

logistic regression

used for classification

training data : $\{(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n)\}$

▢ feature vector $\vec{x}_i \in \mathbb{R}^d \leftarrow$ how many features

▢ label $y_i \in \{0, 1\}$

in logistic regression, we learn from the data a probabilistic model for

$$\Pr(Y=y | \vec{x}) \quad \text{for } Y=0 \text{ and } Y=1.$$

label \leftarrow feature

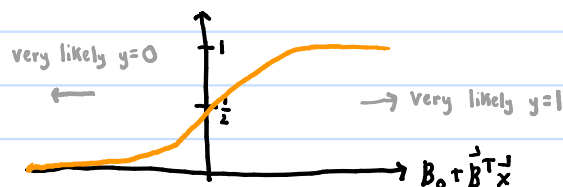
given a data pt. w/ feature vector \vec{x} , what is the probability that its label is $y=1$? $\Pr(Y=1 | \vec{x})$

the logistic model: a linear model for the log odds

$$\log\left(\frac{\Pr(Y=1 | \vec{x})}{1 - \Pr(Y=1 | \vec{x})}\right) = \beta_0 + \vec{\beta}^T \vec{x}$$

odds that $Y=1$ given \vec{x} intercept $\beta_0 \in \mathbb{R}$ coefficients $\vec{\beta} \in \mathbb{R}^d$

$$\Pr(Y=1 | \vec{x}) = \frac{e^{\beta_0 + \vec{\beta}^T \vec{x}}}{1 + e^{\beta_0 + \vec{\beta}^T \vec{x}}} \in (0, 1) \quad \text{"}$$



$$\text{and } \Pr(Y=0 | \vec{x}) = 1 - \Pr(Y=1 | \vec{x}).$$

how to learn $\beta_0, \vec{\beta}$ from data?

under this probabilistic model, write the likelihood, the probability of seeing the data given the probabilistic model & its parameters, then choose the parameters that maximize the likelihood.

$$\mathcal{L}(\beta_0, \vec{\beta}) = \prod_{i=1}^n \Pr(Y=y_i | \vec{x}_i) = \prod_{i=1}^n \Pr(Y=1 | \vec{x}_i)^{y_i} [1 - \Pr(Y=1 | \vec{x}_i)]^{1-y_i}$$

↑
likelihood

assumes independent,
identically distributed data

$$\begin{aligned} \log \mathcal{L} &= \sum_{i=1}^n y_i \log[\Pr(Y=1 | \vec{x}_i)] + (1-y_i) \log[1 - \Pr(Y=1 | \vec{x}_i)] \\ &= \sum_{i=1}^n y_i \log\left[\frac{e^{\beta_0 + \vec{\beta}^T \vec{x}_i}}{1 + e^{\beta_0 + \vec{\beta}^T \vec{x}_i}}\right] + (1-y_i) \log\left[\frac{1}{1 + e^{\beta_0 + \vec{\beta}^T \vec{x}_i}}\right] \\ &= \sum_{i=1}^n y_i (\beta_0 + \vec{\beta}^T \vec{x}_i) - \log(1 + e^{\beta_0 + \vec{\beta}^T \vec{x}_i}) \end{aligned}$$

maximizing $\log \mathcal{L}$ gives same $\beta_0, \vec{\beta}$ as maximizing \mathcal{L} .
 $\log \mathcal{L}$ is monotonic...

$$\rightarrow \vec{\nabla}_{\vec{\beta}} \log \mathcal{L} = \vec{0} \quad \text{to find } \vec{\beta} \text{ that fits the data ("learn", "train")}$$

$$\sum_{i=1}^n \vec{x}_i (y_i - \Pr(Y=1 | \vec{x}_i)) = \vec{0}$$

... non-linear in $\vec{\beta}$

one option to solve for $\vec{\beta}$: the Newton-Raphson algorithm.