2}$

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

$\longrightarrow \phi(x)^T w = 0$   ( scaling $w$ does not change this )



$\longrightarrow$ scale $w$ up by very large values to make $P(y=1 \mid x, w) \simeq 1$

overfitting



$\rightarrow$ test

$\hookrightarrow$ compute $\hat{w}$ here
and scale it $\alpha \hat{w}$
with large $\alpha > 0$

this is what
we want

To prevent overfitting, do regularization on $w$:

$$\min_{w} \; J(w) + \frac{\lambda}{2} \|w\|_2^2 \;\; \triangleq \bar{J}(w)$$

In G.D. $\quad \dfrac{\partial \bar{J}}{\partial w} = \underbrace{\dfrac{\partial J}{\partial w}}_{\text{already computed.}} + \lambda w$

$$= \sum_{i} - \phi(x^i) \left[ y^i - \frac{1}{1 + e^{-\phi(x_i)^T w}} \right] + \lambda w$$