# Logistic Regression (LR)

LR for binary classification.

$$I/P: \ \vec{x} = \vec{x} \in \mathbb{R}^d$$
$$O/P: \ y \in \{0, 1\}$$

We are interested in $P(Y = 1 | X = x)$
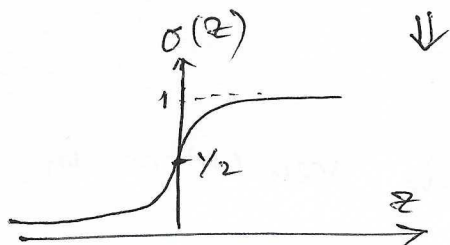
Model assumption: $Y | X = x \sim Bernoulli(\theta_x)$, $\theta_x \in [0, 1]$

Can we regress from $\vec{x} \to \theta_x$? (ie, modelling directly with a linear function)

Q ie. $\hat{\theta} = \vec{w}^T \vec{x}$

$\downarrow$ problem

$\theta$ has to be +ve & less than 1. $w^T x$ will produce some value, & we want something beth. 0 & 1.

$\Downarrow$



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

$$Y | X \sim Bernoulli(\sigma(w^T x)) \quad \leftarrow \underline{Logistic\ regression.}$$

So the posterior probability is,

$$P(Y = 1 | X = x) = \frac{1}{1 + e^{-w^T x}}$$

$\leftarrow$ measure $w$ directly from data.

## Some facts.

$$P(Y = 1 | X = x) = \frac{1}{1 + e^{-w^T x}} \quad \leftarrow complementary$$

$$P(Y = 0 | X = x) = 1 - \frac{1}{1 + e^{-w^T x}} = \frac{1}{1 + e^{w^T x}}$$

$$\hat{y}_{MAP} = \underset{y \in \{0, 1\}}{arg\,max} \ P(Y = y | X = x)$$

So we're literally choosing beth. two numbers:

$$\left.\begin{array}{c} P(Y = 1 | X) \\ P(Y = 0 | X) \end{array}\right\} \text{ so } \underset{y}{arg\,max} \ P(Y = y | x) = 1 \text{ is same as:}$$

$$P(Y = 1 | X) \overset{?}{\geqslant} P(Y = 0 | X)$$

$$\text{or} \quad \frac{P(Y = 1 | X)}{P(Y = 0 | X)} \overset{?}{\geqslant} 1 \quad \left.\begin{array}{c} \text{if yes then } \hat{y} = 1 \\ \text{else } \hat{y} = 0 \end{array}\right\}$$

Also, $\log\left[\dfrac{P(Y=1|X=x)}{P(Y=0|X=x)}\right] \overset{?}{\geq} 0$

$\Rightarrow \log\left[\dfrac{\left(\frac{1}{1+e^{-W^Tx}}\right)}{e^{-W^Tx}/1+e^{-W^Tx}}\right] \geq 0 \Rightarrow \log\left(e^{W^Tx}\right) \overset{?}{\geq} 0$

$\Rightarrow \boxed{W^Tx \overset{?}{\geq} 0}.$  $\leftarrow$ linear classifier.

So if $W^Tx \geq 0 \Rightarrow \hat{y}=1$

$\qquad W^Tx < 0 \Rightarrow \hat{y}=0$

Hence, $W^Tx \equiv$ score of class 1

$\dfrac{1}{1+e^{-W^Tx}} \equiv$ Probability of class 1.

## Estimation of $\vec{w}$

MLE $\infty$  Will be similar to the way we've seen before in coin toss / Bernoulli parameters!

| Coin toss | LR |
|---|---|
| $y \in \{0,1\}$ | $y \in \{0,1\}$ |
| Dataset: $D = \{y_1, \ldots, y_N\}$ | Dataset: $D = \{(x_1,y_1), \ldots (x_N, y_N)\}$ |
| $y \sim Ber(\theta)$ | $y|X=x \sim Ber(\underbrace{\sigma(w^Tx)}_{\theta_x})$ |
| $P(Y=1) = \theta$ | $P(Y=1|X=x) = \theta_x$ |
| $P(Y=0) = 1-\theta$ | $P(Y=0|X=x) = 1-\theta_x$ |
| Likelihood for one sample | Likelihood of one sample: |
| $L(\theta) = \theta^y (1-\theta)^{1-y}$ | $L(w) = \theta_x^y (1-\theta_x)^{1-y}$ |
| Likelihood of dataset: | Likelihood of dataset! |
| $L(\theta) = \prod_{i=1}^{N} \theta^{y_i}(1-\theta)^{1-y_i}$ | $L(w) = \prod_{i=1}^{N} \theta_{x_i}^{y_i}(1-\theta_{x_i})^{1-y_i}$ |
| $\downarrow$ | |
| $\theta^{\alpha_H}(1-\theta)^{\alpha_T}$ | $= \prod_{i=1}^{N}\left[\sigma(w^Tx_i)\right]^{y_i}\left[1-\sigma(w^Tx_i)\right]^{1-y_i}$ |

So the log likelihood for LR is,

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

$$= \sum_{i=1}^{N} \left[ y_i \log\left(\frac{1}{1+e^{-w^T x_i}}\right) + (1-y_i)\log\left(\frac{1}{1+e^{w^T x_i}}\right) \right]$$

$$= \sum_{i=1}^{N} \left[ y_i \log\left(\frac{e^{w^T x_i}}{1+e^{w^T x_i}}\right) + (1-y_i)\log\left(\frac{1}{1+e^{w^T x_i}}\right) \right]$$

$$= \sum_{i=1}^{N} \left[ y_i \log\left(e^{w^T x_i}\right) - y_i \log\left(1+e^{w^T x_i}\right) + \log(1) - \log\left(1+e^{w^T x_i}\right) \right.$$
$$\left. - y_i \log(1) + y_i \log\left(1+e^{w^T x_i}\right) \right]$$

$$= \sum_{i=1}^{N} \left[ y_i \log\left(e^{w^T x_i}\right) - \log\left(1+e^{w^T x_i}\right) \right]$$

$$= \sum_{i=1}^{N} [\underbrace{y_i w^T x_i}_{\substack{\text{Linear} \\ \text{in } w}} - \underbrace{\log(1+e^{w^T x_i})}_{\substack{\text{not linear} \\ \text{in } w}}] \rightarrow \substack{\text{no closed} \\ \text{form solution}}$$

## Vanilla Gradient Descent

Initialize $w^{(0)}$ & then use the update rule:          ← Actually this is gradient ascent

$$w^{t+1} = w^t + \underbrace{\eta}_{\substack{\text{step size/} \\ \text{learning rate}}} \underbrace{\frac{\partial LL(w)}{\partial w}}_{\substack{\text{Gradient (direction of} \\ \text{steepest increase)}}}$$

$$w^{t+1} = w^t + \eta \sum_{i=1}^{N} [y_i x_i^T - \underbrace{\frac{1}{1+e^{w^T x_i}} e^{w^T x_i}}_{P(y_i=1\,|\,x_i, w)} \cdot x_i^T]$$

$$= w^t + \eta \sum_{i=1}^{n} [\underbrace{y_i}_{\substack{\text{truth} \\ \{0,1\}}} - \underbrace{P(y_i=1\,|\,x_i, \vec{w})}_{\substack{\text{what our} \\ \text{model} \\ \text{believes}}}] \underbrace{x_i^T}_{\substack{\uparrow \\ \text{feature vector/} \\ \text{data}}}$$

$$\underbrace{\qquad\qquad\qquad\qquad}_{\text{error}}$$

Geometric interperation of the above equation

Batch gradient descent $\Rightarrow$ moves the entire dataset towards the direction of the gradient

Stochastic " " $\Rightarrow$ perform on sampled points.



GD may or may not work.

Local vs Global minima



GD works always.

---

**MAP estimation of $\vec{w}$**

Model: $Y|X \sim Ber(\sigma(w^T x))$
$w_j \sim N(0, t^2)$ (i.i.d.)

$\hat{W}_{MAP} = \underset{w}{\arg\max} \log P(W|D) = \underset{w}{\arg\max} (\log P(D|W) + \log P(W))$

Now as usual, we'll take derivative w.r.t. $W$ & set to $0$. For the first term, we have already done it. for the second term!

$$P(\vec{W}) = \prod_{j=1}^{d} \frac{1}{\sqrt{2\pi t^2}} e^{-\frac{w^2}{2t^2}} \quad (i.i.d.)$$

$$\log P(\vec{W}) = -\frac{w^2}{2t^2} - \frac{d}{2} \log(2\pi t^2)$$

$$\frac{\partial}{\partial w}[\log P(\vec{W})] = -\frac{2\vec{W}}{2t^2} = -\lambda \vec{W}$$

So the total expression becomes,

$$\frac{\partial}{\partial w}[\cdots] = \sum_{i=1}^{N} [y_i - P(Y=1|\vec{x}_i)]\vec{x}_i^T - \lambda \vec{W}$$

$\underbrace{\qquad\qquad}_{\text{focus on mistakes}}$   $\underbrace{\qquad}_{\text{But try to keep norm of } \vec{w} \text{ small}}$

explain the labels

don't get too confident unless you must.

data wants to increase $||W||$ if it helps in classification

prior resists large weights

$\underbrace{\qquad\qquad\qquad}_{\text{MAP balances the two forces.}}$