

# Logistic Regression

Observation :  $Y \in \{0, 1\}$

e.g. presence ~~of~~ or absence of disease

Feature :  $\underline{X} = \begin{bmatrix} x_1 \\ \vdots \\ x_p \end{bmatrix} \in \mathbb{R}^p$

Regression coefficient:  $\underline{\beta} = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_p \end{bmatrix} \in \mathbb{R}^p$

For data with binary observation, the likelihood function is

$$P(Y = y_i \mid \underline{X}_i = \underline{x}_i) = p(\underline{x}_i, \underline{\beta})^{y_i} (1 - p(\underline{x}_i, \underline{\beta}))^{1 - y_i}$$

Why this likelihood?

$$P(Y = 1 \mid \underline{X}_i = \underline{x}_i) = p(\underline{x}_i, \underline{\beta})$$

$$P(Y = 0 \mid \underline{X}_i = \underline{x}_i) = 1 - p(\underline{x}_i, \underline{\beta})$$

Some function linking  $\underline{x}_i, \underline{\beta}$  and the "probability"  
↑  
need to define this.

Therefore, by independence we have

$$\prod_{i=1}^N P(Y = y_i \mid \underline{X}_i = \underline{x}_i) = \prod_{i=1}^N p(\underline{x}_i, \underline{\beta})^{y_i} (1 - p(\underline{x}_i, \underline{\beta}))^{1 - y_i}$$

Now, we need a model for  $p(\underline{x}_i, \underline{\beta})$ .

Option 1:  $p(\underline{x}_i, \underline{\beta}) = \beta_0 + \underline{\beta}^T \underline{x}_i$ .

This does not work because

$p(\underline{x}_i, \underline{\beta})$  must be probability so that

$0 < p(\underline{x}_i, \underline{\beta}) < 1$ . But if  $p(\underline{x}_i, \underline{\beta}) = \beta_0 + \underline{\beta}^T \underline{x}_i$ , then

$$p(\underline{x}_i, \underline{\beta}) \rightarrow \pm \infty.$$

option 2  $\log P(\underline{x}_i, \beta) = \beta_0 + \beta^T \underline{x}_i$

This is better but still doesn't work

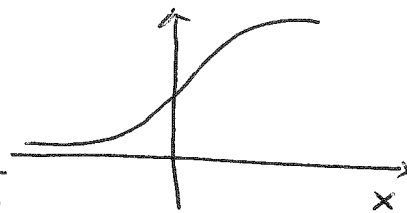
$P(\underline{x}_i, \beta) > 0$  because  $\beta_0 + \beta^T \underline{x}_i \rightarrow -\infty \Rightarrow \log P \rightarrow 0$ .

But unbounded above.

option 3  $\log \frac{P(\underline{x}_i, \beta)}{1 - P(\underline{x}_i, \beta)} = \beta_0 + \beta^T \underline{x}_i$

Then

$$P(\underline{x}_i, \beta) = \frac{1}{1 + e^{-(\beta_0 + \beta^T \underline{x}_i)}}$$



Three properties:

(i)  $0 < P(\underline{x}, \beta) < 1$

(ii) decision boundary is at  $\beta_0 + \beta^T \underline{x} = 0$

(iii) ~~rate of~~ sharpness of decision depends on  $\|\beta\|$ .

How to estimate  $\beta$  from logistic model?

Recall  $P(Y=1 | X_i = \underline{x}_i) = P(\underline{x}_i, \beta) = \frac{1}{1 + e^{-(\beta_0 + \beta^T \underline{x}_i)}}$   
 $= \frac{e^{\beta_0 + \beta^T \underline{x}_i}}{1 + e^{\beta_0 + \beta^T \underline{x}_i}}$

So the negative  $-\log$  likelihood is

$$\begin{aligned} & -\log P(\underline{Y} = \underline{y} | \underline{X} = \underline{x}) \\ &= -\log \left\{ \prod_{i=1}^N P(\underline{x}_i, \beta)^{y_i} (1 - P(\underline{x}_i, \beta))^{1-y_i} \right\} \\ &= -\sum_{i=1}^N \left\{ y_i \log P(\underline{x}_i, \beta) + (1-y_i) \log(1 - P(\underline{x}_i, \beta)) \right\} \\ &= -\sum_{i=1}^N \left\{ y_i \log \frac{P(\underline{x}_i, \beta)}{1 - P(\underline{x}_i, \beta)} + \log(1 - P(\underline{x}_i, \beta)) \right\} \end{aligned}$$

$$= - \sum_{i=1}^N \left\{ y_i (\beta_0 + \beta^T \underline{x}_i) + \log \left( \frac{1}{1 + e^{\beta_0 + \beta^T \underline{x}_i}} \right) \right\}$$

$$= - \sum_{i=1}^N \left\{ y_i (\beta_0 + \beta^T \underline{x}_i) - \log (1 + e^{\beta_0 + \beta^T \underline{x}_i}) \right\}$$

So we can take derivative:

$$\frac{\partial}{\partial \beta} (\cdot) = 0 \implies \text{transcendental equation} \\ (\text{can be solved numerically}).$$

To include sparsity in the logistic regression, we add

$$\min_{(\beta_0, \beta)} \left\{ - \sum_{i=1}^N \left\{ y_i (\beta_0 + \beta^T \underline{x}_i) - \log (1 + e^{\beta_0 + \beta^T \underline{x}_i}) \right\} + \lambda \|\beta\|_1 \right\}.$$

## Poisson Regression

Observation:  $Y = 0, 1, 2, \dots$

The likelihood function is

$$P(Y = y_i \mid X_i = \underline{x}_i) = \frac{\lambda_i^{y_i} e^{-\lambda_i}}{y_i!},$$

where  $\lambda_i$  is a function of  $\underline{x}_i, \beta$ .

(So technically should be ~~not~~

$$\lambda_i = \lambda(\underline{x}_i, \beta).)$$

Therefore, where there are  $N$  observations.

$$\prod_{i=1}^N P(Y = y_i \mid X_i = \underline{x}_i) = \prod_{i=1}^N \frac{\lambda(\underline{x}_i, \beta)^{y_i} e^{-\lambda(\underline{x}_i, \beta)}}{y_i!}$$

What should be the model of  $\lambda(\underline{x}_i, \beta)$ ?

Again, we choose

$$\log\{\lambda(\underline{x}_i, \beta)\} = \beta_0 + \beta^T \underline{x}_i. \quad (1)$$

Note: we do not need to upper bound  $\lambda(\underline{x}_i, \beta)$ , and so we do not need things like  $\frac{\lambda(\underline{x}_i, \beta)}{1 - \lambda(\underline{x}_i, \beta)}$ .

Remark:

$$\lambda(\underline{x}_i, \beta) = \mathbb{E}[Y | X = \underline{x}_i].$$

So (1) can be written as

$$\log \mathbb{E}[Y | X = \underline{x}_i] = \beta_0 + \beta^T \underline{x}_i.$$

How to solve  $\beta$  from Poisson Regression?

The negative log-likelihood is

$$\begin{aligned} & -\log \mathbb{P}(Y = \underline{y} | X = \underline{x}) \\ &= -\log \left\{ \prod_{i=1}^N \frac{\lambda(\underline{x}_i, \beta)^{y_i} e^{-\lambda(\underline{x}_i, \beta)}}{y_i!} \right\} \\ &= -\sum_{i=1}^N \left\{ y_i \log \lambda(\underline{x}_i, \beta) + \lambda(\underline{x}_i, \beta) - \log(y_i!) \right\} \end{aligned}$$

by (1)

$$= -\sum_{i=1}^N \left\{ y_i (\beta_0 + \beta^T \underline{x}_i) + e^{\beta_0 + \beta^T \underline{x}_i} - \log(y_i!) \right\}$$

So  $\beta$  can be found by

$$(\hat{\beta}_0, \hat{\beta}) = \underset{(\beta_0, \beta)}{\operatorname{argmin}} -\sum_{i=1}^N \left\{ y_i (\beta_0 + \beta^T \underline{x}_i) + e^{\beta_0 + \beta^T \underline{x}_i} \right\}.$$

if we want sparsity, then add  $\lambda \|\beta\|_1$ .

## Generalized Linear Model (GLM)

From the above examples, we see that for data that cannot be model by the standard linear model, the GLM could become very useful. In general, the transformation of the conditional mean

$$g(\mathbb{E}[Y | \underline{X} = \underline{x}]) \stackrel{\text{def}}{=} \mu(\underline{x})$$

is called a link function.

$$\begin{aligned} \text{Logit} : \mathbb{E}[Y | \underline{X} = \underline{x}] &= 1 \cdot \mathbb{P}(Y=1 | \underline{X} = \underline{x}) + 0 \cdot \mathbb{P}(Y=0 | \underline{X} = \underline{x}) \\ &= \mathbb{P}(Y=1 | \underline{X} = \underline{x}) = \mathbb{P}(\underline{X} > \beta) \end{aligned}$$

$$\text{Poisson} : \mathbb{E}[Y | \underline{X} = \underline{x}] = \lambda(\underline{x}, \beta).$$

More generally, we have GLM being

$$g(\mu(\underline{x})) = \beta_0 + \beta^T \underline{x}$$

The choice of  $g(\cdot)$  depends on the model:

$$\text{Logit} : g(\mu) = \log\left(\frac{\mu}{1-\mu}\right)$$

$$\text{Poisson} : g(\mu) = \log \mu$$

Why is it called the generalized linear model?

The linear model is a special case:

$$\mathbb{E}[Y | \underline{X} = \underline{x}] = \mu(\underline{x}) = \beta_0 + \beta^T \underline{x}$$

$$g(\mu) = \mu.$$

And if we consider a Gaussian likelihood:

$$\prod_{i=1}^N \mathbb{P}(Y = y_i | \underline{X}_i = \underline{x}_i) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - \mu(\underline{x}_i, \beta))^2}{2\sigma^2}}$$

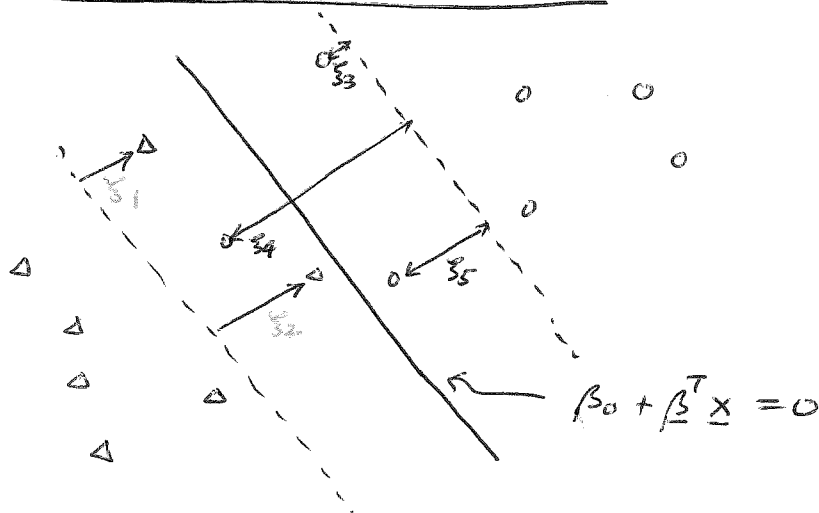
then its negative log-likelihood is

$$-\log(\cdot) = -\left\{ \sum_{i=1}^N \left( \log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) - \frac{(y_i - \mu(x_i, \beta))^2}{2\sigma^2} \right) \right\}$$

So the regression coefficients are

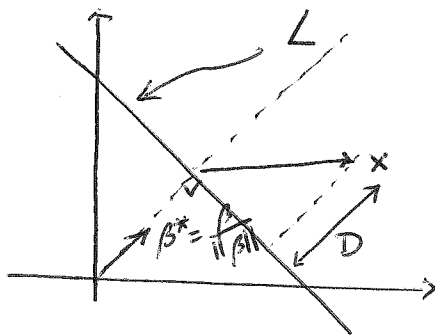
$$\begin{aligned} (\hat{\beta}_0, \hat{\beta}) &= \underset{(\beta_0, \beta)}{\operatorname{argmin}} \sum_{i=1}^N (y_i - \mu(x_i, \beta))^2 \\ &= \underset{(\beta_0, \beta)}{\operatorname{argmin}} \sum_{i=1}^N (y_i - (\beta_0 + \beta^T x_i))^2 \end{aligned}$$

## Support Vector Machine



Idea: Find a separating hyperplane that has the ~~minimum~~ <sup>maximum</sup> margin, and all the data points should be as far from the margin as possible.

But first we need a way to measure the distance between a point  $\underline{x}$  and the line  $L: \beta_0 + \beta^T \underline{x} = 0$ .



(i) if  $\underline{x}_1, \underline{x}_2$  are on  $L$ , the

$$\beta_0 + \beta^T \underline{x}_1 = 0 = \beta_0 + \beta^T \underline{x}_2$$

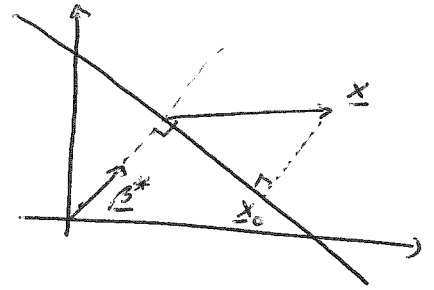
$$\Rightarrow \beta^T (\underline{x}_1 - \underline{x}_2) = 0$$

So  $\underline{\beta}^* \stackrel{\text{def}}{=} \frac{\underline{\beta}}{\|\underline{\beta}\|_2}$  is the normal  $\beta$ .

(2) For any point  $x_0 \in L$ ,

$$\beta_0 + \beta^T x_0 = 0$$

$$\Rightarrow \beta^T x_0 = -\beta_0$$



(3) ~~The~~ The distance between  $x$  and  $L$  is

$$D = \beta^{*T} (x - x_0), \text{ where } x_0 \text{ is a point on } L$$

$$= \frac{\beta^T (x - x_0)}{\|\beta\|}$$

$$= \frac{1}{\|\beta\|} (\beta^T x + \beta_0) \leftarrow \text{this is a signed distance}$$

Going back to SVM formulation, we are interested in solving the problem

$$\begin{aligned} & \max_{\beta_0, \beta} M \\ & \text{s.t. } \frac{1}{\|\beta\|} y_i (\beta_0 + \beta^T x_i) \geq M \end{aligned}$$

Interpretation: (i)  $M = \text{margin}$

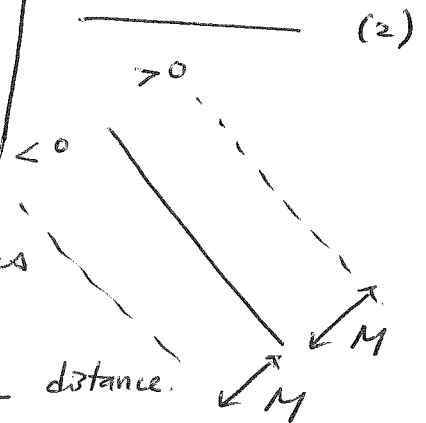
We want  $M$  to be as large as possible

(ii)  $\frac{\beta_0 + \beta^T x_i}{\|\beta\|} \cdot y_i$  is the unsigned distance.

if  $\beta_0 + \beta^T x_i > 0$ , then ideally  $y_i > 0$

if  $\beta_0 + \beta^T x_i < 0$ , then ideally  $y_i < 0$ .

So if  $y_i (\beta_0 + \beta^T x_i) < 0$ , then there is misclassification.



## Modification 1

$$y_i(\beta_0 + \beta^T x_i) \geq M \|\beta\| \quad (3)$$

can be simplified as (4)

$y_i(\beta_0 + \beta^T x_i) \geq 1$ , because if  $(\beta_0, \beta)$  satisfies (3),  $(\frac{1}{M}\beta_0, \frac{1}{M}\beta)$  will satisfy (4). So we can solve this problem:

$$\begin{aligned} \max_{(\beta, \beta_0)} & \frac{1}{2} \|\beta\|^2 \\ \text{s.t.} & y_i(\beta_0 + \beta^T x_i) \geq 1, \quad i=1, 2, \dots, N \end{aligned}$$

## Modification 2

We want to tolerate some misclassification. So we need to relax (4) by

$$y_i(\beta_0 + \beta^T x_i) \geq 1 - \xi_i, \quad i=1, 2, \dots, N$$

and  $\xi_i \geq 0$ ,  $\sum_{i=1}^N \xi_i \leq \text{constant}$ .

So the SVM problem becomes

$$\begin{aligned} \max_{\{\xi_i\}, (\beta, \beta_0)} & \frac{1}{2} \|\beta\|^2 \\ \text{s.t.} & y_i(\beta_0 + \beta^T x_i) \geq 1 - \xi_i \\ & \xi_i \geq 0, \quad \sum_{i=1}^N \xi_i \leq \text{constant} \end{aligned}$$

This can be re-written as

$$\begin{aligned} \max_{(\beta, \beta_0), \{\xi_i\}} & \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^N \xi_i \quad (P_1) \\ \text{s.t.} & y_i(\beta_0 + \beta^T x_i) \geq 1 - \xi_i \\ & \xi_i \geq 0 \end{aligned}$$



Let  $c = \frac{1}{\lambda}$ , then the objective of  $(P_1)$  becomes

$$\min_{(\beta, \beta_0)} \frac{\lambda}{2} \|\beta\|^2 + \sum_{i=1}^N \xi_i$$

Note that by the constraints of  $(P_1)$ , we have

$$\xi_i \geq 1 - y_i (\beta_0 + \beta^T x_i)$$

and  $\xi_i \geq 0$ .

$$\text{So } \sum_{i=1}^N \xi_i \geq \sum_{i=1}^N \left[ 1 - y_i (\beta_0 + \beta^T x_i) \right]_+$$

where  $[\cdot]_+$  ~~denotes~~ returns the positive part of the argument

Therefore,  $(P_1)$  can be solved by ~~finding~~ <sup>minimizing</sup> its lower bound:

$$\min_{(\beta, \beta_0)} \frac{\lambda}{2} \|\beta\|^2 + \sum_{i=1}^N \left[ 1 - y_i (\beta_0 + \beta^T x_i) \right]_+$$

And by absorbing  $\frac{1}{N}$  into the equation, we have

$$\min_{(\beta, \beta_0)} \frac{1}{N} \sum_{i=1}^N \left[ 1 - y_i (\beta_0 + \beta^T x_i) \right]_+ + \frac{\lambda}{2} \|\beta\|^2$$

If we want to enforce sparsity, then we can do

$$\min_{(\beta, \beta_0)} \frac{1}{N} \sum_{i=1}^N \left[ 1 - y_i (\beta_0 + \beta^T x_i) \right]_+ + \lambda \|\beta\|_1 \quad (P_2)$$