

Foundations of Data Science: Logistic regression



THE UNIVERSITY *of* EDINBURGH
informatics

FOUNDATIONS
OF
DATA
SCIENCE

Overview

- Principle of Logistic Regression
- Interpretation of Logistic Regression coefficients
- Multiple Logistic Regression
- Logistic Regression as a classifier
- (Maximum likelihood estimation of Logistic regression coefficients)

The background of the slide features a stylized globe on the left side, partially obscured by a grid of binary code (0s and 1s) that recedes into the distance, creating a sense of depth and digital connectivity. The overall color palette is a mix of light blues and purples.

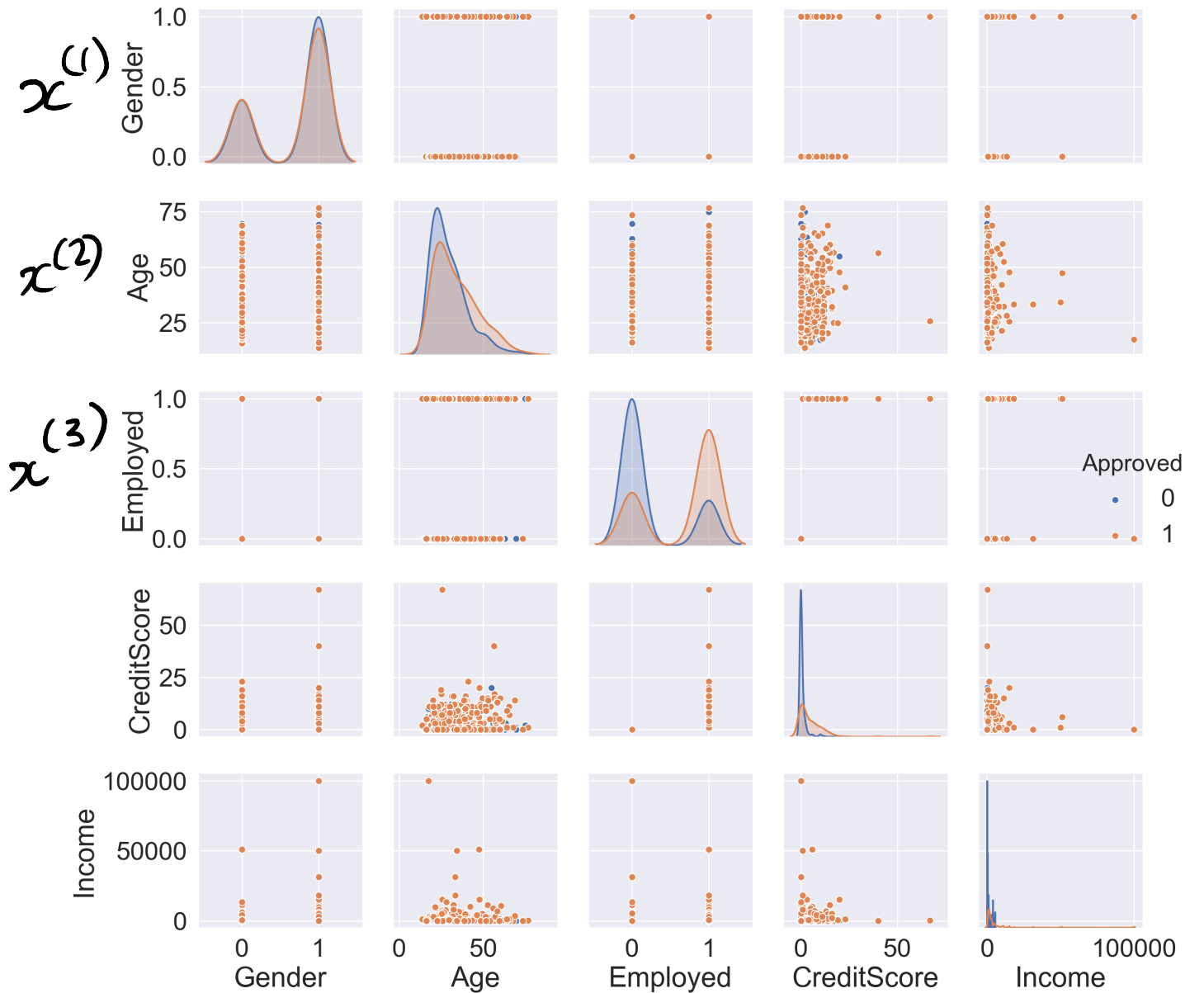
Foundations of Data Science: Logistic regression - Principle of logistic regression

Supervised classification

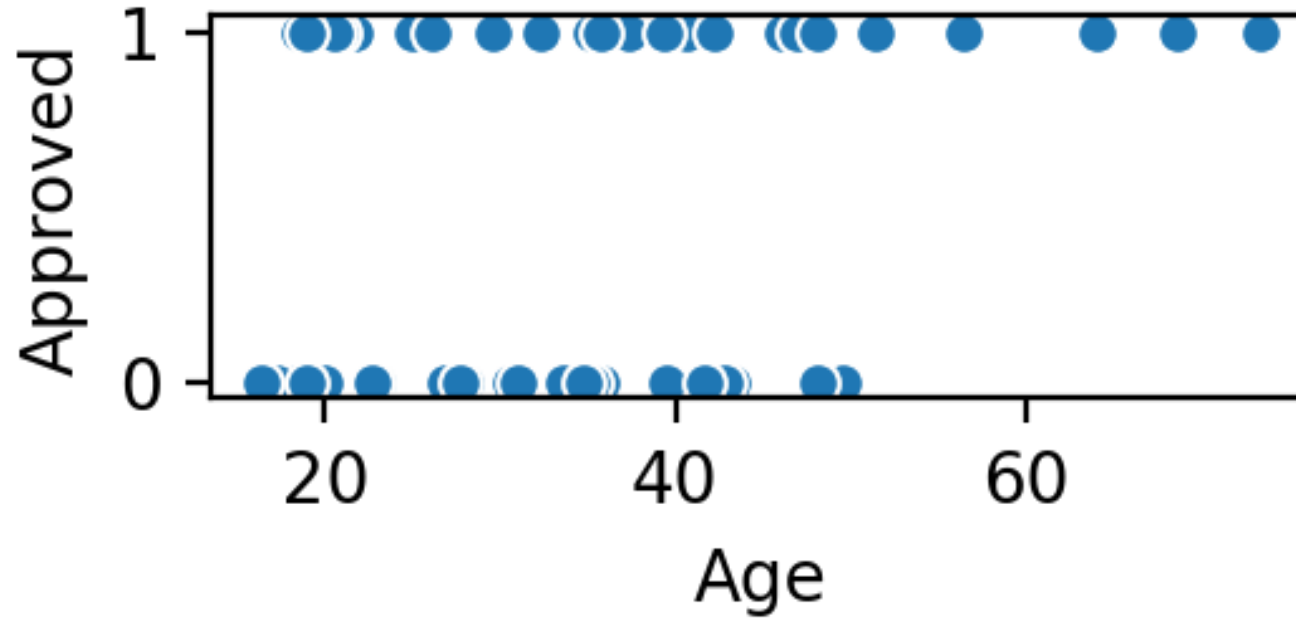
Binary (or dichotomous) response variable:
Credit

Approved

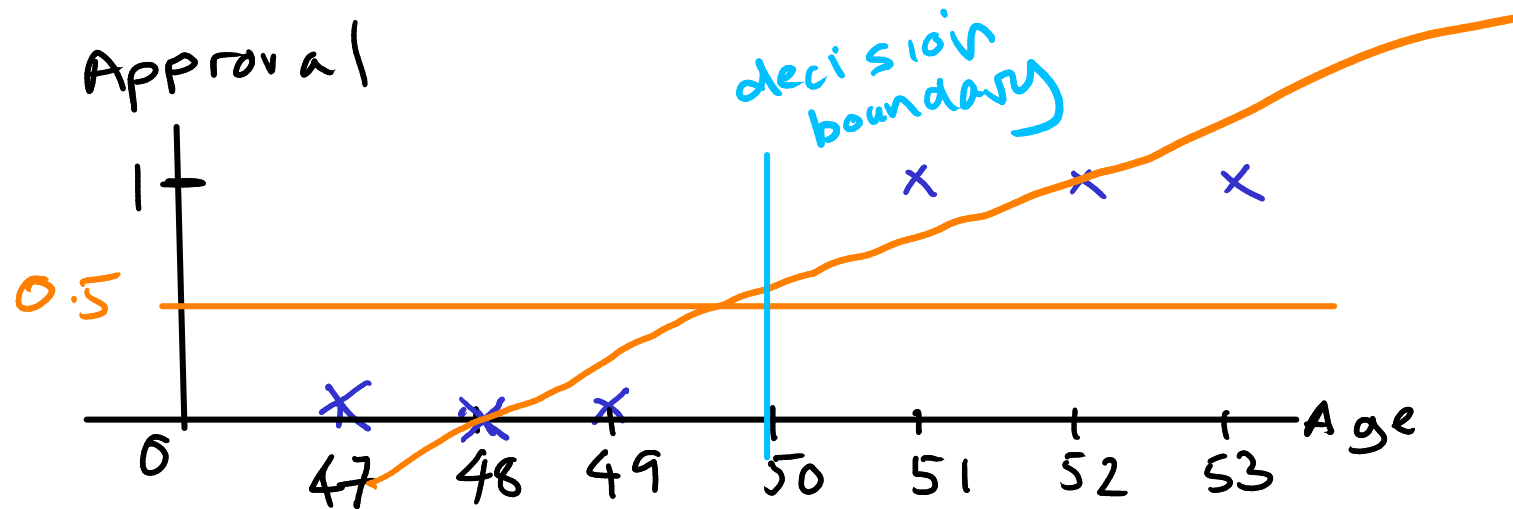
Not approved



Classification task on one continuous variable



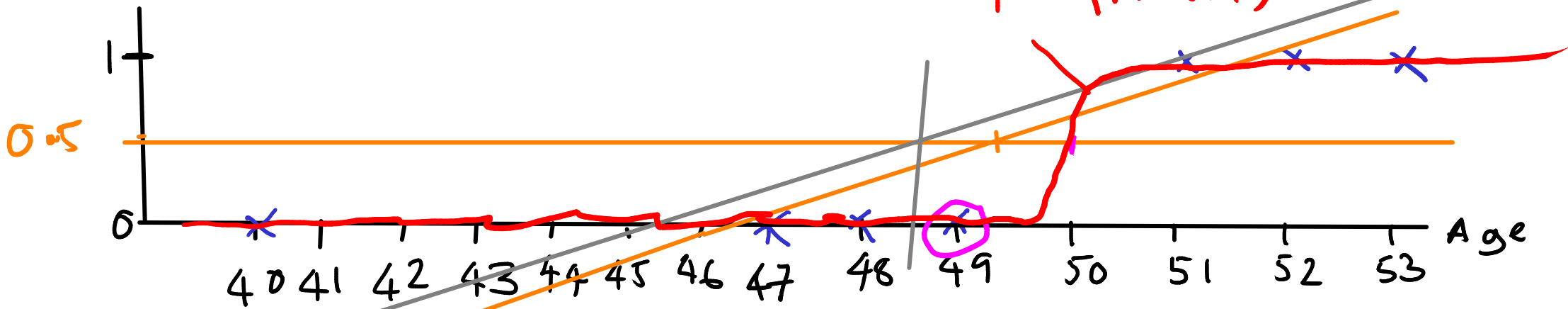
Exercise :



- Draw a Linear regression line through this data
- Convert the linear regression prediction to a predicted class label (0/1)
- Where is the decision boundary?
- How many classification errors are there?

Approval

$P(\text{Approval})$



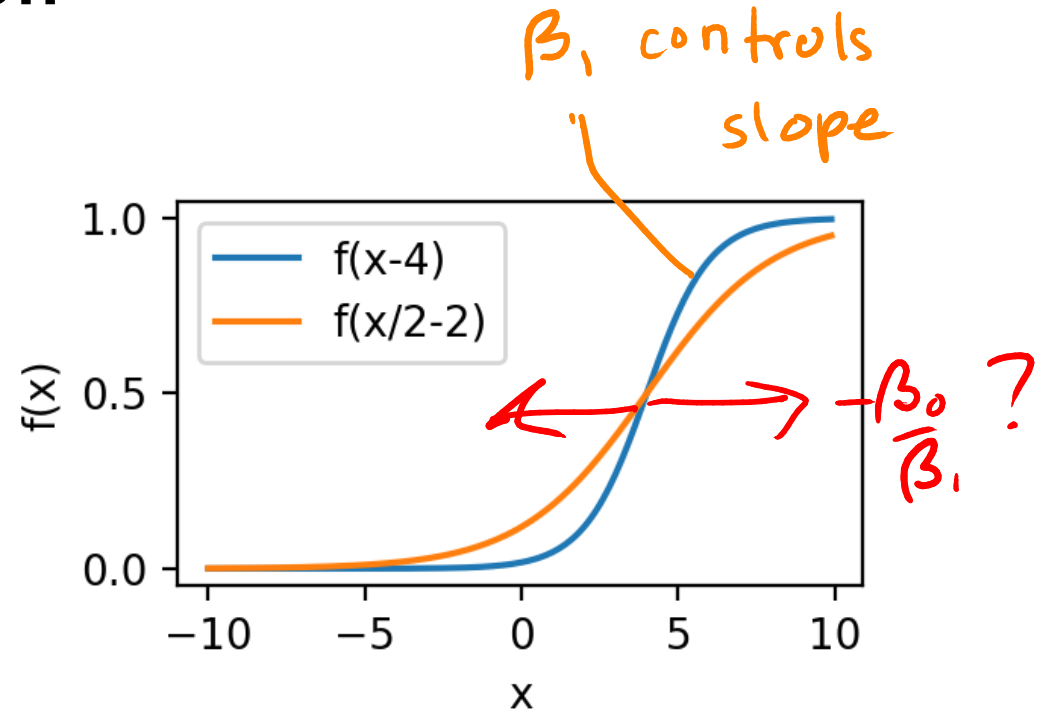
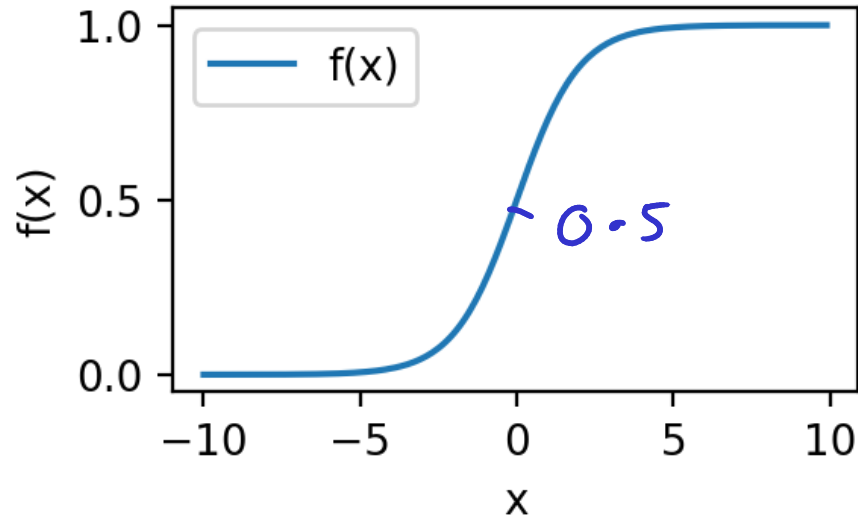
(a) Draw Linear regression line through this toy data

(b) Convert the linear regression prediction to a predicted class label (0/1)

(c) Where is the decision boundary?

(d) How many classification errors are there?

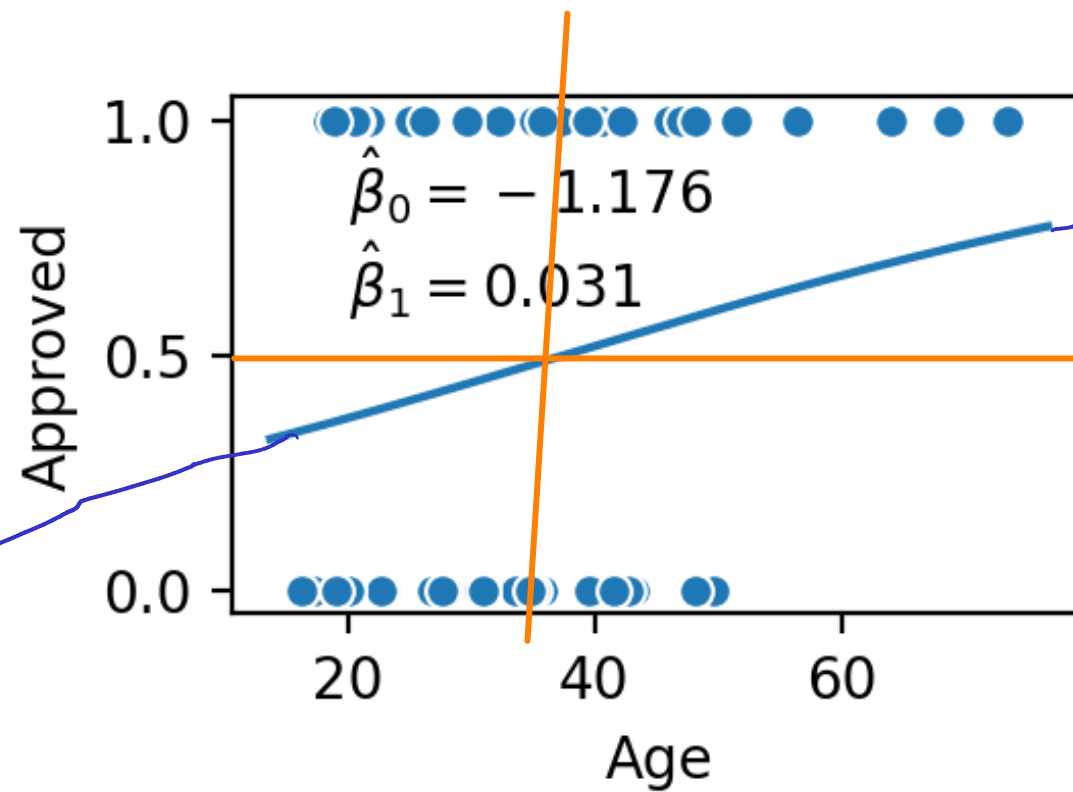
Logistic function



$$f(x) = \frac{e^x}{1 + e^x} = \frac{1}{1 + e^{-x}}$$

$$P(\hat{y} = 1 \mid x = x) = f(\hat{\beta}_0 + \hat{\beta}_1 x) = \frac{1}{1 + e^{-\hat{\beta}_0 - \hat{\beta}_1 x}}$$

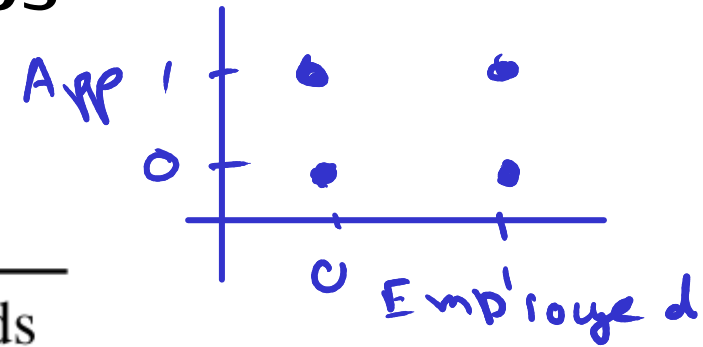
Application to continuous variable in credit example



Binary variables: odds and odds ratios

$$P(Y = y \mid X = x)$$

	Approved	Not approved	Approval odds
Employed			
0	0.25	0.75	0.34
1	0.71	0.29	2.42




$$OR(x) = \frac{2.42}{0.34} = 7.09$$

Effect size
609 %

$y \in \{ \text{"Not approved"}, \text{"Approved"} \}$
 $x \in \{ \text{"Not Emp."}, \text{"Emp."} \}$

$$\text{Odds (Success)} = \frac{P(\text{Success})}{P(\text{Failure})} = \frac{P(\text{Success})}{1 - P(\text{Success})}$$

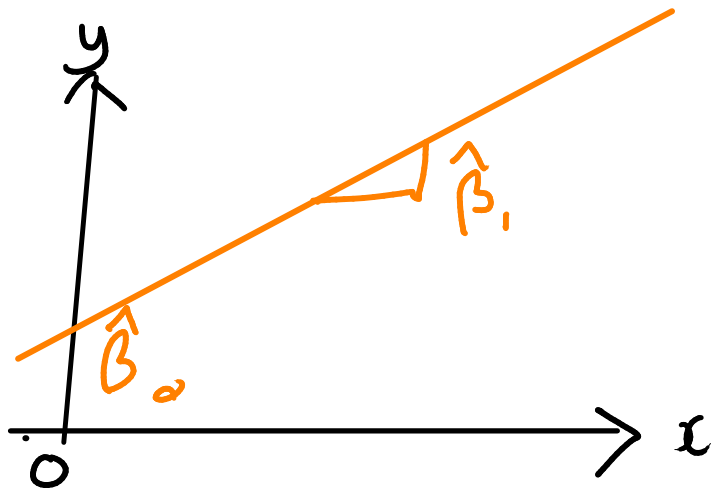
$$\text{Odds ratio } OR(x) = \frac{\text{Odds (Success) } | x = \text{True}}{\text{Odds (Success) } | x = \text{False}}$$

The background of the slide features a stylized globe on the left side, partially obscured by a grid of binary code (0s and 1s) that recedes into the distance, creating a sense of depth. The overall color palette is a mix of light blues and purples.

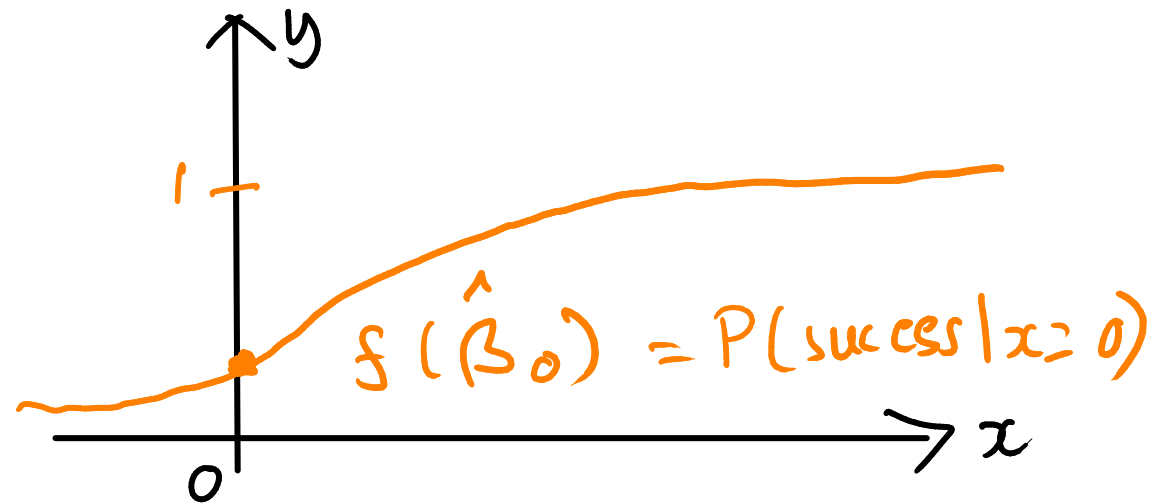
**Foundations of Data Science:
Logistic regression -
Interpretation of logistic regression
coefficients**

Interpretation of $\hat{\beta}_0$

Lin reg

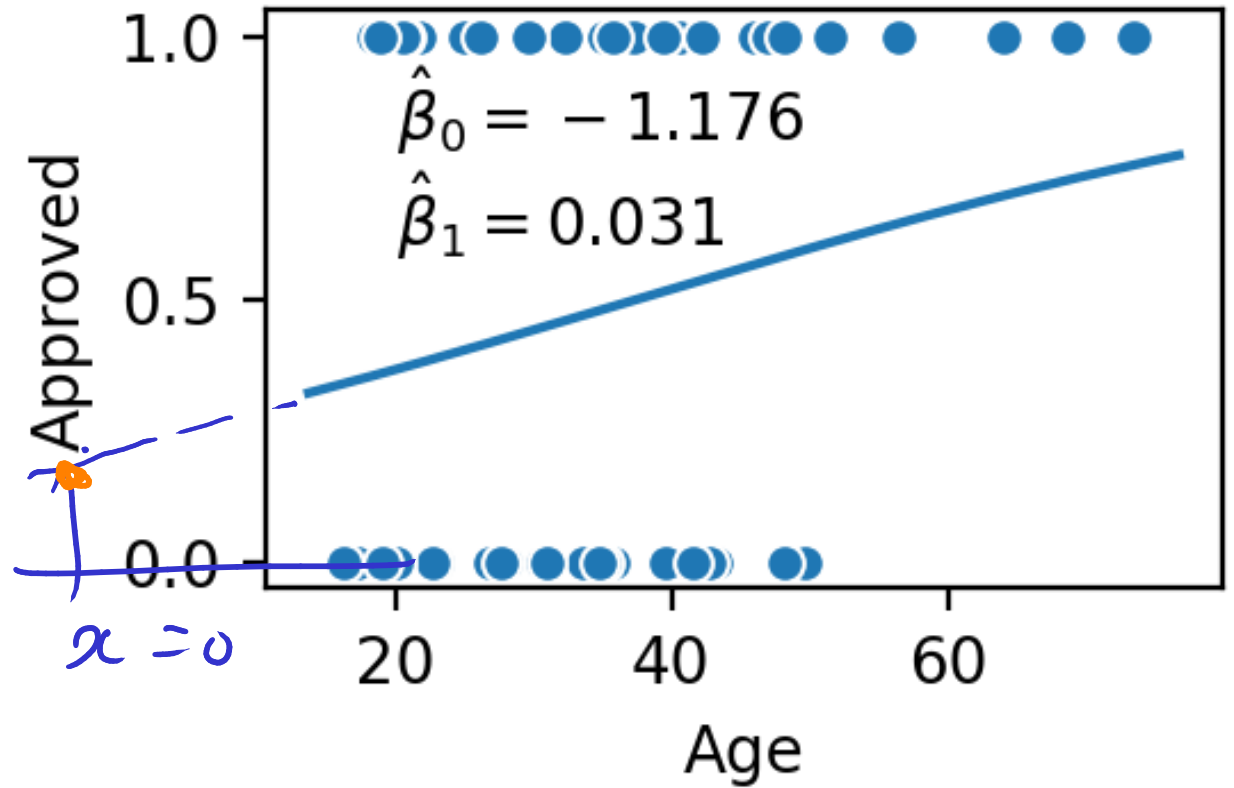


Log. reg



$$\begin{aligned} f(\hat{\beta}_0 + 0 \cdot \hat{\beta}_1) &= f(\hat{\beta}_0) \\ &= \frac{1}{1 + e^{-\hat{\beta}_0}} \end{aligned}$$

$$f(\hat{\beta}_0) = f(-1.176) \\ = 0.236$$



Log odds

$$\text{Log Odds (Success)} = \ln \frac{P(\text{Success})}{P(\text{Failure})}$$

$\swarrow \log_e$

$$\begin{aligned} \text{Log odds} + 1 & \\ \Rightarrow \text{odds} \times e & \end{aligned}$$

$$\ln \frac{P(\text{Success})}{1 - P(\text{Success})} = \ln \frac{p}{1-p}$$

p	odds	Log odds
0.5	1	0
> 0.5	> 1	> 0
< 0.5	< 1	< 0
1	$\rightarrow \infty$	$\rightarrow \infty$
0	0	$\rightarrow -\infty$

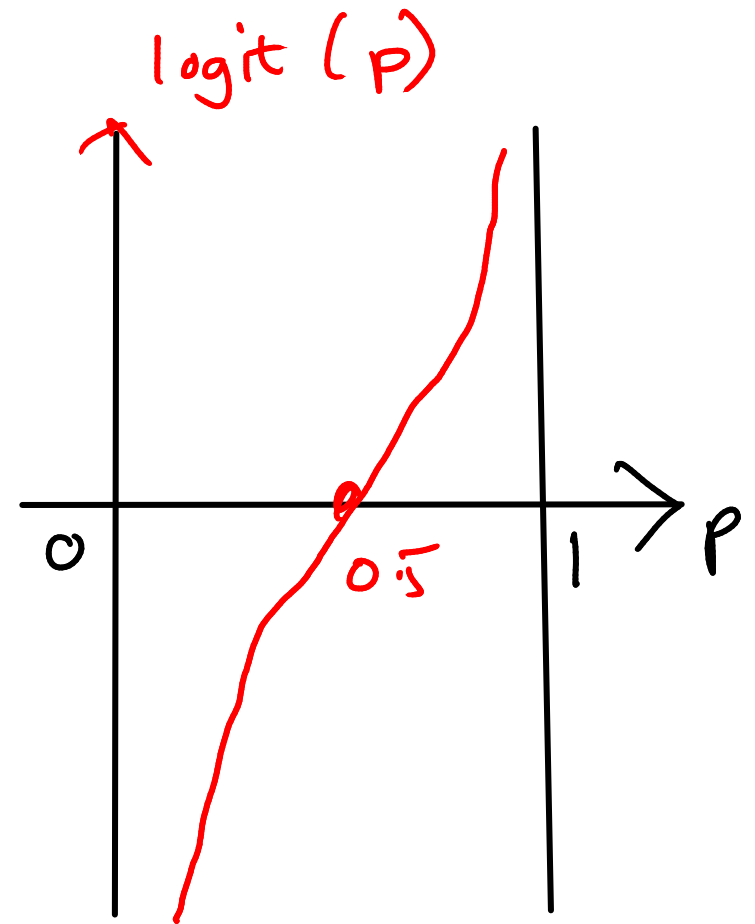
Logit scale

Log odds +1 \Rightarrow Odds increase by factor e

logistic unit = logit

$$\hat{\beta}_0 = -1.176 \text{ logits}$$

$$\text{logit}(p) = \ln \frac{p}{1-p}$$



Logistic Regression in terms of log odds

Success $P(Y=1|x) = f(\beta_0 + \beta_1 x) = \frac{1}{1 + e^{-\beta_0 - \beta_1 x}}$ ①

Failure $P(Y=0|x) = 1 - f(\beta_0 + \beta_1 x) = 1 - \frac{1}{1 + e^{-\beta_0 - \beta_1 x}}$

$$= \frac{e^{-\beta_0 - \beta_1 x}}{1 + e^{-\beta_0 - \beta_1 x}}$$
 ②

① / ②

Odds $\frac{P(Y=1|x)}{P(Y=0|x)} = \frac{1}{e^{-\beta_0 - \beta_1 x}} = e^{\beta_0 + \beta_1 x}$

Log odds $\ln \frac{P(Y=1|x)}{P(Y=0|x)} = \beta_0 + \beta_1 x = \text{logit}(P(Y=1|x))$

Interpretation of $\hat{\beta}_1$

$$\begin{aligned}\text{Odds}(x) &= e^{\hat{\beta}_0 + \hat{\beta}_1 x} \\ &= e^{\hat{\beta}_0} e^{\hat{\beta}_1 x}\end{aligned}$$

$$x = \{0, 1\} \quad \text{OR}(x) = \overbrace{\text{Odds}(1)}^{\text{Odds}(0)}$$

$$\text{OR}(x) = e^{\hat{\beta}_1}$$

$$\log \text{OR}(x) = \hat{\beta}_1$$

credit e.g. $\text{OR}(\text{Age}) = e^{0.03} \approx 1.03$



**Foundations of Data Science:
Logistic regression -
Multiple logistic regression**

Principle of multiple logistic regression

Predictor variables $x^{(1)}$: Age
 $x^{(2)}$: Employment
= {0, 1}

$$P(Y=1 \mid x^{(1)}, x^{(2)}, \dots) \\ = f(\beta_0 + \beta_1 x^{(1)} + \beta_2 x^{(2)} + \dots)$$

Multiple logistic regression applied to the credit example

	Variable	Coefficient	Odds or OR
$\hat{\beta}_0$	Intercept	-1.969	0.140 ← odds
$\hat{\beta}_1$	Age	0.029	1.030 ← OR
$\hat{\beta}_2$	Employed	1.881	6.562 ← OR

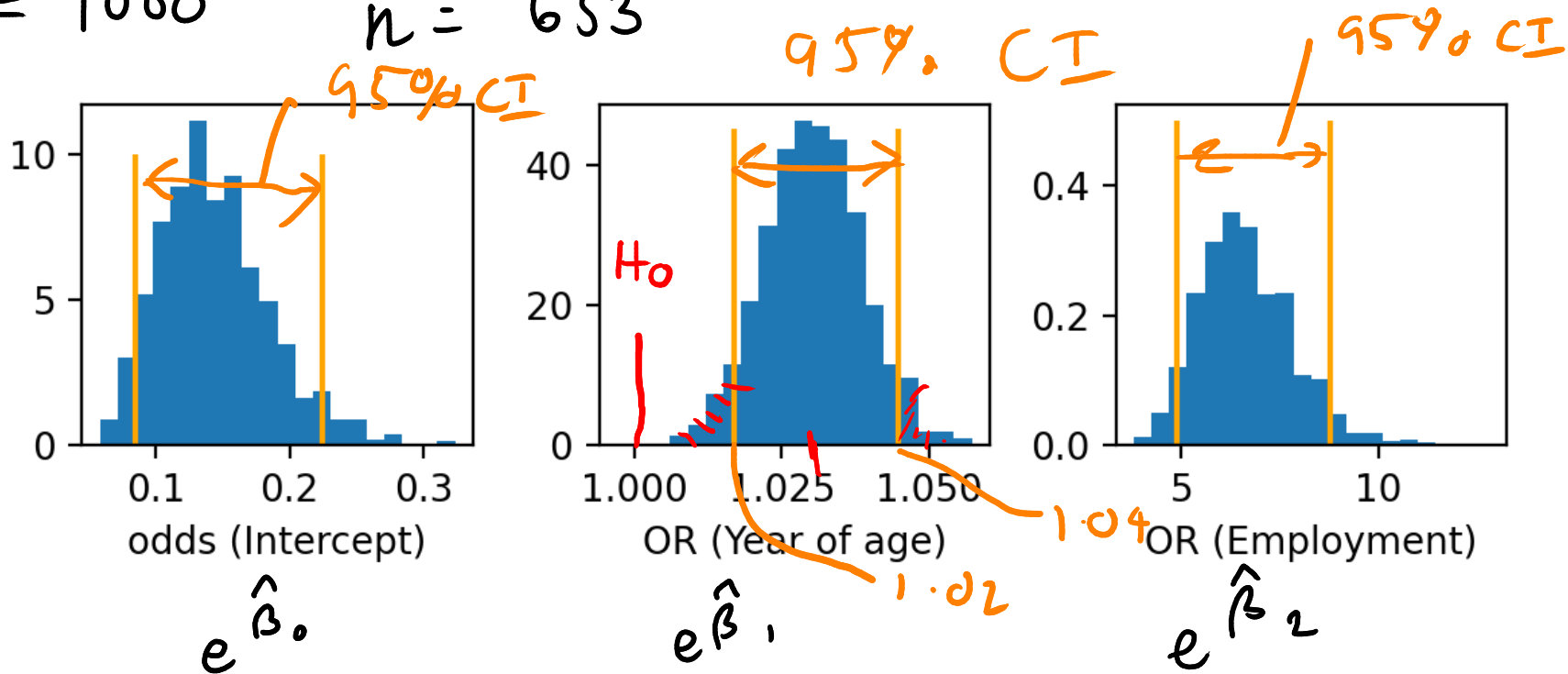
log[↑] odd
logits

$e^{+\hat{\beta}}$

Bootstrap confidence intervals

$B = 1000$

$n = 653$



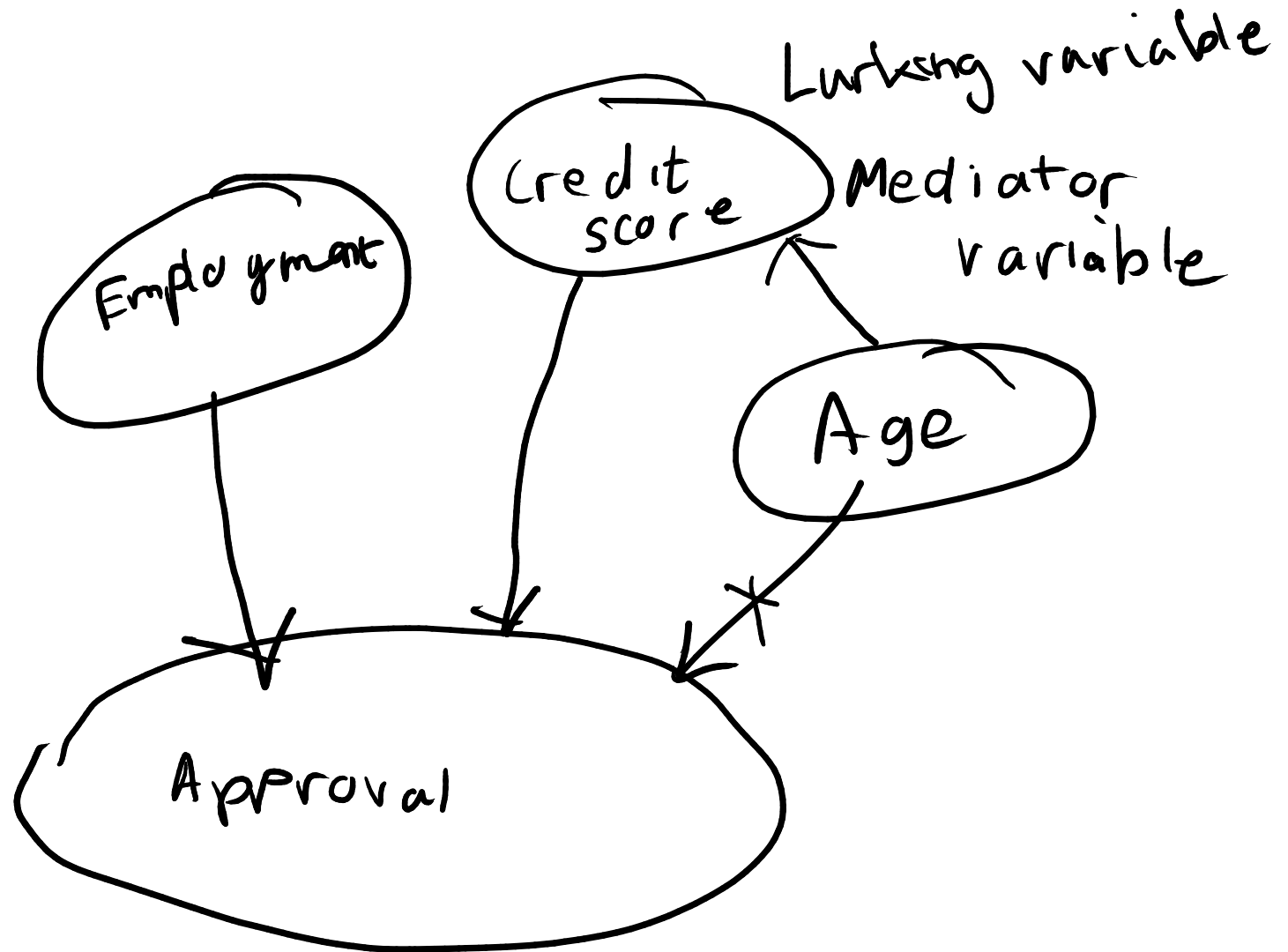
Does age affect credit approval?

H_0 : age does not affect credit approval $\Rightarrow e^{\hat{\beta}_1} = 1$

H_a : " " affect credit approval in some way.

Discussion question

Can you think of any problems in the reasoning that we've used to suggest age and credit approval are related?





**Foundations of Data Science:
Logistic regression -
The logistic regression classifier**

Converting logistic regression to a classifier

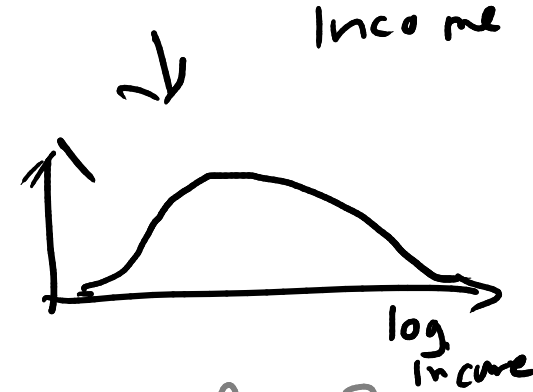
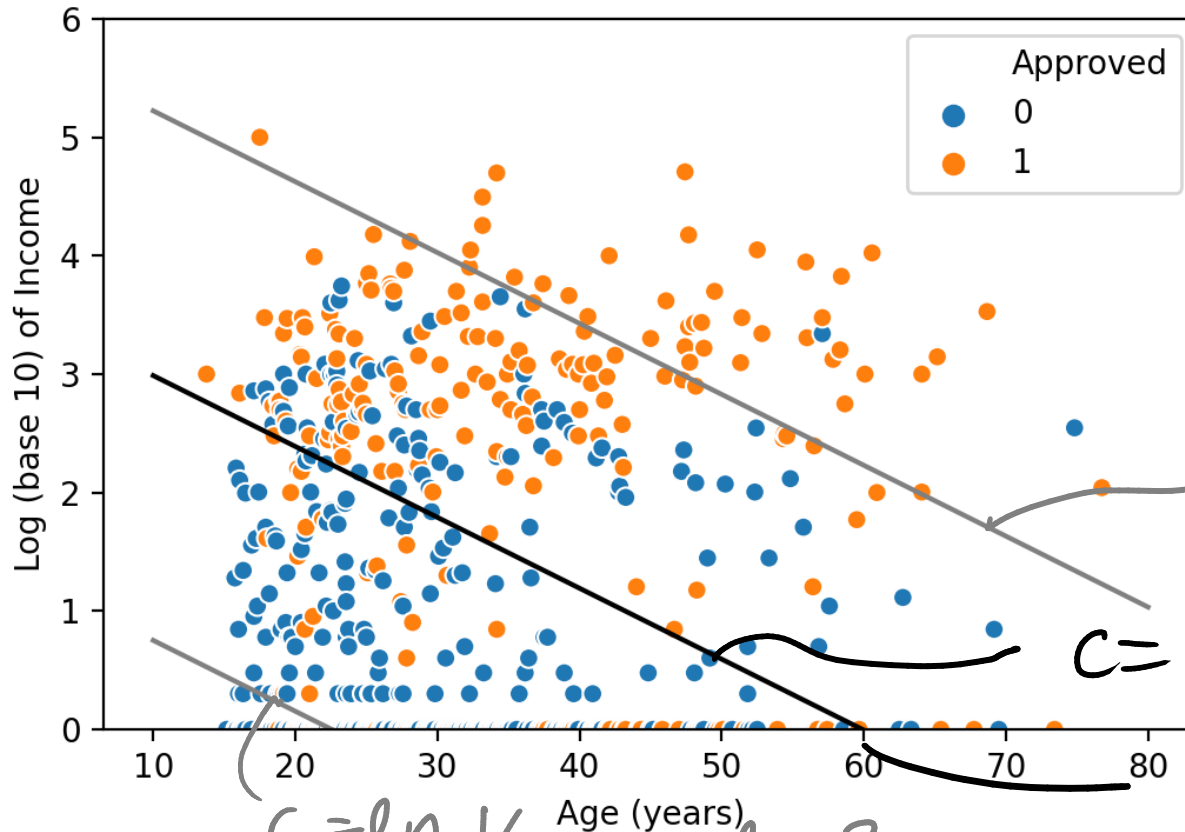
- Fit logistic regression model
- Set threshold c in terms of log odds and apply to predicted log odds:

$$\hat{\beta}_0 + \hat{\beta}_1 x^{(1)} + \hat{\beta}_2 x^{(2)} + \dots \geq c \Rightarrow \hat{y} = 1$$

$$\hat{\beta}_0 + \hat{\beta}_1 x^{(1)} + \hat{\beta}_2 x^{(2)} + \dots < c \Rightarrow \hat{y} = 0$$

$$c = 0 \Rightarrow \text{odds of 1} \Rightarrow p = 0.5$$

Decision boundary



Handwritten notes:

$$c = \ln 3$$

odds e^3

$$p = \text{logit}(\ln^3)$$

Handwritten note:

$$c = \ln \frac{1}{3} = -\ln 3$$

Decision boundary

$$\beta_0 + \beta_1 x^{(1)} + \beta_2 x^{(2)} + \dots > c \Rightarrow \hat{y} = 1$$

$$\beta_0 + \beta_1 x^{(1)} + \beta_2 x^{(2)} + \dots \leq c \Rightarrow \hat{y} = 0$$

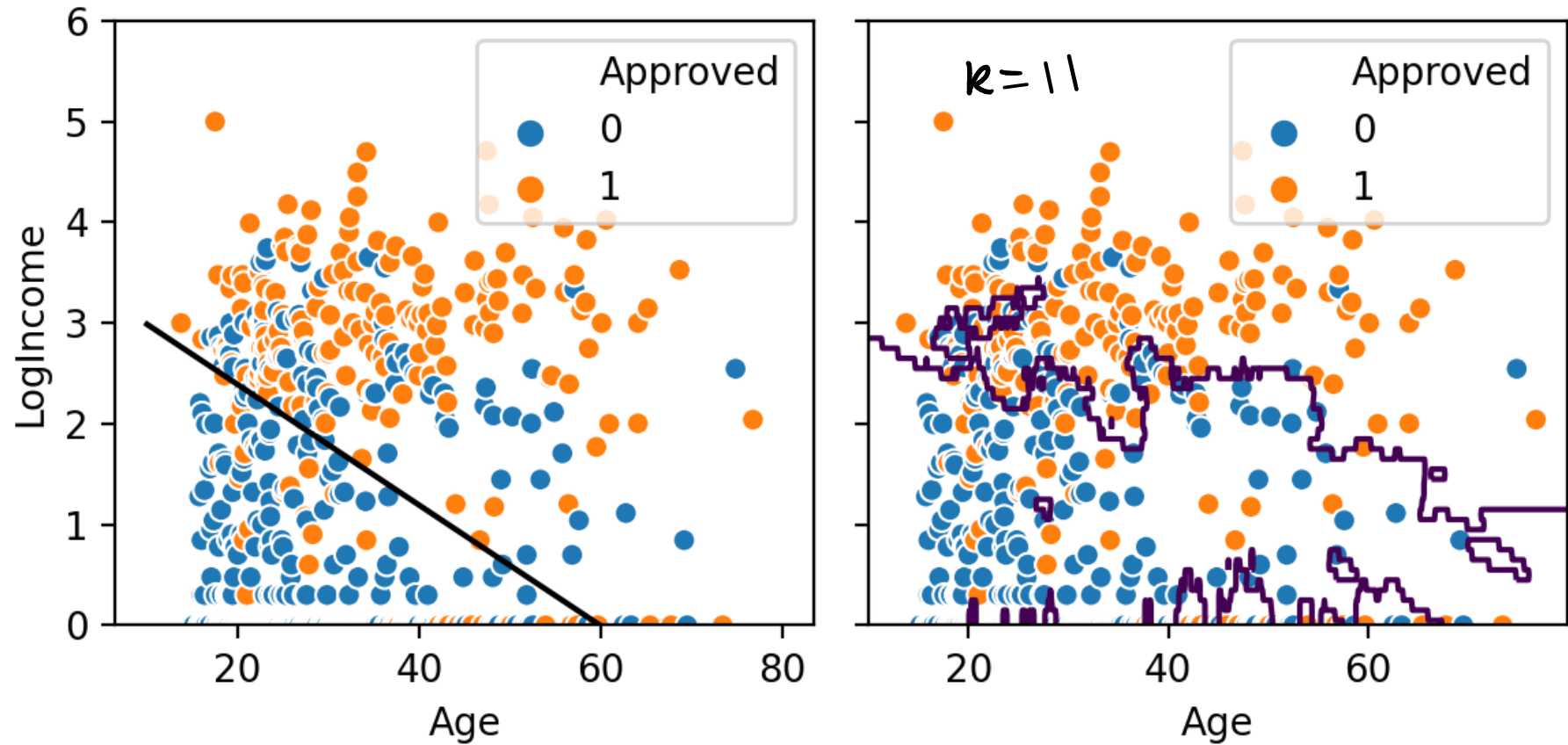
Ethics: logistic regression can be transparent

Credit scoring system:

- If you are in employment you score 1.625, if not you score 0
- Multiply your age by 0.029 and add the result to your score
- Round your income to the nearest 1000.
Multiply the number of zeros in this figure by 0.320
and add the result to your score
- If you scored more than 2.246, your credit will be approved

Cf. "Promote Values of Transparency, Autonomy and Trustworthiness" (Vallor, 2018)

Logistic regression versus k-NN




Decision boundary, flexibility/over-fitting, transparency

Standardised input variables

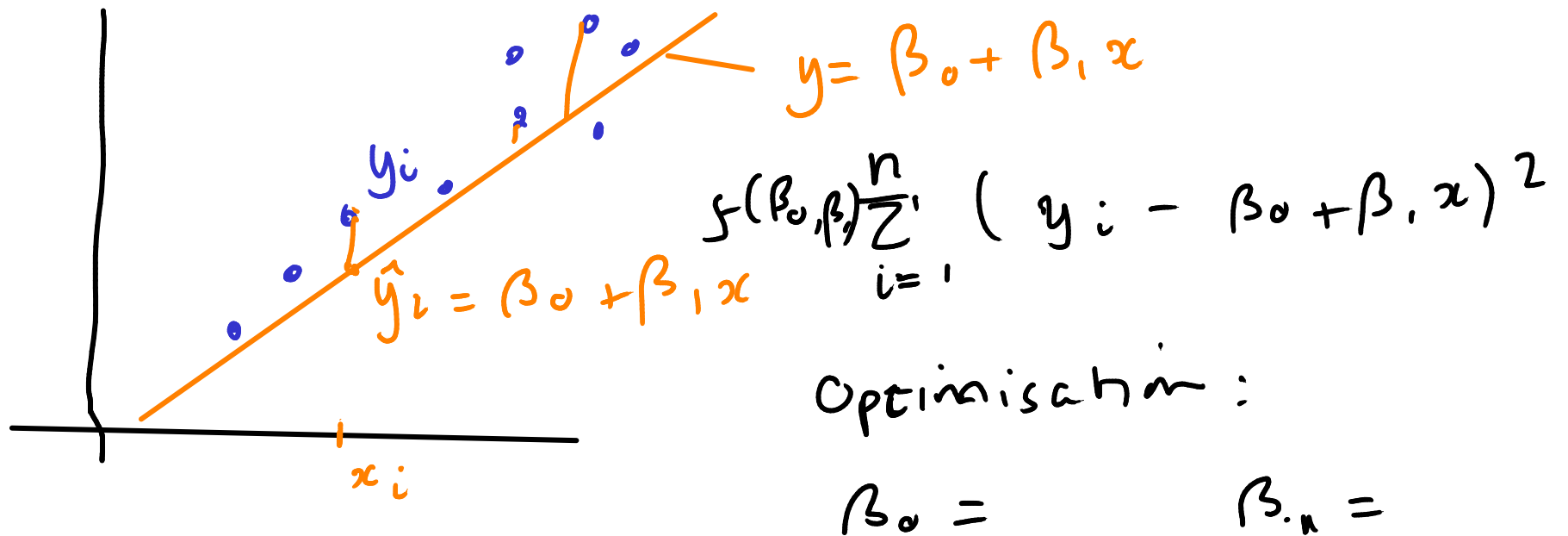
Summary

- Interpret $\hat{\beta}_0$ and $\hat{\beta}_1$ in terms of log odds
- Extend logistic regression to multiple variables
- Use logistic regression as a classifier
- Not covered (yet): derivation of logistic regression from principle of max likelihood

The background of the slide features a blue-toned digital aesthetic. It includes a faint, glowing globe on the left side and a pattern of binary code (0s and 1s) that appears to be receding into the distance, creating a sense of depth and data flow.

**Foundations of Data Science:
Logistic regression -
Maximum likelihood estimation of
logistic regression coefficients**

Principle of Maximum Likelihood



Adjust coefficients so as to maximise the likelihood of the data.

⇒ Expression for max. likelihood
Optimise w.r.t β_0, β_1, \dots

Likelihood of one point

$$P(Y_i = 1 \mid X_i = x_i) = f(\beta_0 + \beta_1 x_i) \quad \textcircled{1}$$

$$P(Y_i = 0 \mid X_i = x_i) = 1 - f(\beta_0 + \beta_1 x_i) \quad \textcircled{2}$$

$$\Rightarrow P(Y_i = y_i \mid X_i = x_i) =$$

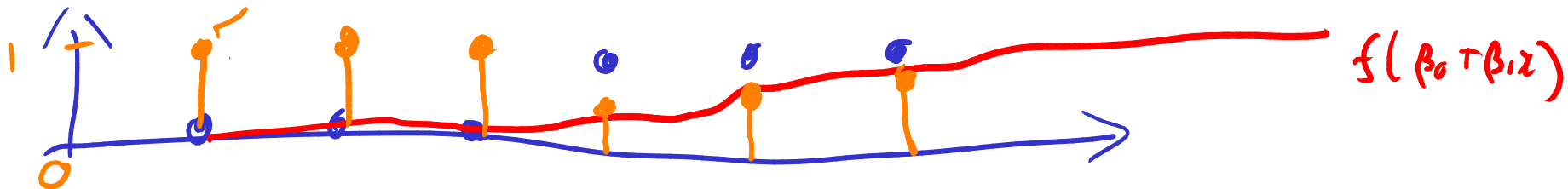
$$y_i \underset{\textcircled{1}}{f(\beta_0 + \beta_1 x_i)} + (1 - y_i) \underset{\textcircled{2}}{(1 - f(\beta_0 + \beta_1 x_i))}$$

Likelihood of data given model

Assumption: responses are independent, given value of predictor variables

$$\begin{aligned} P(Y=y \mid X=\underline{x}) &= P(Y=y_1 \mid X=x_1) P(Y=y_2 \mid X=x_2) \dots \\ &= \prod_{i=1}^n P(Y=y_i \mid X=x_i) \end{aligned}$$

$$= \prod_{i=1}^n \left\{ y_i f(\beta_0 + \beta_1 x_i) + (1-y_i) (1-f(\beta_0 + \beta_1 x_i)) \right\}$$



Likelihood of data given model

$$\ln ab = \ln a + \ln b$$

$$\ln \prod_{i=1}^n a_i = \sum_{i=1}^n \ln a_i$$

Log likelihood:

$$\ln P(Y=y | X=\underline{x}) =$$

$$\sum_{i=1}^n \ln \left\{ y_i f(\beta_0 + \beta_1 x_i) + (1 - y_i) (1 - f(\beta_0 + \beta_1 x_i)) \right\}$$

UPDATE

Log your COVID vaccination in the app

39,014

Total number of new daily cases
across the UK

[VIEW THE LATEST DATA >](#)

You can help fight COVID-19
by aiding research

<https://covid.joinzoe.com/>