

Lecture 5 - 10/09/24

Logistic Regression

Classification Problems :

- So far regression problems were \rightarrow Given X to output Y
- Many times we wanna know which class a data point belongs to...
i.e.

Based on:

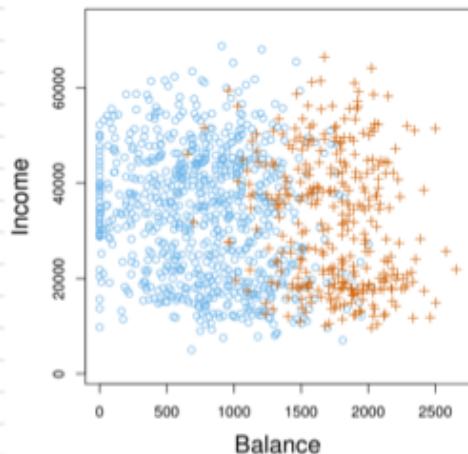
- Annual Income
- Monthly credit Card Balance

} Predict whether the individual will go broke or not!

The outcome of this is binary, either yes or no... this is exactly what logistic regression predicts.

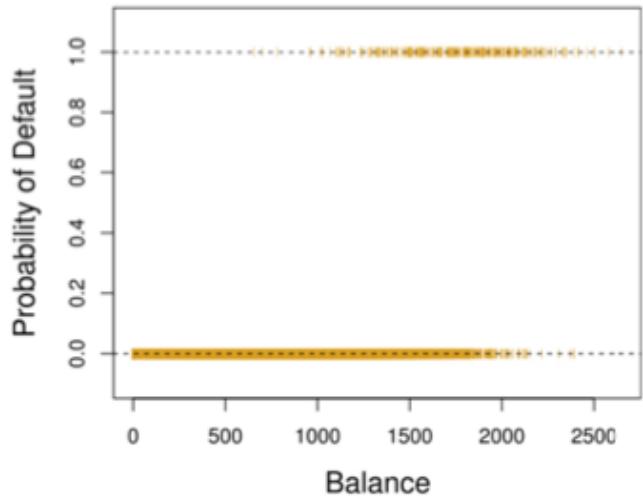
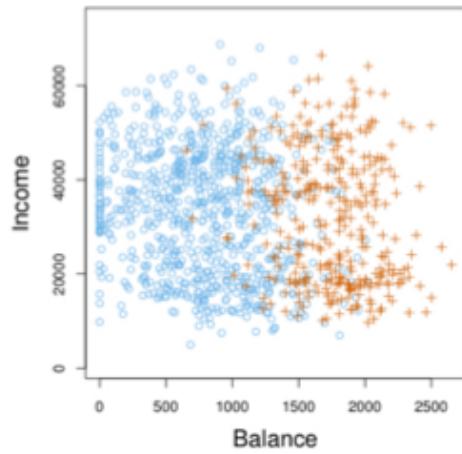
- Linear regression predicts continuous values
- Logistic regression predicts discrete values (binary values) !!!
- We then want our model to classify data points

↳ Classification Problem



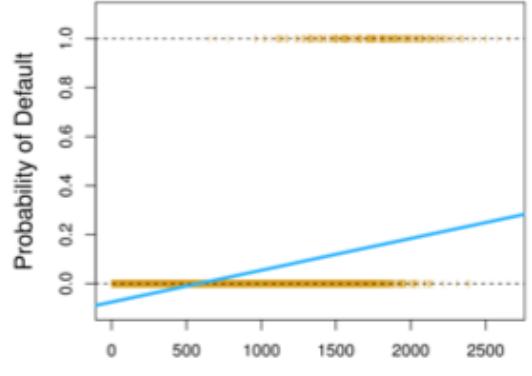
- We want a model that guesses the color of the datapoints
- We immediately see we are trying to classify!
- This is a binary classification

↳ A popular classification model is **Logistic Regression**



- Using Linear Regression in this example is bad!
- Modeling Linear Regression could give us:

$$\rightarrow P(x) \leq 0 \quad \text{or} \quad P(x) \geq 1$$

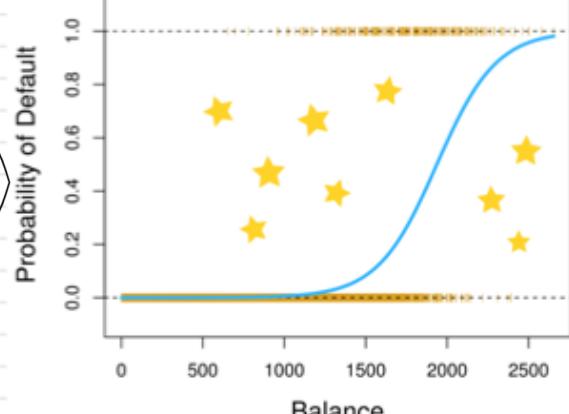


- We are instead seeking a model for $P(x)$:

such that:

Asymptotically gives $P(x) \rightarrow 1$

Asymptotically gives $P(x) \rightarrow 0$



Logistic Function

- What we are looking for is the **logistic function**:

$$p(X) = \frac{e^{\beta_0 + \beta_1 X}}{1 + e^{\beta_0 + \beta_1 X}}$$

↳ Sometimes we refer to it as the **Sigmoid Function**

↳ the **exponent** looks like our **Linear Regression**

↳ But there are nice properties to it...

WHAT IF $\beta_0 + \beta_1 X \rightarrow \infty$? $\frac{e^{\infty}}{1 + e^{\infty}} = 1$

WHAT IF $\beta_0 + \beta_1 X \rightarrow -\infty$? $\frac{e^{-\infty}}{1 + e^{-\infty}} = 0$

WHAT IF $\beta_0 + \beta_1 X = 0$? $\frac{1}{1+1} = \frac{1}{2}$



↳ Given this we need that the **logistic function always produces S-shape!!**

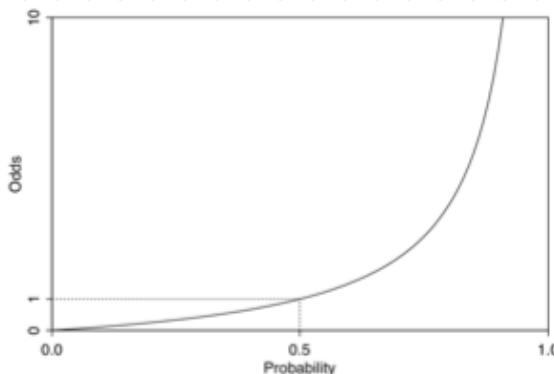
↳ Sometimes **odds** is thought as:

$$\frac{P(x)}{1 - P(x)}$$

↖ PROB OF X HAPPENING
↘ PROB OF X NOT HAPPENING

Probability of an odd

- Odds are positive real numbers
- $\text{odds}(A) > \text{odds}(B)$ exactly when $P(A) > P(B)$



Odds of 1 is a bit special in the sense that it has the same prob of that event happening as not happening.

It is a 50-50 chance.

- The odds in two groups can be compared by their odds ratio
 - Often use ratio to compare chance of event happening [under 2 ≠ odds]
 - $\text{odds} = 5 \rightarrow \text{odds of having cancer is 5 times higher for a smoker than a non-smoker}$
- COND 1: PATIENT SMOKES
COND 2: PATIENT DOESN'T SMOKES

- For the logistic function:

$$\text{odds} \rightarrow \frac{p(x)}{1-p(x)} = e^{\beta_0 + \beta_1 x}$$

- Taking the log gives us the log odds / logit:

$$\log\left(\frac{p(x)}{1-p(x)}\right) = \beta_0 + \beta_1 x$$

Logistic Regression has a logit that is linear in x .

Probabilities as function of feature X

RECALL:

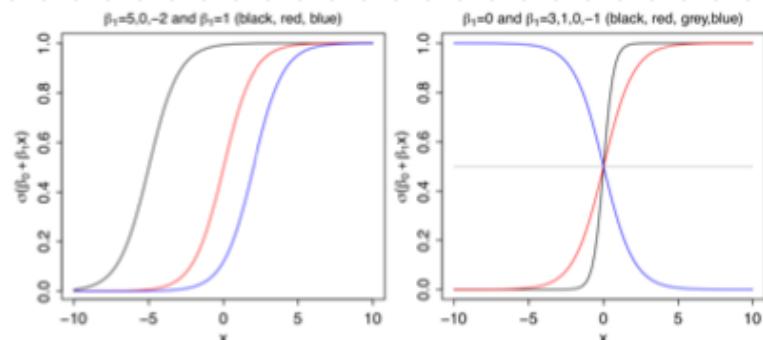
$$p(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

- we see that β_0 / β_1 determine where $p(x) = 0.5$

- INTERPRETATION OF COEFFICIENTS:

↳ A unit change in feature x_i means:

- Change β_i in log-odds
- Multiplicative change in odds by a factor e^{β_i} WHERE:
 - e^{β_i} is an odds ratio.
 - β_i is a log odds ratio.



↳ Given the nonlinear translation between odds & probability, it is hard to communicate the actual change in $p(x)$.

BUT:

- ! LOGIT-TRANSFORMATION IS MONOTONE!
- ! POSITIVE β_i GIVES POSITIVE CHANGE IN $P(x)$!
- ! A HIGHER β_i GIVES A STEEPER CURVE!

LOGISTIC REGRESSION WITH MULTIPLE FEATURES

- GIVEN x_1, \dots, x_n , Logistic Reg. models conditional prob. $y=1$ as:

$$P(y_i = 1 | x_i = x) = \frac{e^{x\beta}}{1 + e^{x\beta}}$$

WHICH THEN LEADS TO:

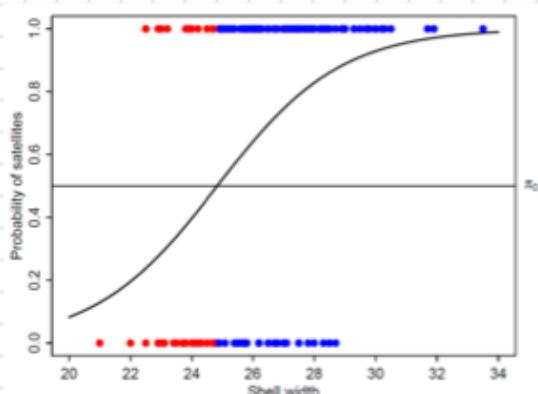
$$\frac{P(x)}{1 - P(x)} = e^{x\beta} \quad \xrightarrow{\text{LOG-LOGS}} \quad \log\left(\frac{P(x)}{1 - P(x)}\right) = x\beta$$

PREDICTION

- TO MAKE A PREDICTION $\hat{y} =$

WE CLASSIFY TO THE CLASS WITH HIGHEST PROBABILITY

$$\hat{y}_i = \begin{cases} 1, & P(y_i = 1 | x_i) > 0.5 \\ 0, & P(y_i = 0 | x_i) \leq 0.5 \end{cases}$$



- We see that our plotted function can output probs between 0 & 1
- We draw our decision line at 0.5 so that:
 - $P(x)$'s ≥ 0.5 will be classified as 1
 - $P(x)$'s < 0.5 will be classified as 0

		$\hat{Y} = 0$	$\hat{Y} = 1$	Total
$Y = 1$	16	95	111	
$Y = 0$	27	35	62	
Total	43	130	173	

THIS IS A CONFUSION MATRIX:

- TELLS PERFORMANCE OF OUR ALGORITHM

How to Read It?

- 16 points were classified as 0 by our model but instead they were 1's
- 35 points the other way around

While 95 & 27 were classified correctly

We can also retrieve the accuracy of the model by doing:

$$\frac{\text{Correctly Classified}}{\text{Total Points}} = \frac{27+95}{27+95+16+35} = 0.705$$

THE DECISION BOUNDARY WAS 0.5, BY CHANGING IT WE'LL GET DIFFERENT RESULT BUT HOW DO WE CHOOSE THIS DECISION BOUNDARY?

Training error vs Test error

- What about **quality of our model**?
- **Test error rate** is one measure:

PROBABILITY OF MAKING A WRONG PREDICTION!

$$\frac{1}{n} \sum_{i=1}^n I(y_i \neq \hat{y}_i)$$

- By computing this on **training data** → **training error rate**
- By computing this on **new data** → **test error rate**

Estimating the regression coefficients

- Params are estimated by **maximum likelihood**

BINOMIAL LIKELIHOOD: $L(\beta, y) = \prod_{i:y_i=1} p(x_i) \prod_{j:y_j=0} (1 - p(x_j))$

LOG LIKELIHOOD: $L(\beta, y) = \sum_{i=1}^n (y_i \log p(x_i) + (1 - y_i) \log(1 - p(x_i)))$

- WHAT WE WANNA DO IS **MAXIMISE LIKELIHOOD**:

ITERATIVELY REWEIGHTED LEAST SQUARES METHOD

- Model Summary output usually reports a **Wald-Test** for testing $\beta_j = 0$:

$$\frac{\hat{\beta}_j}{SE(\hat{\beta}_j)} \sim N(0,1)$$

THE **LIKELIHOOD-RATIO TEST** IS A BETTER, MORE GENERAL, TEST FOR COMPARING MODELS:

$$-2 \log \frac{L(\hat{\beta}_0)}{L(\hat{\beta}_1)} = -2 \left[\log L(\hat{\beta}_1) - \log L(\hat{\beta}_0) \right] \sim \chi^2_{df_1 - df_0}$$

DEFIN CAUCED DEVANCE TESTS AS THEY **COMPARE DEVANCE**.

DIAGNOSIS PLOTS FOR MODEL CHECKIN ARE SIMILAR TO LINEAR REGRESSION.

Model Probability

- SINCE Y IS BINARY,

$$\mathbb{E}(Y_i = 1 | x) = 1 \cdot P(Y_i = 1 | x) + 0 \cdot P(Y_i = 0 | x) = P(Y_i = 1 | x)$$

SO **WHAT WE SEE IS THAT BY MODELLING THE PROBS, WE ALSO MODEL THE EXPECTATION OF Y** .

Generalised Linear Models

E = EXPECTATION

- In linear regression, we model $E(y|x)$ directly as linear combination
- In logistic regression, we model the logit of $E(y|x)$ as linear combination
- We can then say:

BOTH ARE GENERALISED LINEAR MODELS \rightarrow WHERE FUNCTION OF THE MEAN IS LINEAR

$$g(E(y|x)) = X\beta$$

GIVEN THAT g = MONOTONE FUNCTION

$$E(y|x) = g^{-1}(X\beta)$$

Multinomial Logistic Regression

- Imagine y_i taking values among $1, \dots, K$
- Using one-hot encoding:

$$y_i = (0, 1, 0, \dots, 0) \quad \text{WHERE THE ONE STANDS FOR THE RIGHT CLASS}$$

It is clear that y_i follows a multinomial distribution

Its probability vector is (p_{i1}, \dots, p_{iK})

So now we have $K-1$ probabilities.

- Generalization to multiple classes:

- Select an arbitrary class as baseline (K here)

- Consider for another class k the odds of being k rather than K

K = BIG K

k = small k

$$\log = \frac{P(y=k|X=x)}{P(y=K|X=x)} = X\beta^k$$

THE PARAMS, $\beta^1, \dots, \beta^{K-1}$ ARE EACH A VECTOR OF LENGTH $p+1$

PROBS ARE:

$$P(y=k|X=x) = \frac{e^{X\beta^k}}{1 + \sum_{c=1}^{K-1} e^{X\beta^c}}$$

This is for small k

$$P(y=K|X=x) = \frac{1}{1 + \sum_{c=1}^{K-1} e^{X\beta^c}}$$

This is for big K

They sum to 1

- Most common choice = Predict that y belongs to class with HIGHEST PROB:

$$\hat{y} = \operatorname{argmax}_k P(y=k|X=x)$$

- This is obviously NOT THE ONLY WAY TO GENERALIZE GENERAL CATEGORIES.