

Inf2 - Foundations of Data Science: Multiple logistic regression for explanation and prediction



THE UNIVERSITY *of* EDINBURGH
informatics

F O U N D A T I O N S
O F
D A T A
S C I E N C E

Announcements

- Week 4 workshop - we'll look at the paper that we'll be refer to in the exam
- Uses concepts from today's lecture!
- Solutions for Week 3 Workshop now available
- Solutions for this Week 4 Workshop will be available later in the week
- Badges on order!

Where we're at in the Maximum Likelihood Principle and Regression

Week 4: Logistic regression

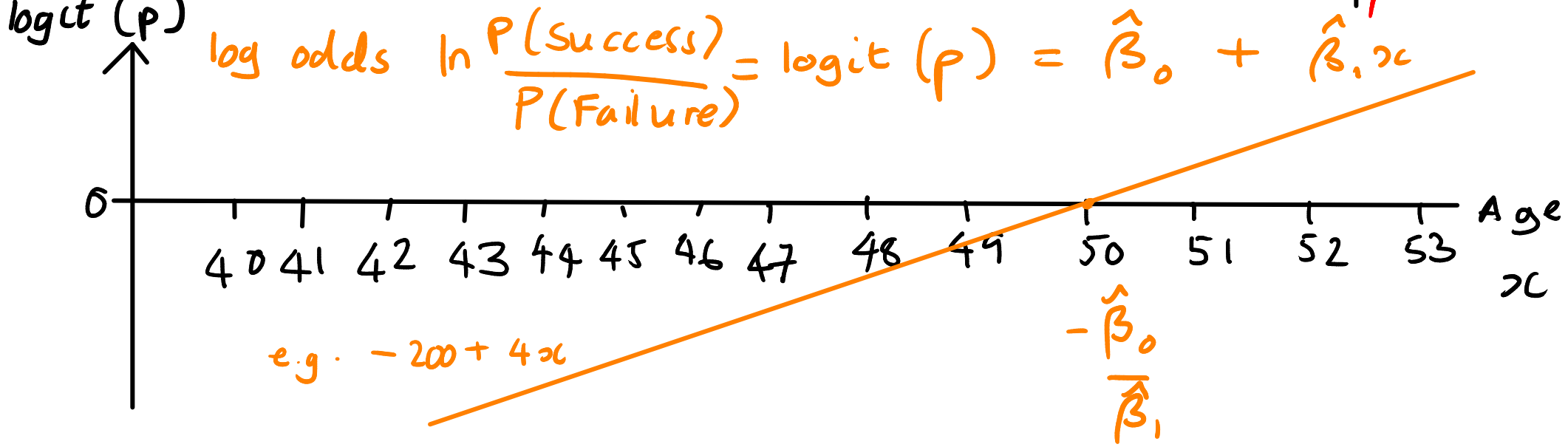
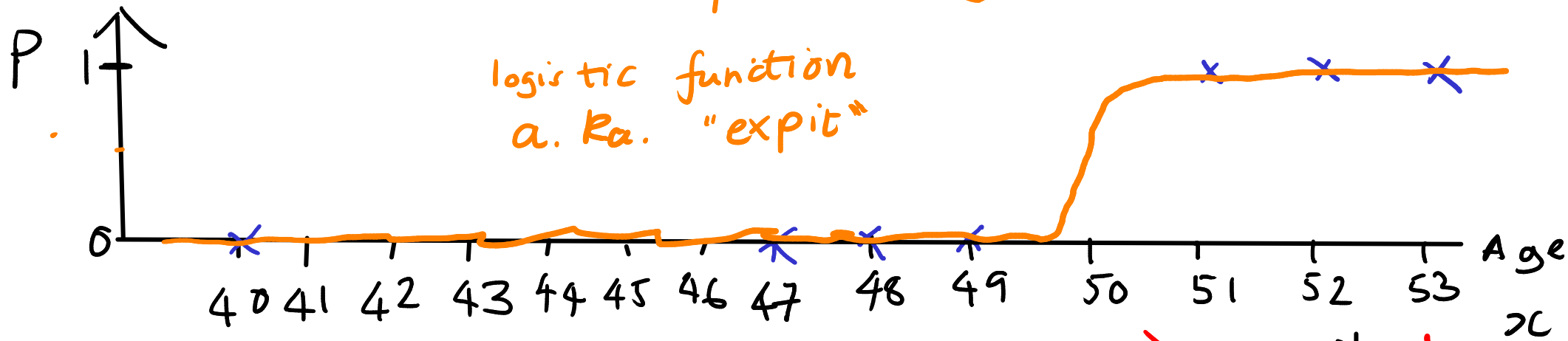
Week 5: The maximum likelihood principle, and how we can use it to derive linear, logistic and other types of regression

Today

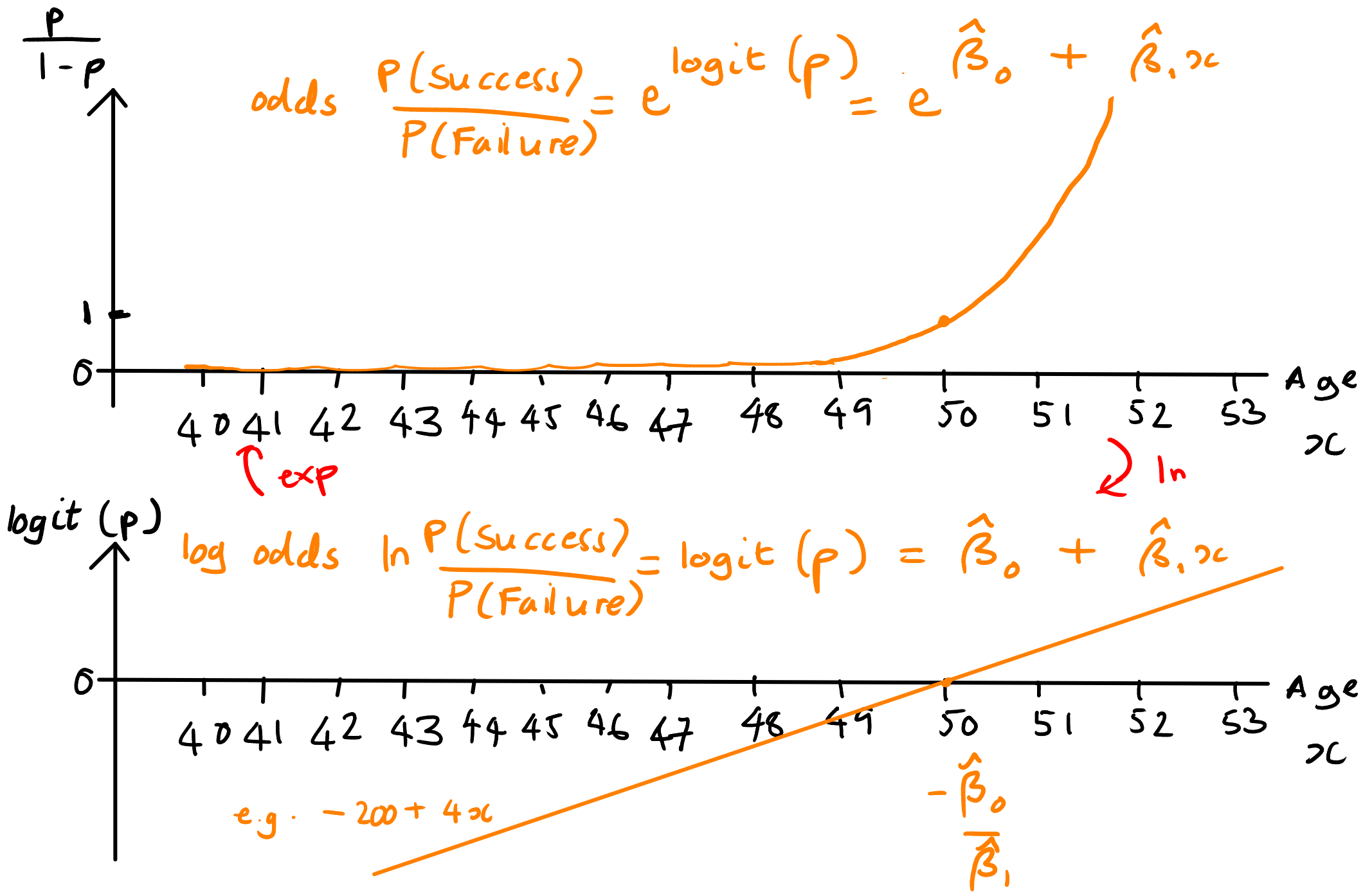
- Recap
- Multiple Logistic Regression
- Confidence intervals on coefficients
- Machine learning: Logistic Regression as a classifier
- Ethics of logistic regression

Probability and log odds views of logistic regression

Approval $P(\text{success}) = P = f(\hat{\beta}_0 + \hat{\beta}_1 x)$



Odds and log odds views of logistic regression



Binary variables: odds and odds ratios

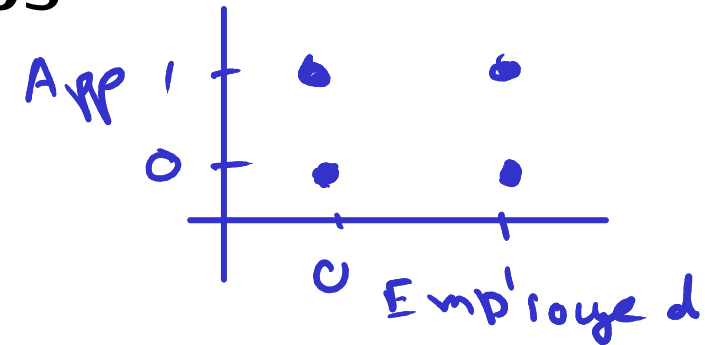
$$P(Y = y \mid X = x)$$

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

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

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

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



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

Effect size
609 %

Inf2 - Foundations of Data Science: Multiple logistic regression



THE UNIVERSITY *of* EDINBURGH
informatics

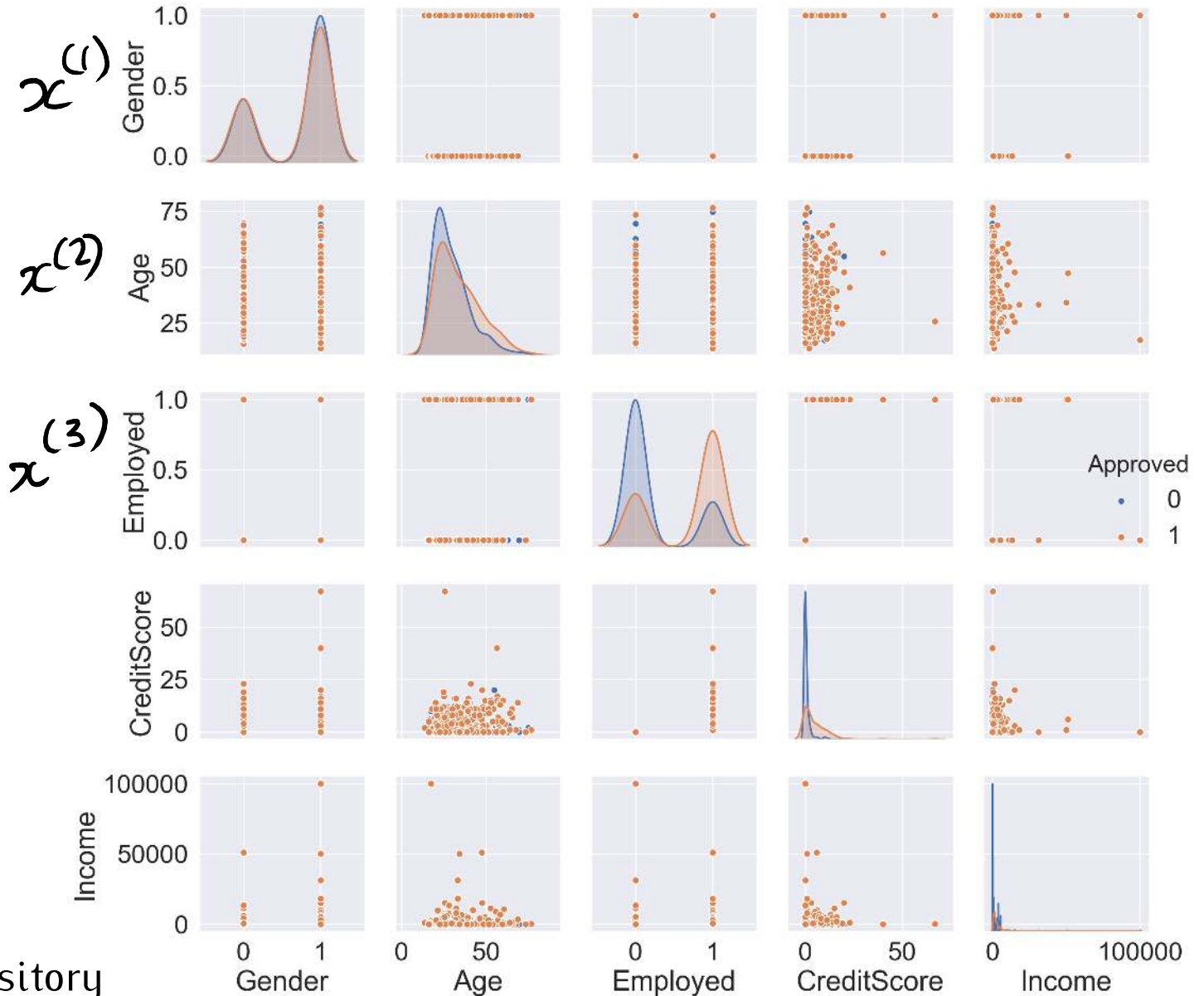
F **O** **U** **N** **D** **A** **T** **I** **O** **N** **S**
O **F**
D **A** **T** **A**
S **C** **I** **E** **N** **C** **E**

Supervised classification

Binary (or dichotomous) response variable:
Credit

Approved

Not approved



UCI Machine Learning Repository

<https://archive.ics.uci.edu/dataset/27/credit+approval>

Principle of multiple logistic regression

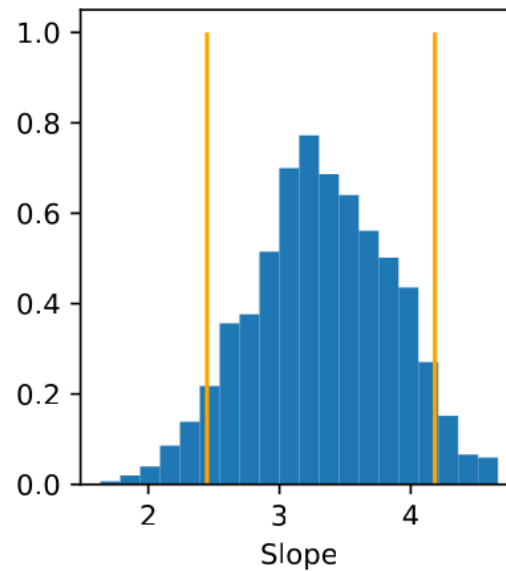
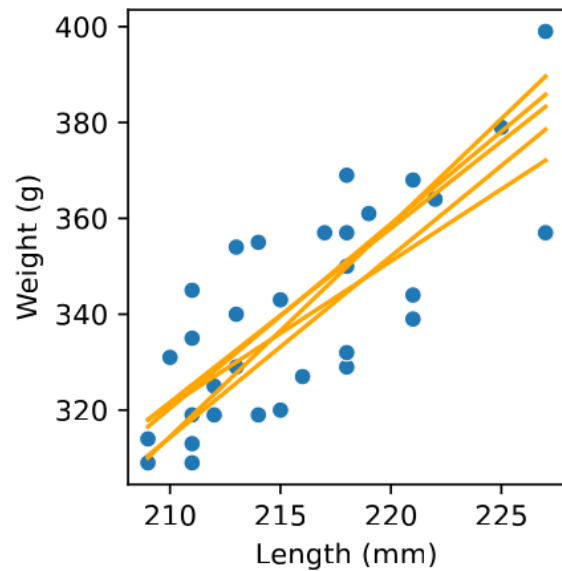
Predictor variables $x^{(1)}$: Age
 $x^{(2)}$: Employment

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

Multiple logistic regression applied to the credit example

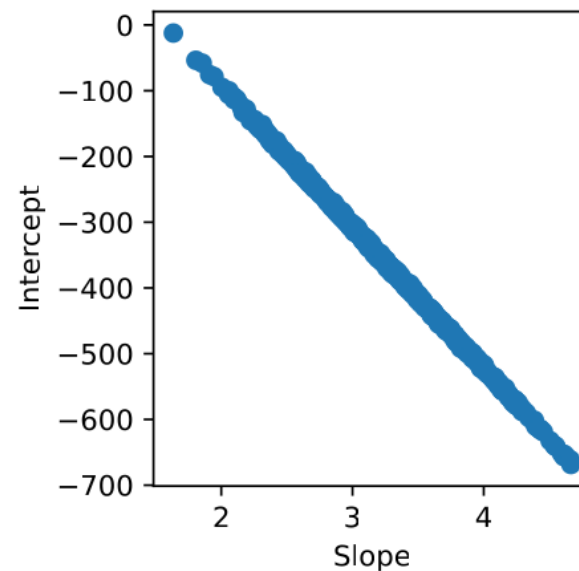
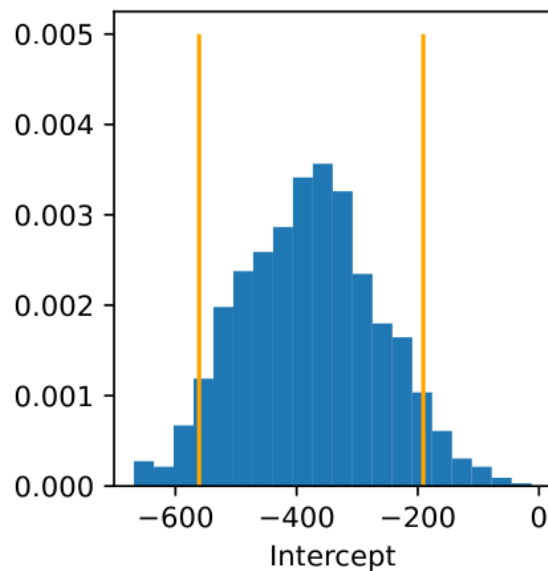
	Variable	Coefficient	Odds or OR
$\hat{\beta}_0$	Intercept	-1.969	0.140
$\hat{\beta}_1$	Age	0.029	1.030
$\hat{\beta}_2$	Employed	1.881	6.562

Bootstrap confidence intervals for regression coefficients



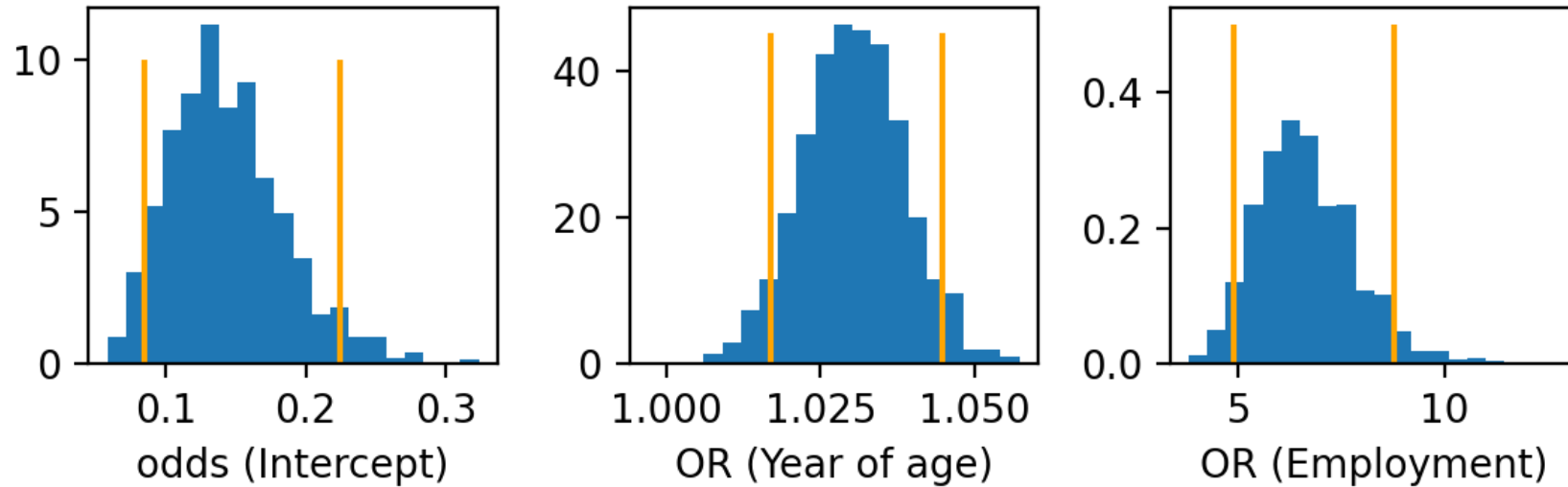
Demo

Code this for Logistic Regression in the lab!



Bootstrap confidence intervals

$B = 1000$ $n = 653$



Does age affect credit approval?

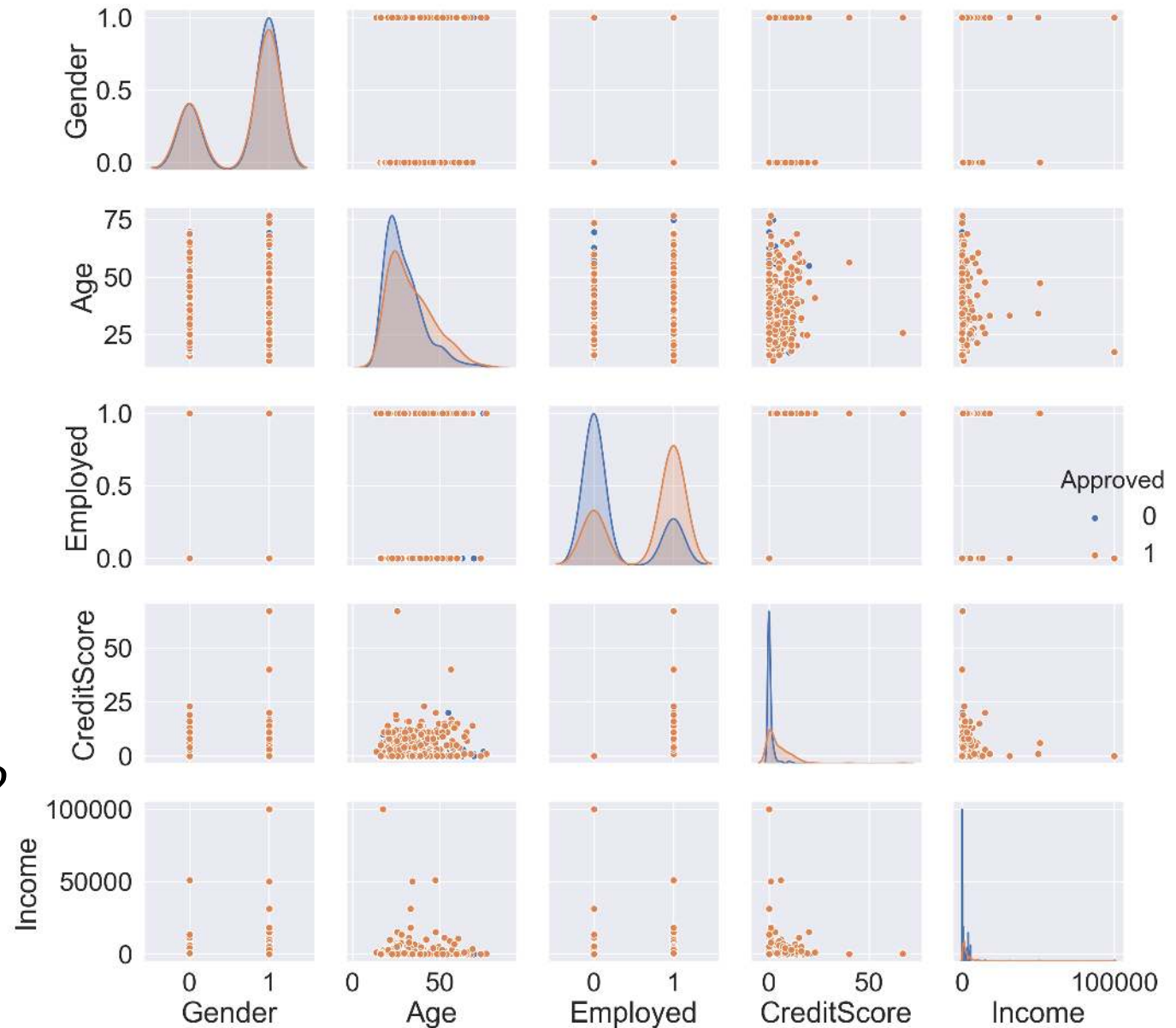
Discussion question

Our analysis so far shows that age and credit approval are related.

So all other things being equal, a 20 year old is less likely to have credit approved than a 50 year old.

Do we believe this yet?

What further analysis should we do?



Explanation - "controlling for", "adjusting for"

This week's lab

Multiple logistic regression on fuller set of variables

Using Logistic Regression as a Machine Learning algorithm

Inf2 - Foundations of Data Science: The logistic regression classifier



THE UNIVERSITY *of* EDINBURGH
informatics

F **O** **U** **N** **D** **A** **T** **I** **O** **N** **S**
O **F**
D **A** **T** **A**
S **C** **I** **E** **N** **C** **E**

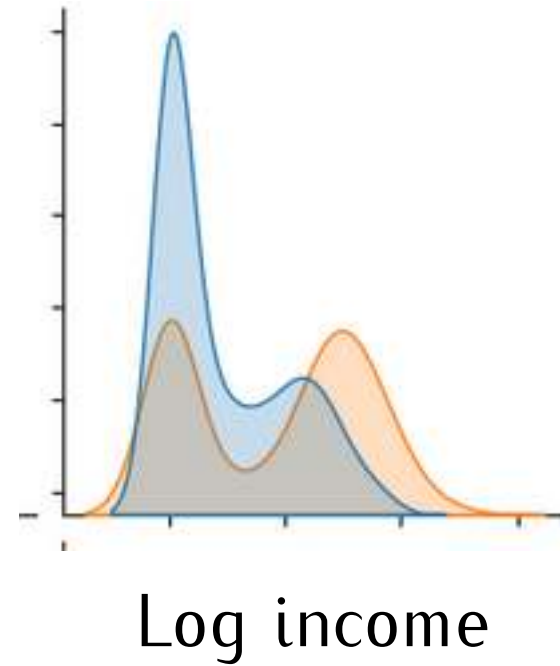
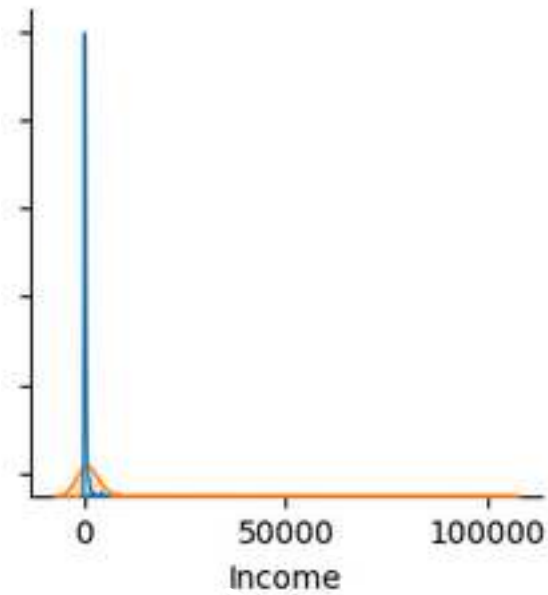
Converting logistic regression to a classifier

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

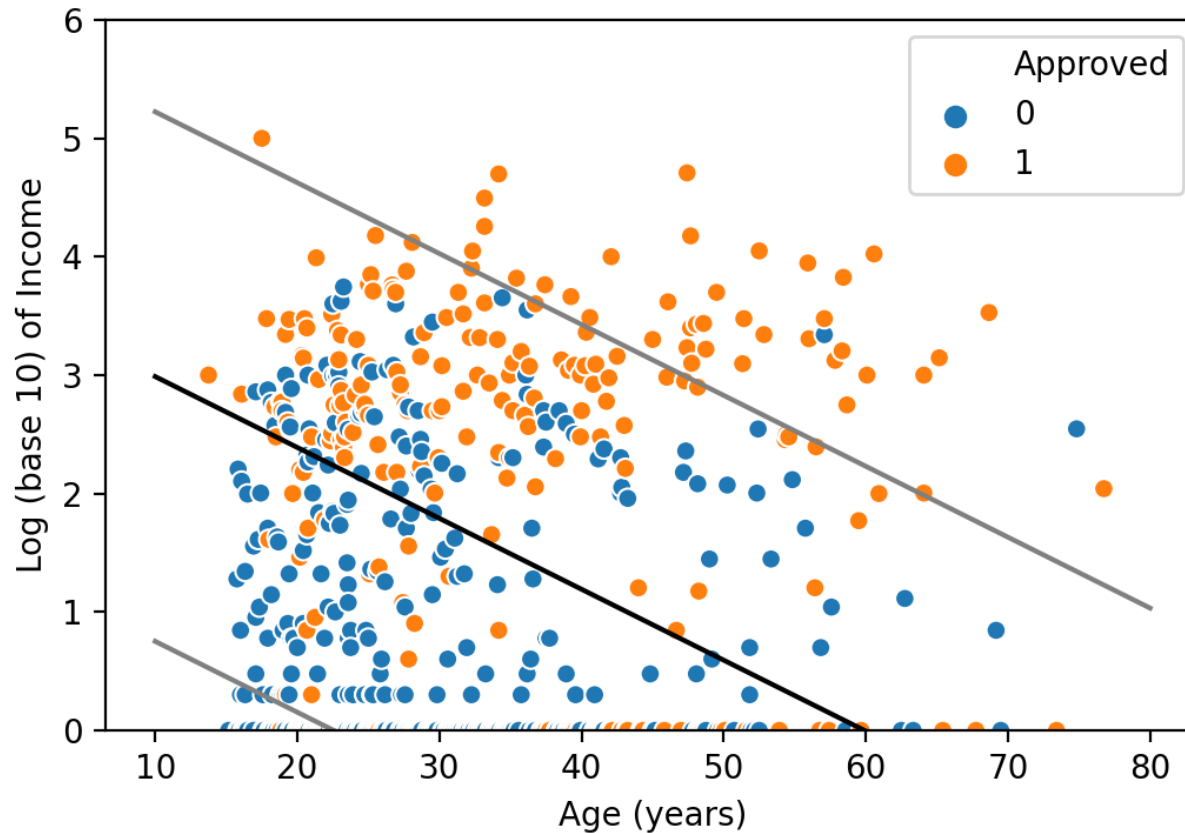
$$\hat{\beta}_0 + \hat{\beta}_1 x^{(1)} + \hat{\beta}_2 x^{(2)} + \dots \geq c \Rightarrow \hat{y} = 1$$
$$\hat{\beta}_0 + \hat{\beta}_1 x^{(1)} + \hat{\beta}_2 x^{(2)} + \dots < c \Rightarrow \hat{y} = 0$$

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

Machine learning trick: make marginal distributions more normal



Decision boundary



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

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

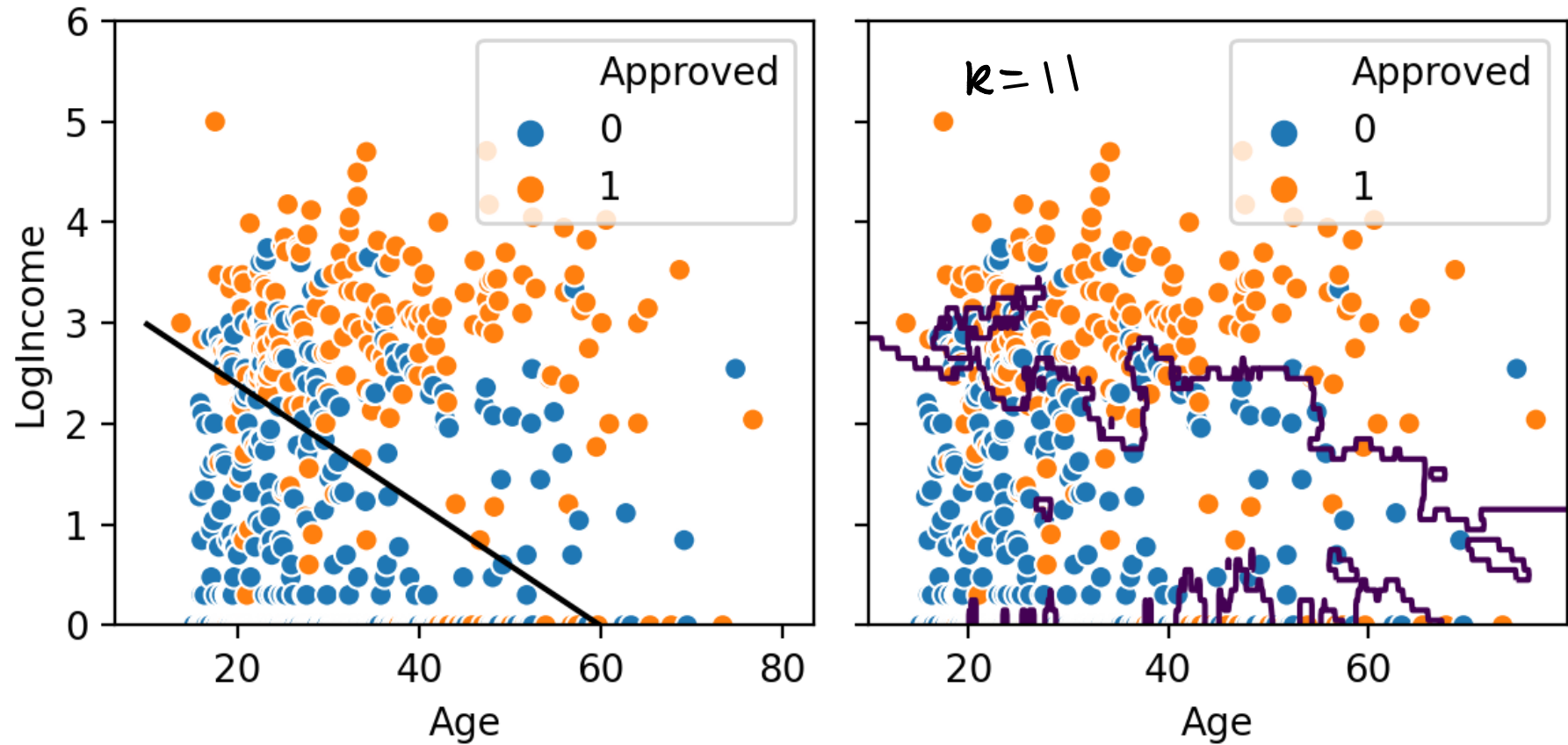
Ethics: logistic regression can be transparent

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

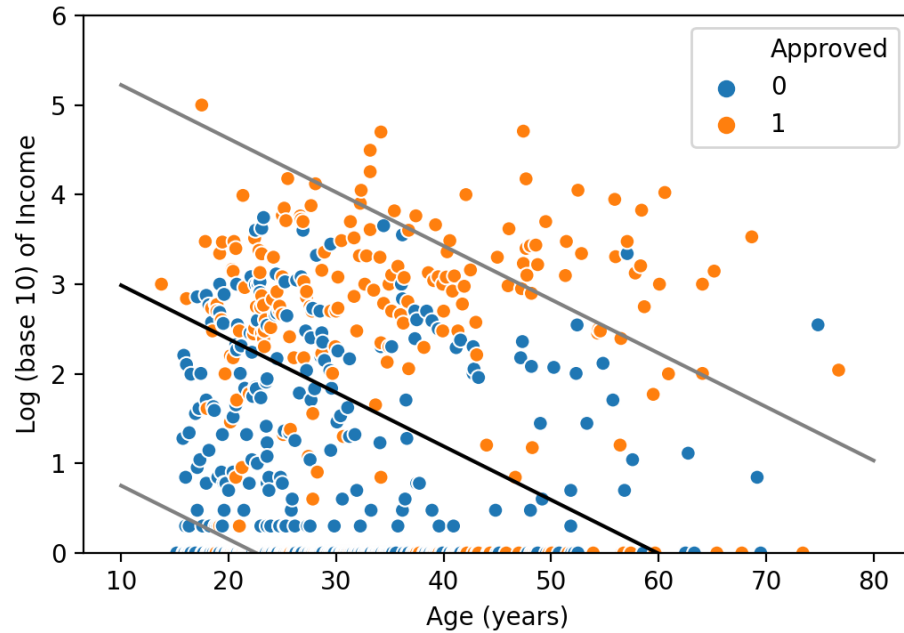
Standardised input variables

Receiver-operator characteristic

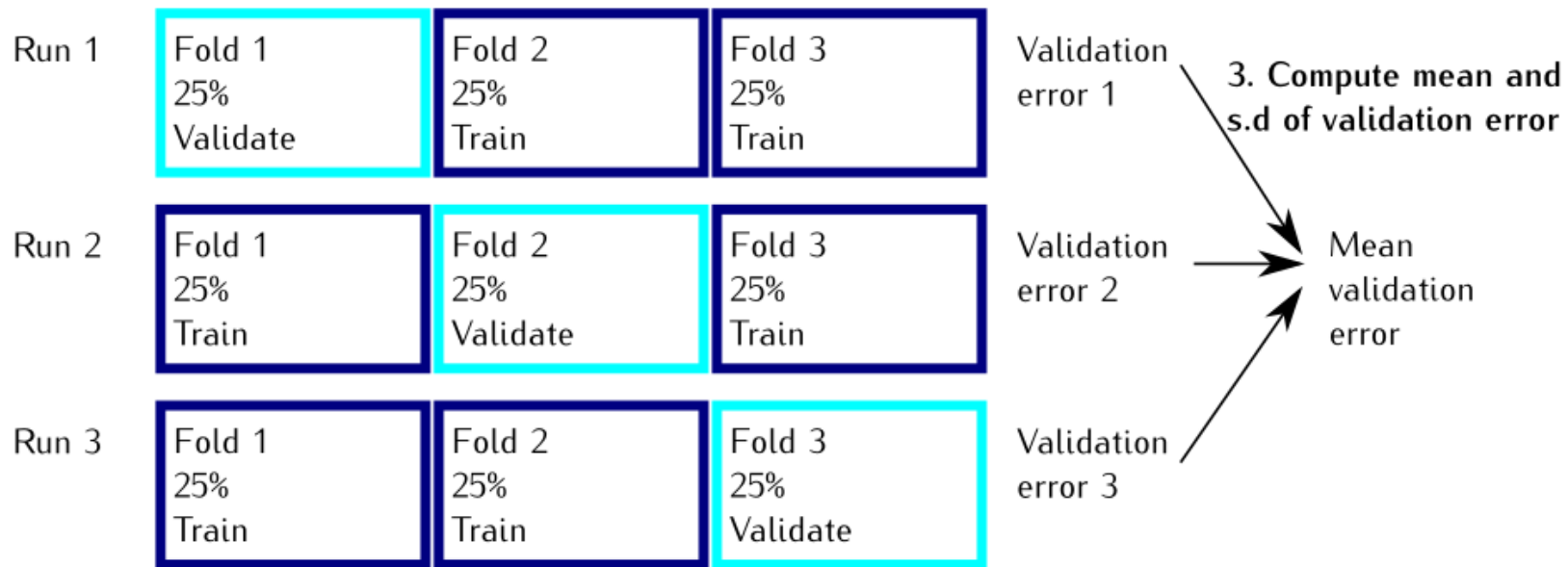
Sensitivity

Selectivity/
Specificity

False positive
rate



Cross validation for predicting metrics



c.f. Chapter 12 of the lecture notes

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
- Practical and ethical pros and cons of logistic regression versus other methods