

Lecture 9 - Feb 7

Logistic Regression Assessing Classifiers

Reading

Week 5 Notes in GitHub

[*Data Mining and Machine Learning*](#)

24 Logistic Regression

[Elements of Statistical Learning](#)

4.4 Logistic Regression

Upcoming Deadlines

Homework 2 (Feb 18)

Logistic Regression

Given independent predictors X_1, \dots, X_d + categorical output Y ,
logistic regression tries to predict a probability distribution for Y
 $P(c_i | x_j)$

Binary Logistic Regression

$$X = \begin{pmatrix} 1 & x_{11} & \cdots & x_{1d} \\ 1 & x_{21} & \cdots & x_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{nd} \end{pmatrix}^{\textcolor{blue}{X_1}}_{\textcolor{blue}{X_n}}$$

with $Y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$ where $y_i \in \{0, 1\}$

$\textcolor{blue}{x_1, \dots, x_d}$ independent predictors

$\textcolor{blue}{y_1, \dots, y_n}$ Bernoulli RV predictors

Let X_1, \dots, X_d be independent predictors

$Y \in \{0, 1\}$ Bernoulli RV.

D - dataset of $x_i \in \mathbb{R}^{d+1}$, $x_i = (1, x_{i1}, \dots, x_{id})$, y_i labels

$$\pi(x) = P(Y=1 | X=x) \Rightarrow 1 - \pi(x) = P(Y=0 | X=x)$$

\uparrow
unknown true probability

Goal: predict $\pi(x)$

A linear regression model $\pi(x) = \Theta^T x$ is not appropriate since $\pi(x)$ is a probability in $[0, 1]$.

\Rightarrow we feed $\Theta^T x$ into a function squishing it into $[0, 1]$

The logistic (sigmoid) function $\sigma: \mathbb{R} \rightarrow (0, 1)$ is

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

Properties :

$$\textcircled{1} \quad 0 < G(z) < 1 \text{ since } 0 < \frac{1}{1+e^z} < 1$$

$$\textcircled{2} \quad G'(z) = \frac{e^z(1+e^z) - e^z(e^z)}{(1+e^z)^2} = \frac{e^z}{1+e^z} \cdot \frac{1}{1+e^z}$$

$$= G(z) \cdot G(-z) = G(z)(1-G(z)) > 0 \quad (G \text{ is increasing})$$

$$\textcircled{3} \quad \lim_{z \rightarrow -\infty} G(z) = \lim_{z \rightarrow -\infty} \frac{1}{1+e^z} = 0$$

$$\textcircled{4} \quad \lim_{z \rightarrow \infty} G(z) = \lim_{z \rightarrow \infty} \frac{1}{1+e^z} = \frac{1}{1+0} = 1$$

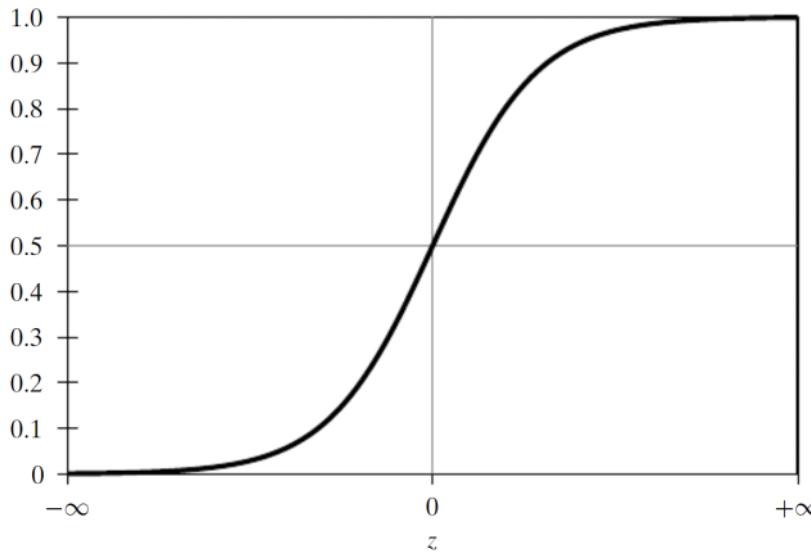


Figure 24.1. Logistic function.

Logistic regression model:

$\hat{\alpha}^T x$

LOGISTIC REGRESSION

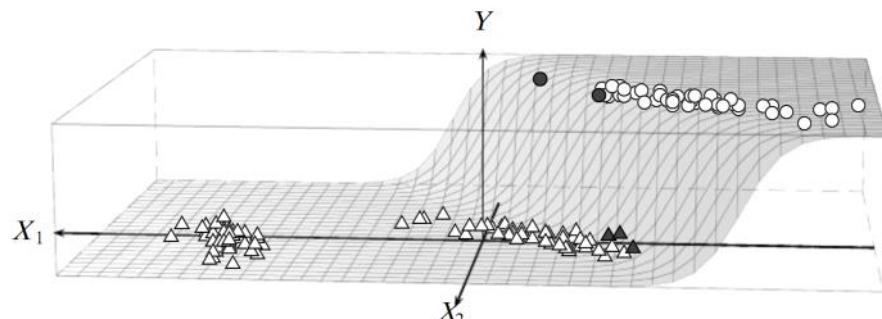
$$\left. \begin{array}{l} P(Y=1|X=x) = \sigma(\theta^T x) = \frac{e^{\theta^T x}}{1+e^{\theta^T x}} \\ P(Y=0|X=x) = \sigma(-\theta^T x) = \frac{1}{1+e^{\theta^T x}} \end{array} \right\} P(Y|X=x) = \sigma(\theta^T x)^y \sigma(\theta^T x)^{1-y}$$

The log-odds ratio for $Y=1$ is

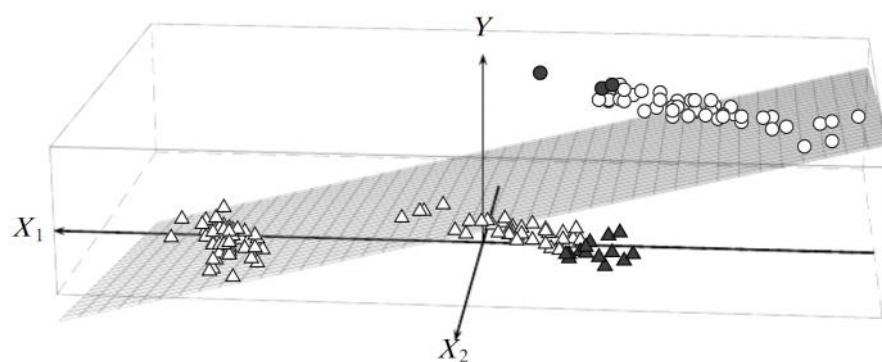
$$\log\left(\frac{\hat{\pi}(x)}{1-\hat{\pi}(x)}\right) = \log\left(\frac{e^{\theta^T x}}{1+e^{\theta^T x}}\right) = \theta^T x$$

\Rightarrow logistic regression assumes log-odds are linear in θ .

prediction rule: $\hat{y} = \begin{cases} 1, & \text{if } \sigma(\theta^T x) \geq 0.5 \\ 0, & \text{if } \sigma(\theta^T x) < 0.5 \end{cases}$



(a) Logistic Regression



(b) Linear Regression

Figure 24.2. Logistic versus linear regression: Iris principal components data. Misclassified point are shown in dark gray color. Circles denote Iris-virginica and triangles denote the other two Iris types.

Likelihood is the probability of the observed data given the estimated parameters

$$L(\theta) = P(Y|\theta) = \prod_{i=1}^n P(y_i|x_i) = \prod_{i=1}^n \sigma(\theta^T x_i)^{y_i} \sigma(-\theta^T x_i)^{1-y_i}$$

We will solve $\max_{\theta} L(\theta)$ or, equivalently $\max_{\theta} \underbrace{\log(L(\theta))}_{\text{log-likelihood}}$

$$\begin{aligned} \log(L(\theta)) &= \log \left(\prod_{i=1}^n \sigma(\theta^T x_i)^{y_i} \sigma(-\theta^T x_i)^{1-y_i} \right) \\ &= \sum_{i=1}^n y_i \log(\sigma(\theta^T x_i)) + (1-y_i) \log(\sigma(-\theta^T x_i)) \end{aligned}$$

Equivalently, minimize the cross-entropy error function

$$E(\theta) = -\log(L(\theta)) = -\sum_{i=1}^n \left(y_i \log \underbrace{\sigma(\theta^T x_i)}_{\pi(x_i)} + (1-y_i) \log \underbrace{\sigma(-\theta^T x_i)}_{1-\pi(x_i)} \right)$$

Typically, we solve $\nabla E(\theta) = 0$ to try to solve $\min_{\theta} E(\theta)$, but there is no closed-form solution for θ , so we use gradient descent. An explicit gradient ∇E speeds it up, so let's find it.

$$\frac{\partial}{\partial z_i} \left\{ \log(\sigma(z_i)) \right\} = \frac{1}{\sigma(z_i)}$$

$$\frac{\partial}{\partial z_i} \left\{ \log(\sigma(-z_i)) \right\} = \frac{\partial}{\partial z_i} \left\{ \log(1 - \sigma(z_i)) \right\} = \frac{-1}{1 - \sigma(z_i)}$$

$$\frac{\partial \sigma(z_i)}{\partial z_i} = \sigma(z_i)(1 - \sigma(z_i)) = \sigma(z_i)\sigma(-z_i)$$

$$\nabla_{\theta}(z_i) = \nabla_{\theta}(\theta^T x_i) = x_i$$

$$\Rightarrow \nabla_{\theta}(\log(\sigma(\theta^T x_i))) = \frac{\partial \log(\sigma(\theta^T x_i))}{\partial \sigma(\theta^T x_i)} \cdot \frac{\partial \sigma(\theta^T x_i)}{\partial \theta^T x_i} \cdot \frac{\partial \theta^T x_i}{\partial \theta}$$

↙

$$= \frac{1}{\sigma(\theta^T x_i)} \cancel{\sigma(\theta^T x_i)\sigma(-\theta^T x_i)} \cdot x_i$$

$$= \sigma(-\theta^T x_i) x_i$$

$$\nabla_{\theta}(\log(\sigma(-\theta^T x_i))) = \frac{\partial \log(\sigma(-\theta^T x_i))}{\partial \sigma(\theta^T x_i)} \cdot \frac{\partial \sigma(\theta^T x_i)}{\partial \theta^T x_i} \cdot \frac{\partial \theta^T x_i}{\partial \theta}$$

↙

$$= \frac{-1}{1 - \sigma(\theta^T x_i)} \cdot \sigma(\theta^T x_i) \cancel{\sigma(-\theta^T x_i)} \cdot x_i$$

↙

$$= -\sigma(\theta^T x_i) x_i$$

$$\Rightarrow \nabla_{\theta} E(\theta) = - \sum_{i=1}^n y_i \sigma(-\theta^T x_i) x_i - (1 - y_i) \sigma(\theta^T x_i) x_i \quad \star$$

$$= - \sum_{i=1}^n y_i \underbrace{\left(\underline{\sigma(-\theta^T x_i)} + \underline{\sigma(\theta^T x_i)} \right)}_1 x_i - \sigma(\theta^T x_i) x_i$$

$$= \sum_{i=1}^n \left(\underline{g(\theta^T x_i)} - \underline{y_i} \right) \underline{x_i}$$

Descent

Algorithm 24.1: Logistic Regression: Stochastic Gradient Ascent

LOGISTICREGRESSION-SGA ($\mathbf{D}, \eta, \epsilon$):

$w \rightarrow \theta$

$\theta \rightarrow g$

```

1 foreach  $x_i \in \mathbf{D}$  do  $\tilde{x}_i^T \leftarrow (1 \ x_i^T)$  // map to  $\mathbb{R}^{d+1}$ 
2  $t \leftarrow 0$  // step/iteration counter
3  $\theta \leftarrow (0, \dots, 0)^T \in \mathbb{R}^{d+1}$  // initial weight vector
4 repeat
5    $\theta \leftarrow \theta^t$  // make a copy of  $\theta^t$ 
6   foreach  $\tilde{x}_i \in \tilde{\mathbf{D}}$  in random order do
7      $\nabla(\theta, \tilde{x}_i) \leftarrow (y_i - g(\theta^T \tilde{x}_i)) \cdot \tilde{x}_i$  // compute gradient at  $\tilde{x}_i$ 
8      $\theta \leftarrow \theta + \eta \cdot \nabla(\theta, \tilde{x}_i)$  // update estimate for  $\theta$ 
9    $\theta^{t+1} \leftarrow \theta$  // update  $\theta^{t+1}$ 
10   $t \leftarrow t + 1$ 
11 until  $\|\theta^t - \theta^{t-1}\| \leq \epsilon$ 

```

} SGD
batch = 1

Suppose $Y \in \{c_1, \dots, c_k\}$

We model Y as a k -dim Multivariate Bernoulli R.V. taking values

$$\begin{aligned} c_1 \rightarrow e_1 &= (1, 0, \dots, 0) \\ c_2 \rightarrow e_2 &= (0, 1, \dots, 0) \\ \vdots &\quad \vdots \quad \vdots \\ c_k \rightarrow e_k &= (0, 0, \dots, 1) \end{aligned} \quad \left. \right\} \text{one-hot encoding}$$

Denote $\Pr(Y=e_i | X=x) = \underbrace{\pi_i(x)}_{\text{unknown probability}} \text{ for } i=1, \dots, k$

of class c_i given
datapoint x_i

$$\Rightarrow \sum_{i=1}^k \pi_i(x) = \sum_{i=1}^k \Pr(Y=e_i | X=x) = 1$$

$$\text{And, } \Pr(Y|X=x) = \prod_{j=1}^k (\pi_j(x))^{y_j}$$

In multiclass logistic regression, select c_k as a reference class and let log-odds of others w.r.t. c_k be linear in X

$$\log \left(\frac{P(Y=c_i | X=x)}{P(Y=c_k | X=x)} \right) = \log \left(\frac{\hat{\pi}_i(x)}{\hat{\pi}_k(x)} \right) = \Theta_i^T x = \Theta_{i0} + \Theta_{i1}x_1 + \dots + \Theta_{id}x_d$$

$$\Rightarrow \hat{\pi}_{i,}(x) = \exp(\Theta_i^T x) \hat{\pi}_k(x), \quad i=1, \dots, k-1$$

$$\sum_{j=1}^{k-1} \hat{\pi}_j(x) + \hat{\pi}_k(x) = 1$$

$$\sum_{j=1}^{k-1} \exp(\Theta_j^T x) \hat{\pi}_k(x) + \hat{\pi}_k(x) = 1$$

$$\hat{\pi}_k(x) \left(1 + \sum_{j=1}^{k-1} \exp(\Theta_j^T x) \right) = 1$$

$$\hat{\pi}_k(x) = \frac{1}{1 + \sum_{j=1}^{k-1} \exp(\Theta_j^T x)}$$

$$\hat{\pi}_i(x) = \frac{\exp(\Theta_i^T x)}{1 + \sum_{j=1}^{k-1} \exp(\Theta_j^T x)}$$

Setting $\Theta_k = 0$ so $\exp(\Theta_k^T x) = 1$, so

$$\hat{\pi}_i(x) = \frac{\exp(\Theta_i^T x)}{\sum_{j=1}^k \exp(\Theta_j^T x)}$$

Softmax function

Prediction criteria $\hat{y} = \arg \max_{c_i} \{\hat{\pi}_i(x)\} = \arg \max_{c_i} \left\{ \frac{\exp(\Theta_i^T x)}{\sum_{j=1}^k \exp(\Theta_j^T x)} \right\}$

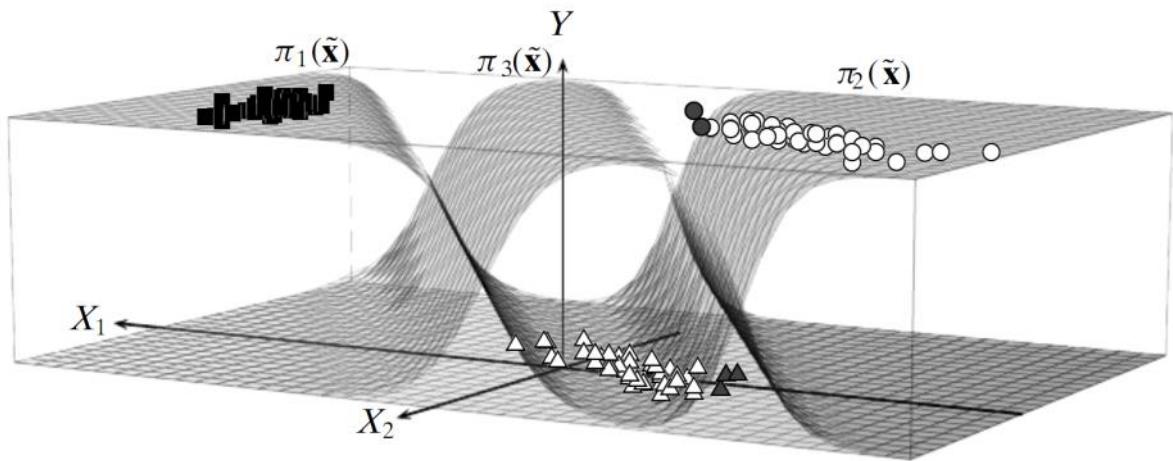


Figure 24.3. Multiclass logistic regression: Iris principal components data. Misclassified point are shown in dark gray color. All the points actually lie in the (X_1, X_2) plane, but c_1 and c_2 are shown displaced along Y with respect to the base class c_3 purely for illustration purposes.

Suppose $M: \mathbb{R}^d \rightarrow \{c_1, \dots, c_k\}$ is a classifier

To build a classifier requires a training set.
We assess its performance on a testing set.

22.1 Classification Performance Measures

D - testing set of n points $x_i \in \mathbb{R}^d$ with labels y_i

$\hat{y}_i = M(x_i)$ - predicted label in $\{c_1, \dots, c_k\}$

Error Rate : $\frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{y_i \neq \hat{y}_i\}}$

Accuracy : $\frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{y_i = \hat{y}_i\}}$

These metrics are global and do not consider how different classes contribute to the error.

Let $D_j = \{x_i \mid y_i = c_j\}$ points in class j

$$n_j = |D_j|$$

$D_j = \{x_i \mid \hat{y}_i = c_j\}$ points with predicted class j

$R_j = \{x_i \mid \hat{y}_i = c_j\}$ points with predicted class c_j

$$M_j = |R_j|$$

A confusion matrix is a $k \times k$ table N with

$N_{ij} = M_{ij}$ - number of points with predicted class c_i with true label c_j

$$= |R_i \cap D_j| = \{x_a \in D \mid \hat{y}_a = c_i \text{ and } y_a = c_j\}$$

		true class		
		c_1	c_2	c_3
predicted class	c_1	M_{11}	M_{12}	M_{13}
	c_2	M_{21}	M_{22}	M_{23}
	c_3	M_{31}	M_{32}	M_{33}

Precision: the class-specific precision of M for class c_i is

$$\text{prec}_i = \frac{M_{ii}}{M_i} = \frac{\text{true class } i \text{ predictions}}{\text{class } i \text{ predictions}} \quad (i = \text{max})$$

i.e. the fraction of our class i predictions that are correct.

Global precision is

$$\text{precision} = \sum_{i=1}^k \left(\frac{m_i}{n} \right) \text{prec}_i = \sum_{i=1}^k \frac{n_{ii}}{n} \quad (= \text{accuracy})$$

Recall:

$$\text{recall}_i = \frac{n_{ii}}{n_i} \quad (1 = \max)$$

the fraction of class i points we correctly identify

There is a trade-off between precision and recall...
for example, if $M(x_i) = y_2$ for all x_i , then recall_2 is

$$\frac{n_2}{n_2} = 1$$

but precision is $\frac{n_2}{n}$, which is likely low.

The class-specific F-measure tries to balance them

$$F_i = \frac{2}{\frac{1}{\text{prec}_i} + \frac{1}{\text{recall}_i}} = \frac{2 \text{prec}_i \cdot \text{recall}_i}{\text{prec}_i + \text{recall}_i} = \frac{2 \cdot \frac{n_{ii}}{m_i} \cdot \frac{n_{ii}}{n_i}}{\frac{n_{ii}}{m_i} + \frac{n_{ii}}{n_i}}$$

$\cancel{2n_{ii}^2}$ $\cancel{2n_{ii}}$

$$= \frac{Zn_{ii}^2}{n_{ii}n_i + n_{ii}m_i} = \frac{Zn_{ii}}{n_i + m_i}$$

for F_i to be near 1, we need both to be near 1.

Global F-score = $\frac{1}{k} \sum_{i=1}^k F_i$ AKA \bar{F} -Score

Binary Classification: Suppose there are only 2 classes, we can call

		True Class	
		pos.	neg.
Predicted class	pos	n_{11} = True Positives	n_{12} = False Positives
	neg	n_{21} = False Negatives	n_{22} = True Negatives

$$\text{prec}_p = \frac{TP}{TP+FP}$$

% of pos. prediction that are correct

$$\text{recall}_p = \frac{TP}{TP+FN}$$

% of pos. points that are correct

Perceiver Operating Characteristic (ROC)

Receiver Operating Characteristic (ROC)

Popular strategy for assessing a binary classifier which outputs a score value for the positive class for each point in the test set, e.g. Bayes classifier outputs a posterior probability $P(c_i|x_i)$.

Typically, $M(x_i) = \begin{cases} C_1, & \text{if score} \geq g \\ C_2, & \text{else} \end{cases}$ but g is a somewhat arbitrary threshold

ROC plots false positive rate vs. true positive rate

for each value of g from g^{\min} to g^{\max}

$$\frac{FP}{\text{neg. points}} \quad \frac{TP}{\text{pos. points}}$$

$\min S(x_i)$ $\max S(x_i)$

g -score for pos. class

At g^{\min} ...

		truth	
		P	N
pred.	P	O	O
	N	FN	TN

At g^{\max} ...

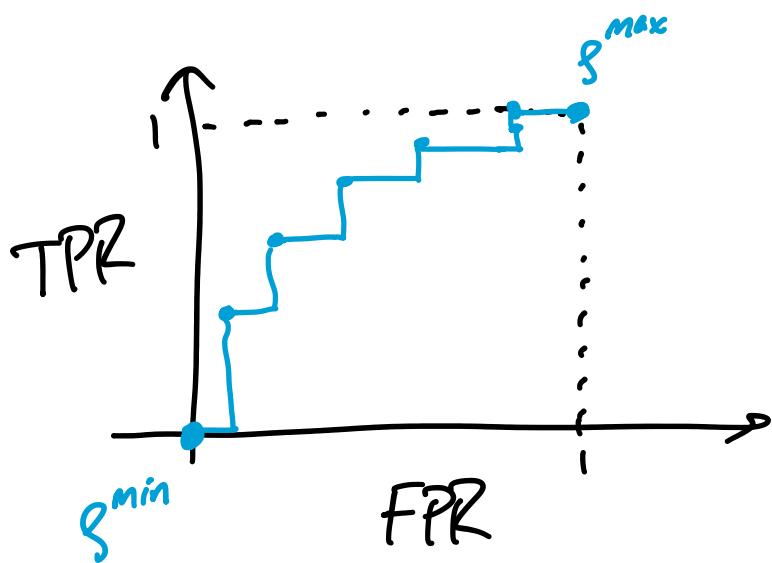
		truth	
		P	N
pred	P	TP	FP
	N	O	O

Ideal Classifier

	P	N
P	TP	O
N	O	TN

between...

$$R_s(\gamma) = \{x_i \in D \mid S(x_i) > \gamma\} - \text{points classified pos. for threshold } \gamma$$



Ideally, the curve should grow to 1 as quickly as possible



Area under ROC curve (AUC) is a one-number performance

metric in $[0,1]$ - the prob. classifier will rank a random pos. test point higher than a random neg. point.

imp: ROC/AUC is not sensitive to class imbalance