

# Analysis of high throughput antibody data for better understanding of immunogenetics and epidemiology of malaria

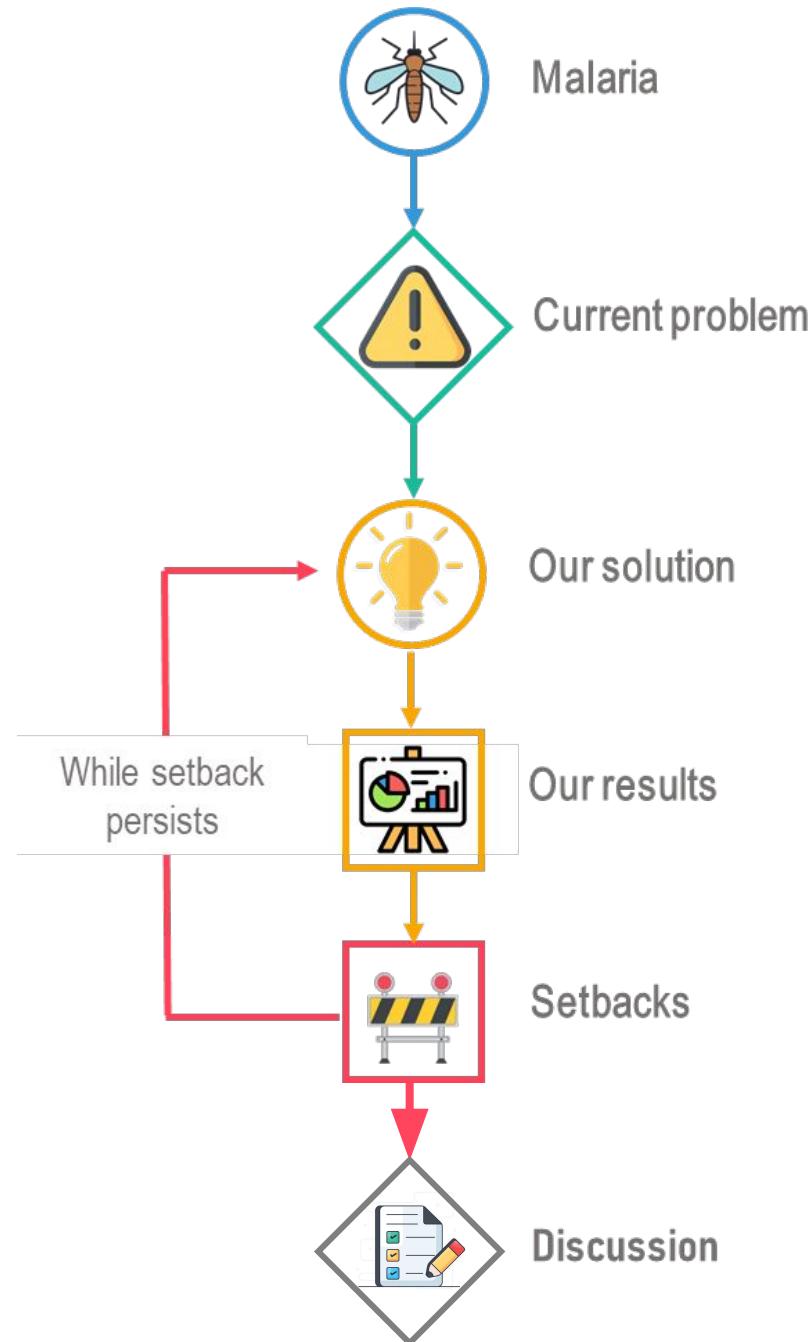
Ref FCT : SFRH/BD/147629/2019

Supervisor: Dr Nuno Sepúlveda  
Co-supervisor: Dr Clara Cordeiro

André  
Fonseca  
06/12/2021

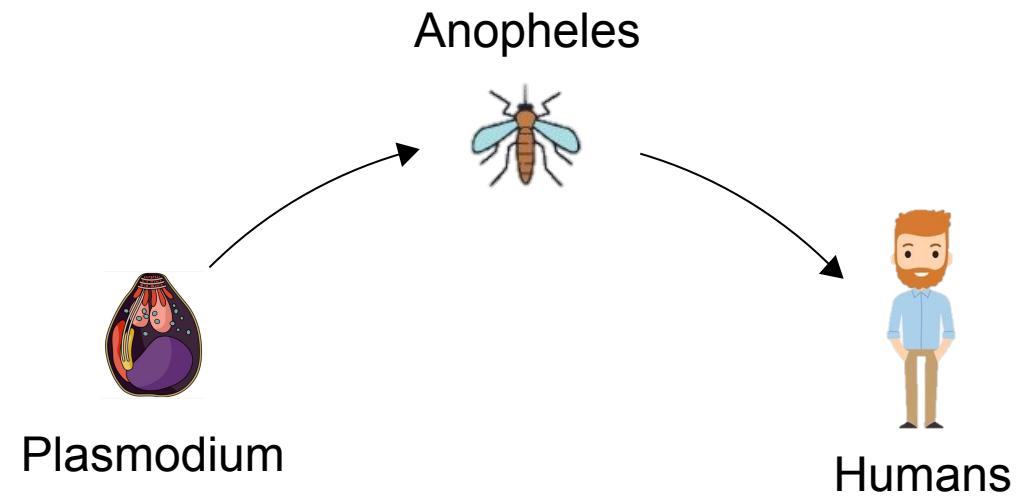


# Agenda



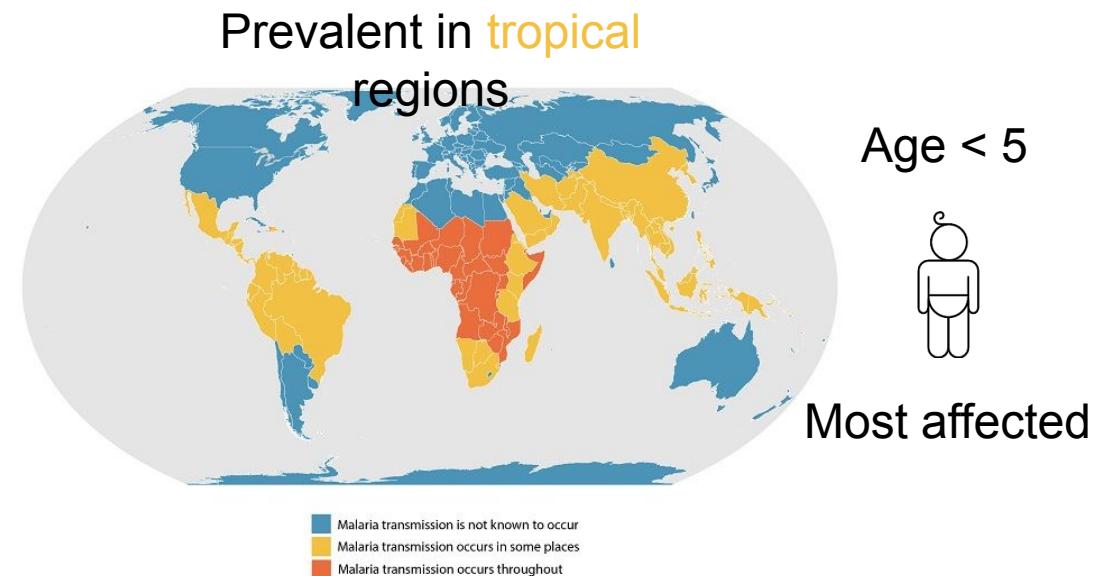
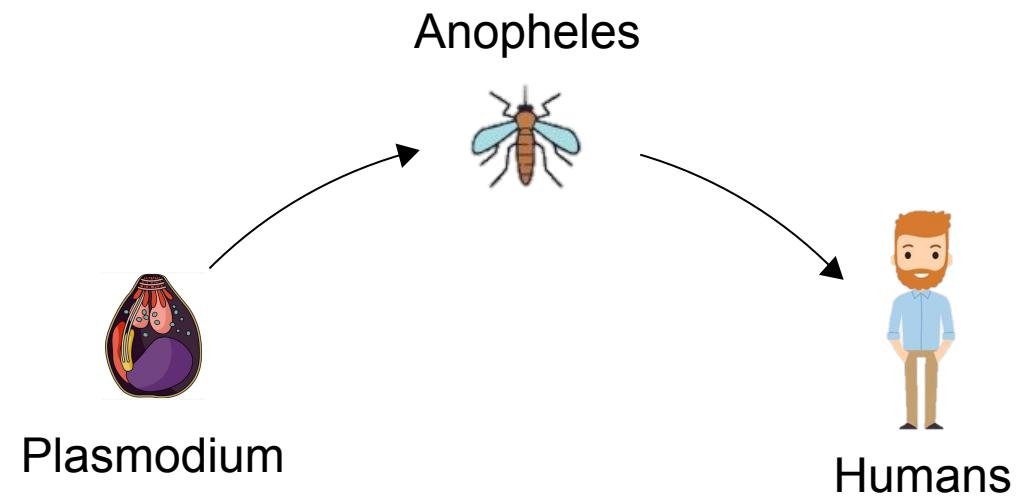


# Malaria



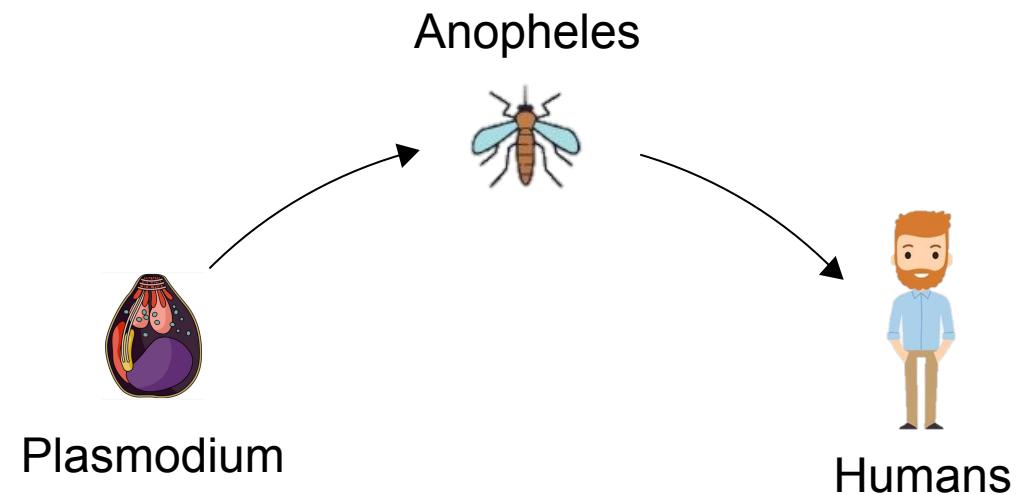


# Malaria

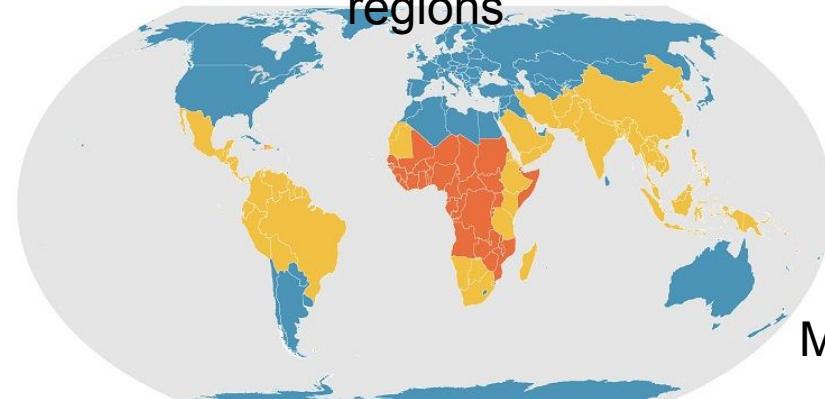




# Malaria



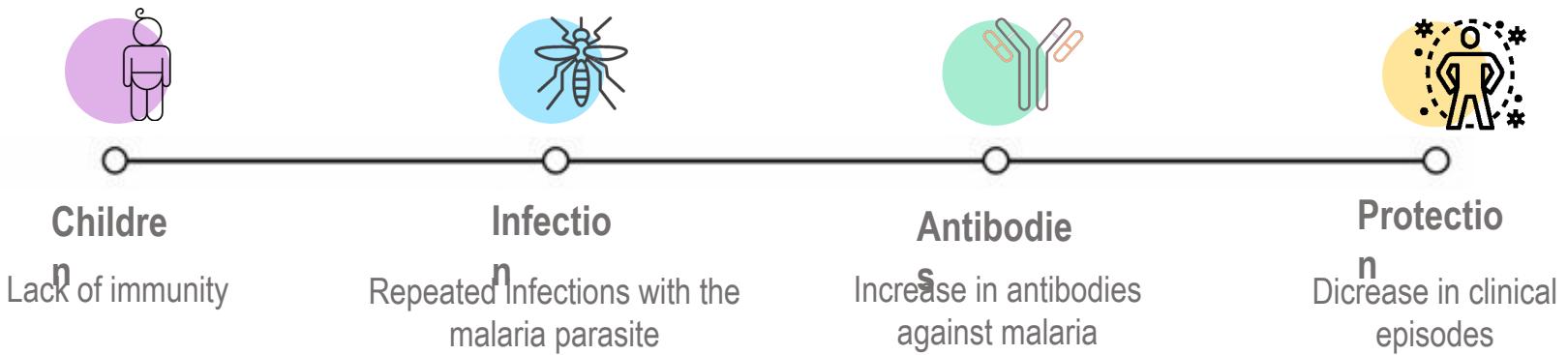
Prevalent in **tropical**  
regions



Age < 5

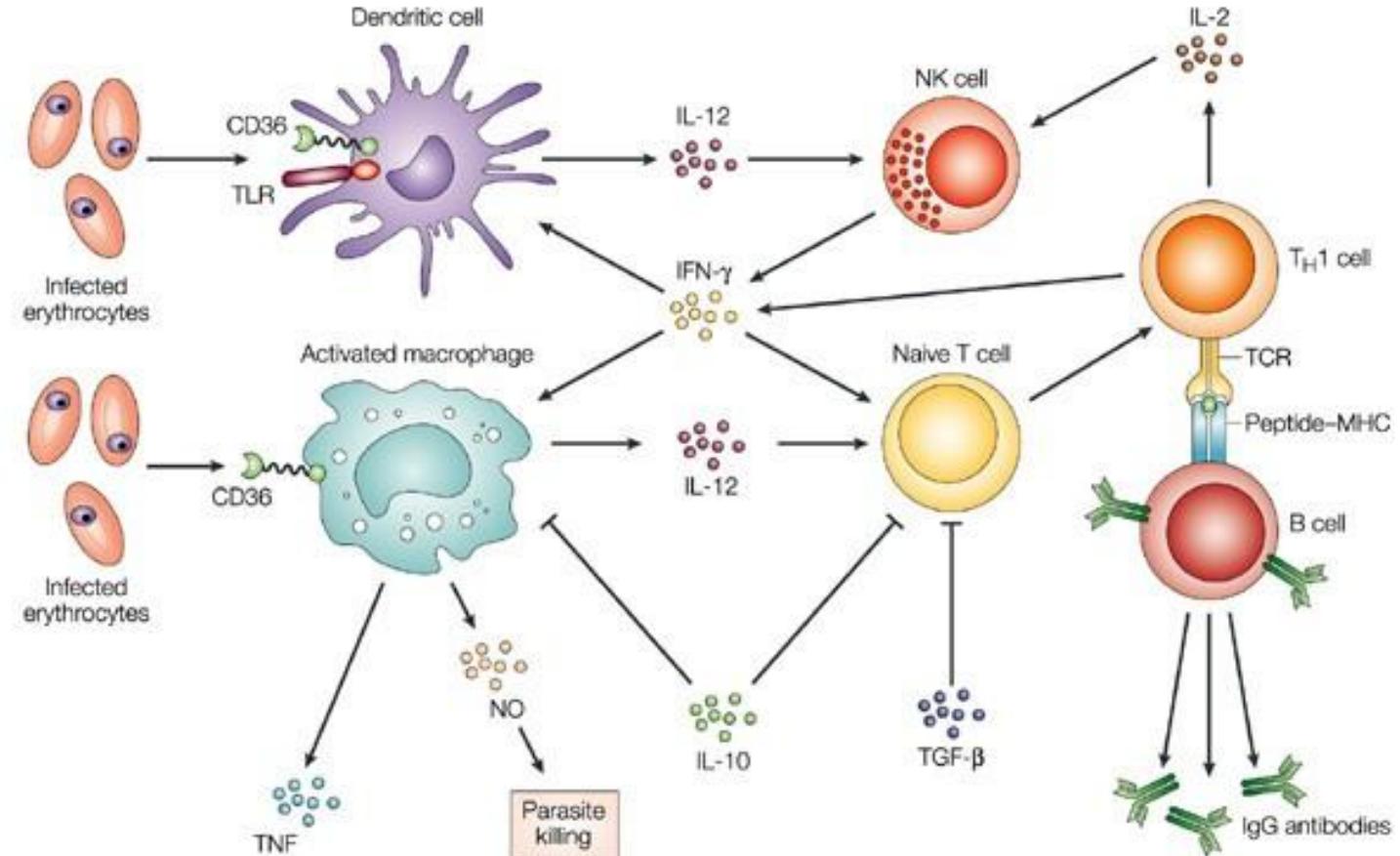


Most affected





# Antibody protection





## Problem



Unable to identify antibodies for  
**vaccine** development against  
malaria



## Problem



Unable to identify antibodies for  
**vaccine** development against  
malaria



### RTS,S/AS01

Efficacy on children (36%)



## Problem



Unable to identify antibodies for **vaccine** development against malaria



RTS,S/AS01

Efficacy on children (36%)



Difficulties in identifying vaccine candidate



Biological factors



Technical factors



## Problem



Unable to identify antibodies for **vaccine** development against malaria



RTS,S/AS01

Efficacy on children (36%)



Difficulties in identifying vaccine candidate



**Biological factors**

- Evolution of the parasite
- Complex life cycle
- Individual immune response



**Technical factors**



## Problem



Unable to identify antibodies for **vaccine** development against malaria



RTS,S/AS01

Efficacy on children (36%)



Difficulties in identifying vaccine candidate



**Biological factors**

- Evolution of the parasite
- Complex life cycle
- Individual immune response



**Technical factors**

- Experimental design



## Problem



Unable to identify antibodies for **vaccine** development against malaria



RTS,S/AS01

Efficacy on children (36%)



Difficulties in identifying vaccine candidate



**Biological factors**

- Evolution of the parasite
- Complex life cycle
- Individual immune response



**Technical factors**

- Experimental design
- Sub-optimal analysis



**Inconsistent results**



failure of the statistical assumptions (**normality**)



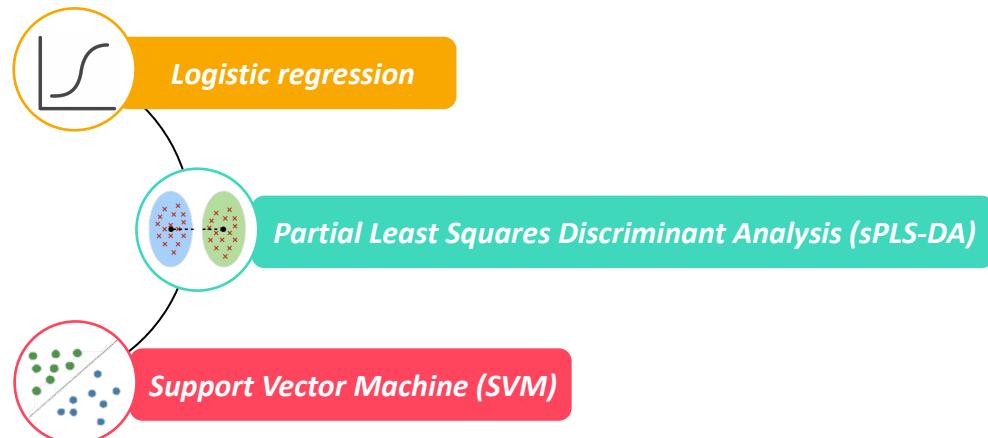
# Problem

Research

## Immune Signature Against *Plasmodium falciparum* Antigens Predicts Clinical Immunity in Distinct Malaria Endemic Communities

### Authors

Carla Proletti, Lutz Krause, Angela Trieu, Daniel Dodoo, Ben Gyan, Kwadwo A. Koram, William O. Rogers, Thomas L. Richie, Peter D. Crompton, Phillip L. Felgner, and Denise L. Doolan





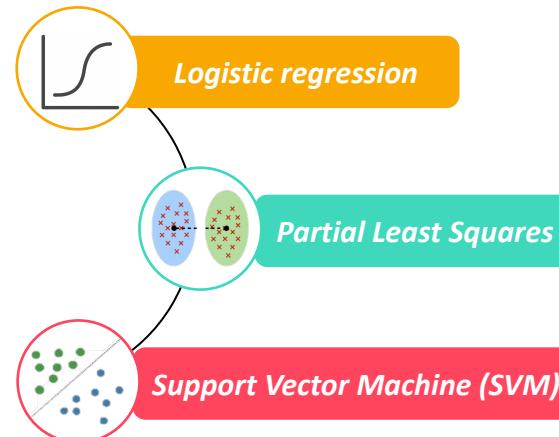
# Problem

Research

## Immune Signature Against *Plasmodium falciparum* Antigens Predicts Clinical Immunity in Distinct Malaria Endemic Communities

### Authors

Carla Proletti, Lutz Krause, Angela Trieu, Daniel Dodoo, Ben Gyan, Kwadwo A. Koram, William O. Rogers, Thomas L. Richie, Peter D. Crompton, Phillip L. Felgner, and Denise L. Doolan



PLOS COMPUTATIONAL BIOLOGY

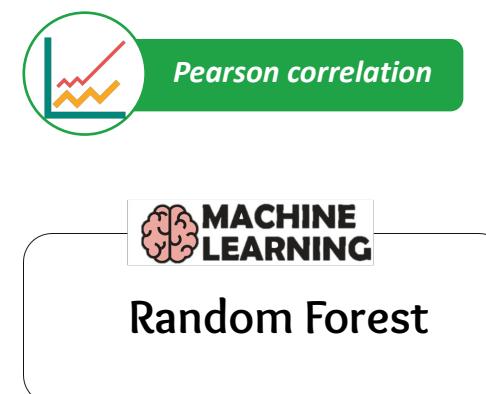
### RESEARCH ARTICLE

## Identification of immune signatures predictive of clinical protection from malaria

John Joseph Valletta, Mario Recker\*

Centre for Mathematics and the Environment, University of Exeter, Penryn Campus, Penryn, United Kingdom

\* [m.recker@exeter.ac.uk](mailto:m.recker@exeter.ac.uk)





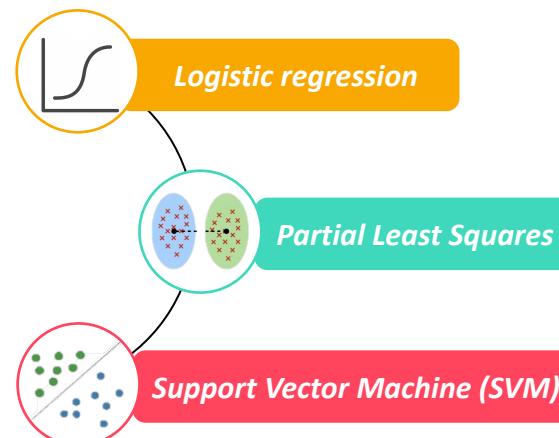
# Problem

Research

## Immune Signature Against *Plasmodium falciparum* Antigens Predicts Clinical Immunity in Distinct Malaria Endemic Communities

### Authors

Carla Proletti, Lutz Krause, Angela Trieu, Daniel Dodoo, Ben Gyan, Kwadwo A. Koram, William O. Rogers, Thomas L. Richie, Peter D. Crompton, Phillip L. Felgner, and Denise L. Doolan



PLOS COMPUTATIONAL BIOLOGY

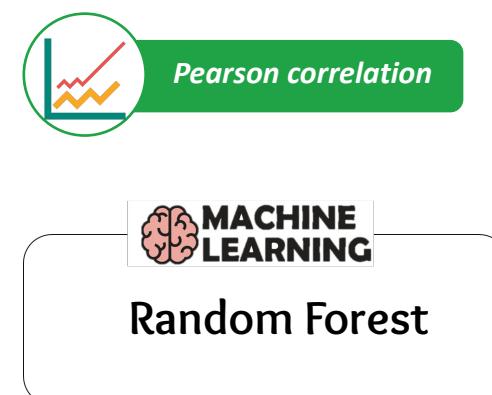
### RESEARCH ARTICLE

## Identification of immune signatures predictive of clinical protection from malaria

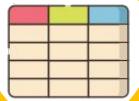
John Joseph Valletta, Mario Recker\*

Centre for Mathematics and the Environment, University of Exeter, Penryn Campus, Penryn, United Kingdom

\* [m.recker@exeter.ac.uk](mailto:m.recker@exeter.ac.uk)



Statistical pipelines to analyze immunological data consistently and reliably



## Our dataset



IgG antibodies against **36 Plasmodium Falciparum antigens**



Antibody quantity information for **121 Kenian children** (aged 1-10 years)



The **response variable** “Status” was binary (Protected vs. Susceptible)

Numeric    Binary

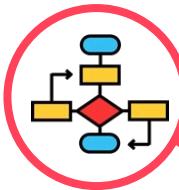
MSP	CSP	Pf34	Status
1			Protected
			Susceptible
			Susceptible
			Protected

36 antigens

121 children



# Solution



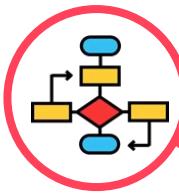
## Parametric pipeline

**Normal distribution and flexible mixture models for seropositivity determination**



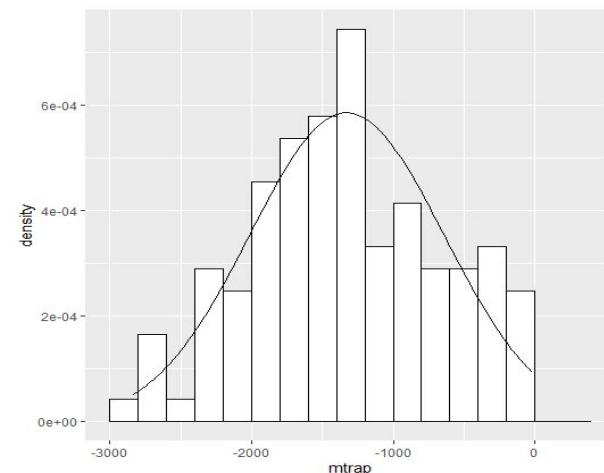
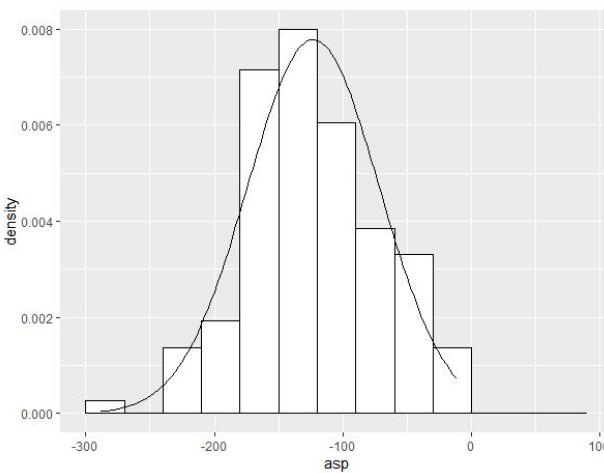


# Solution



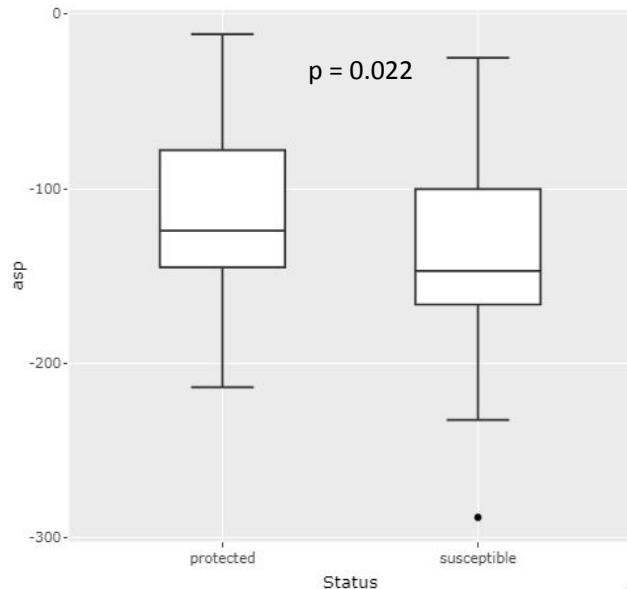
## Parametric pipeline

**Normal distribution and flexible mixture models for seropositivity determination**

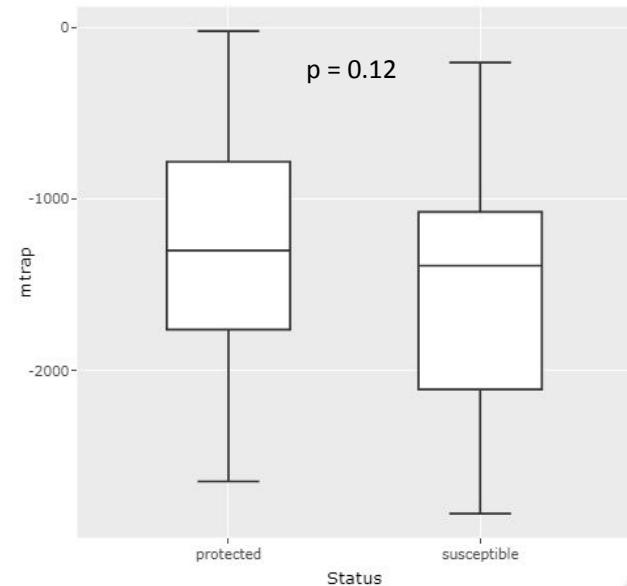




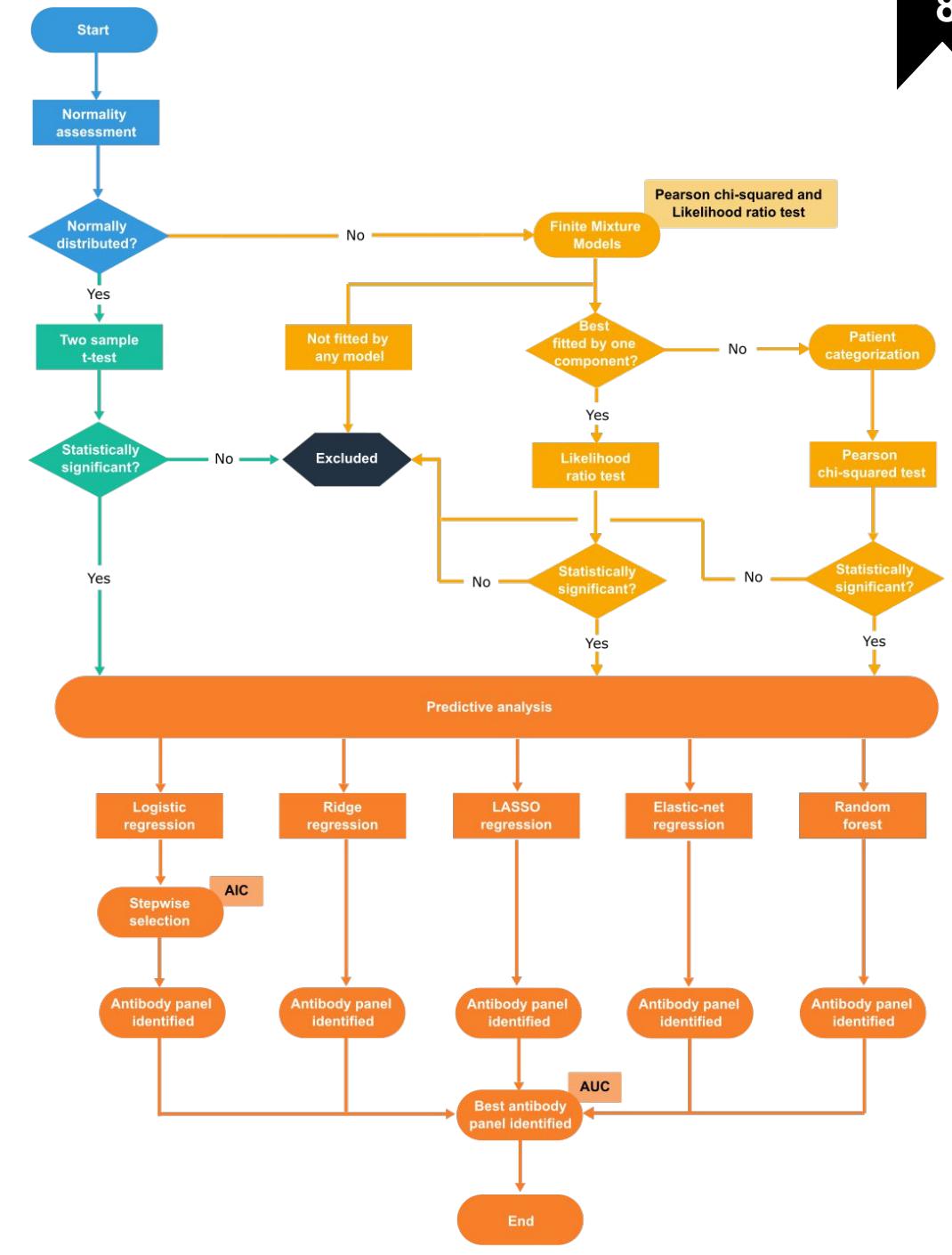
# Solution



Predictive analysis  
P-value < 0.05

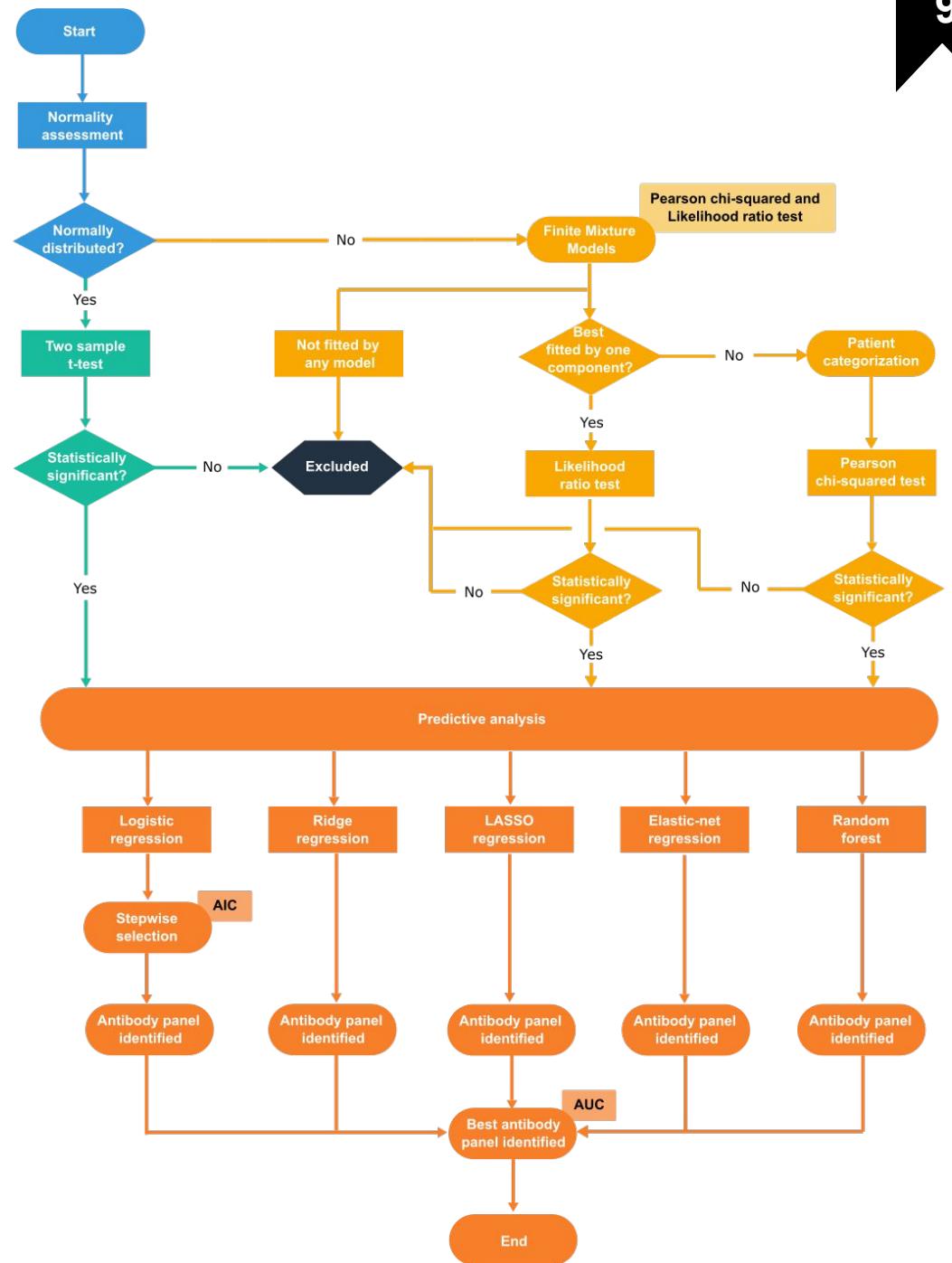


Excluded  
P-value > 0.05



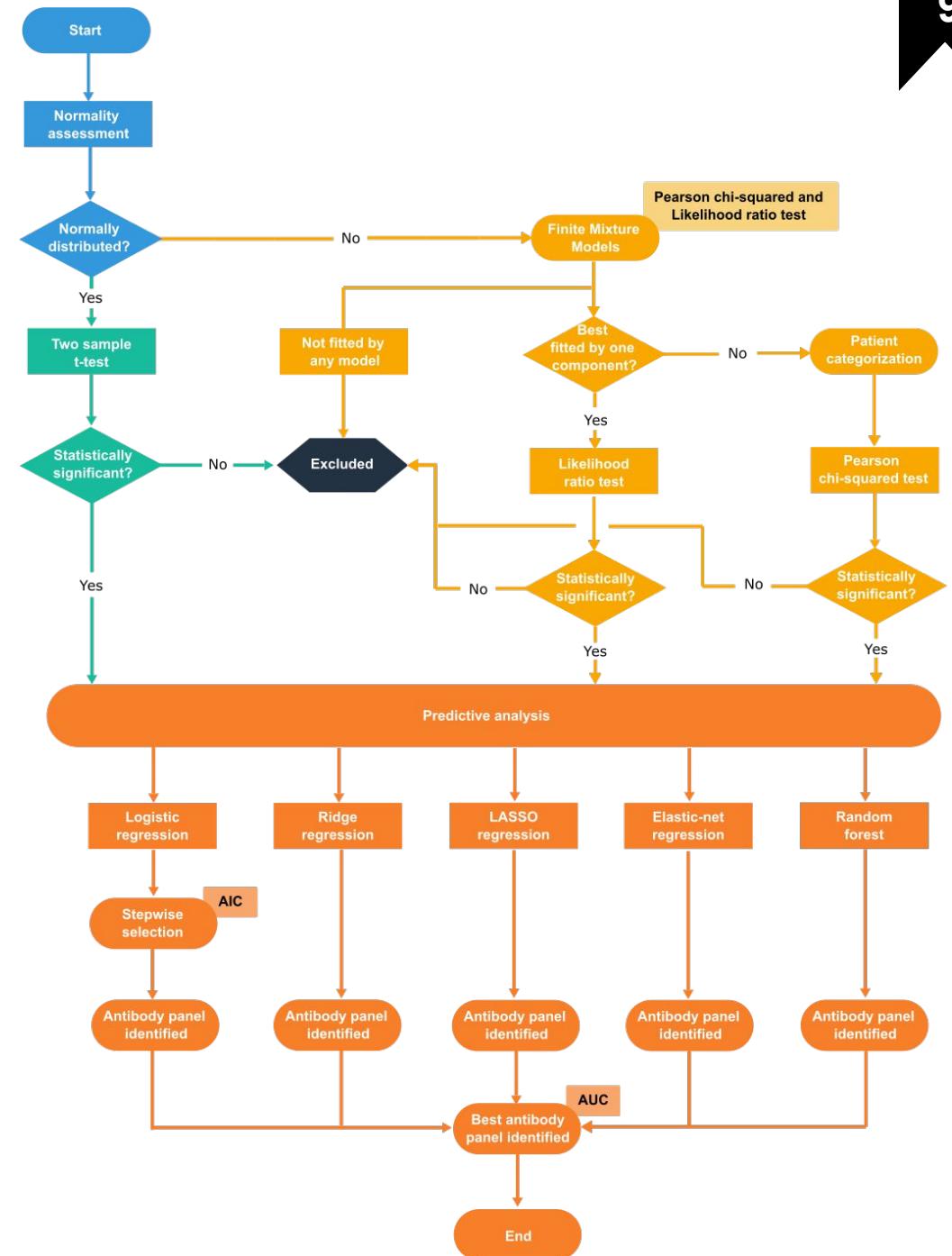
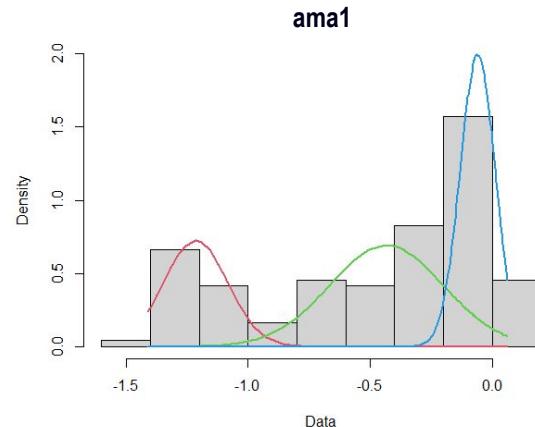
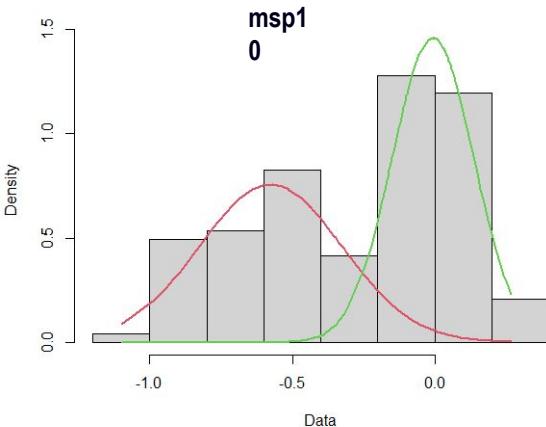


# Solution



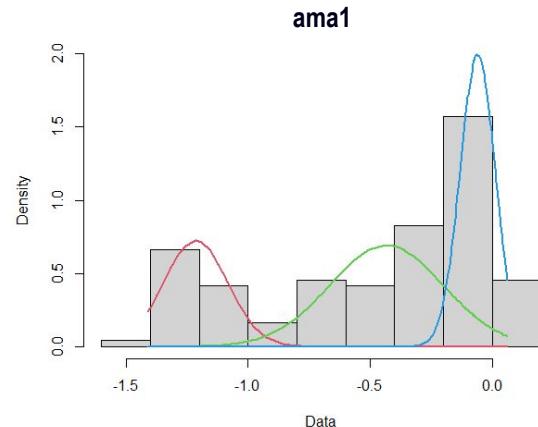
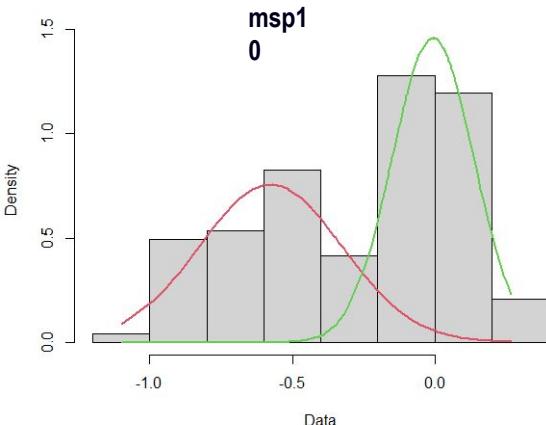


# Solution





## Solution

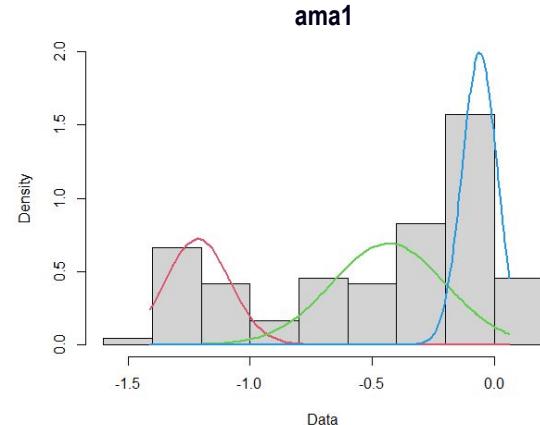
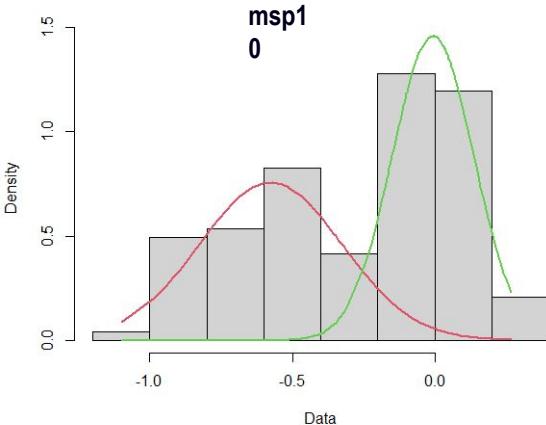


Mixture models are often used in seroepidemiologic studies

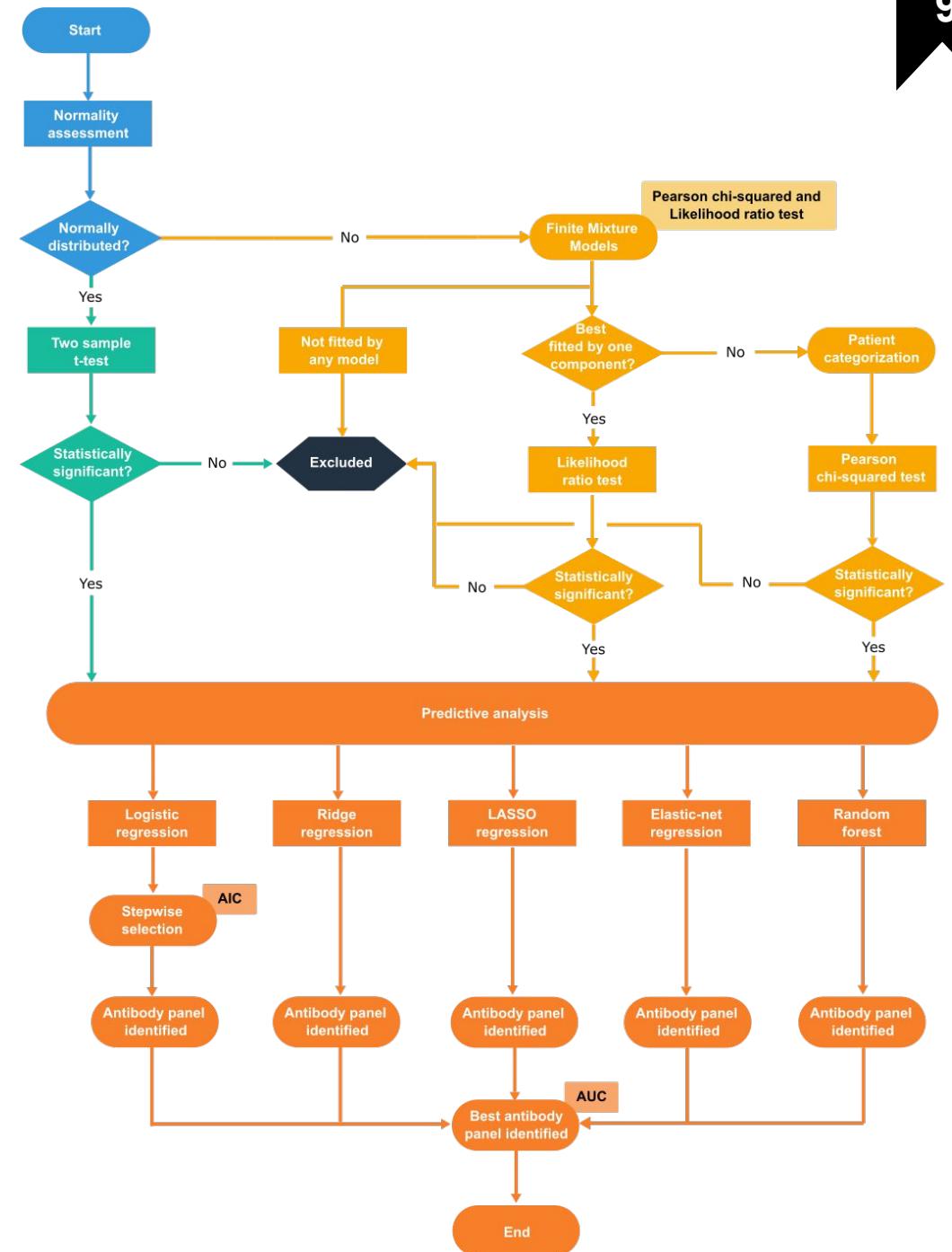
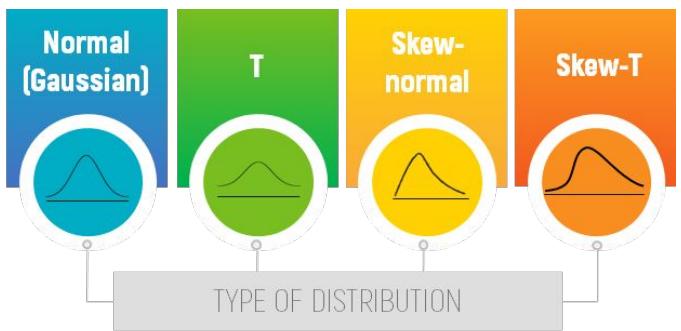




## Solution

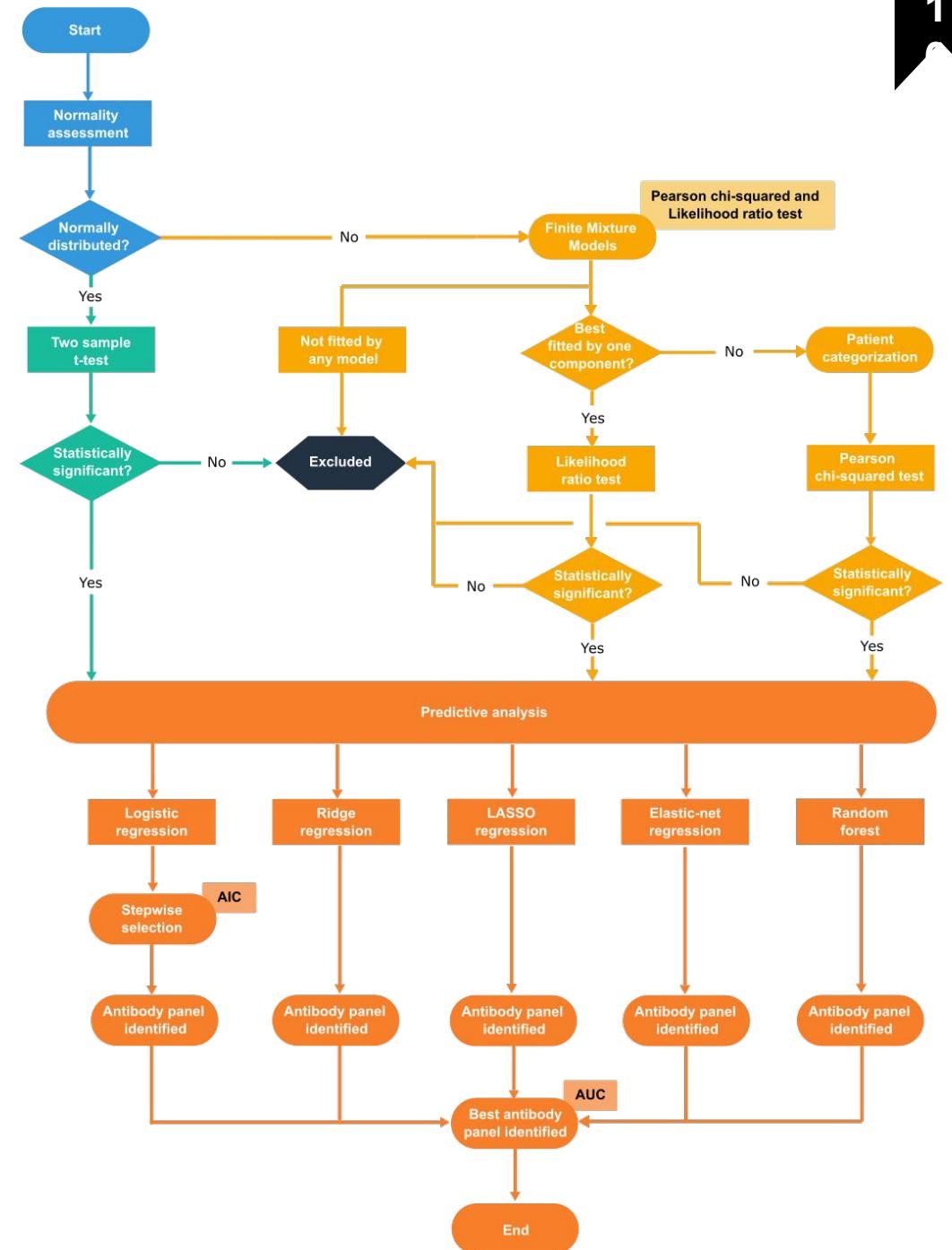
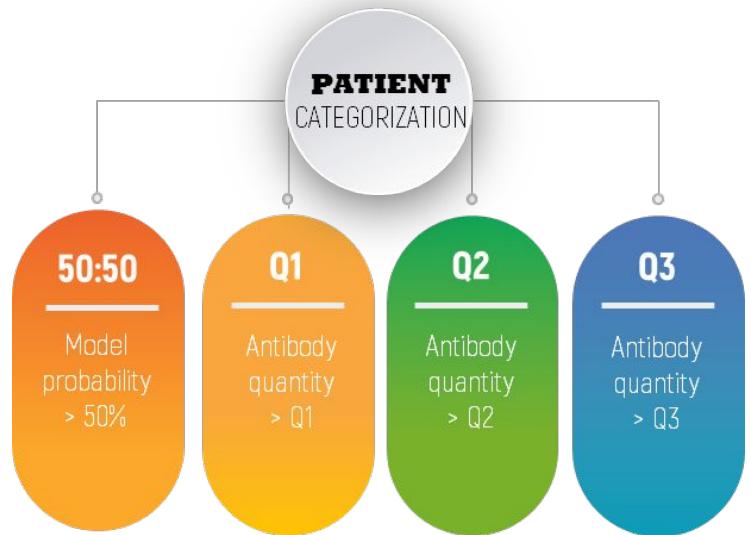


**Mixture models** are often used in seroepidemiologic studies



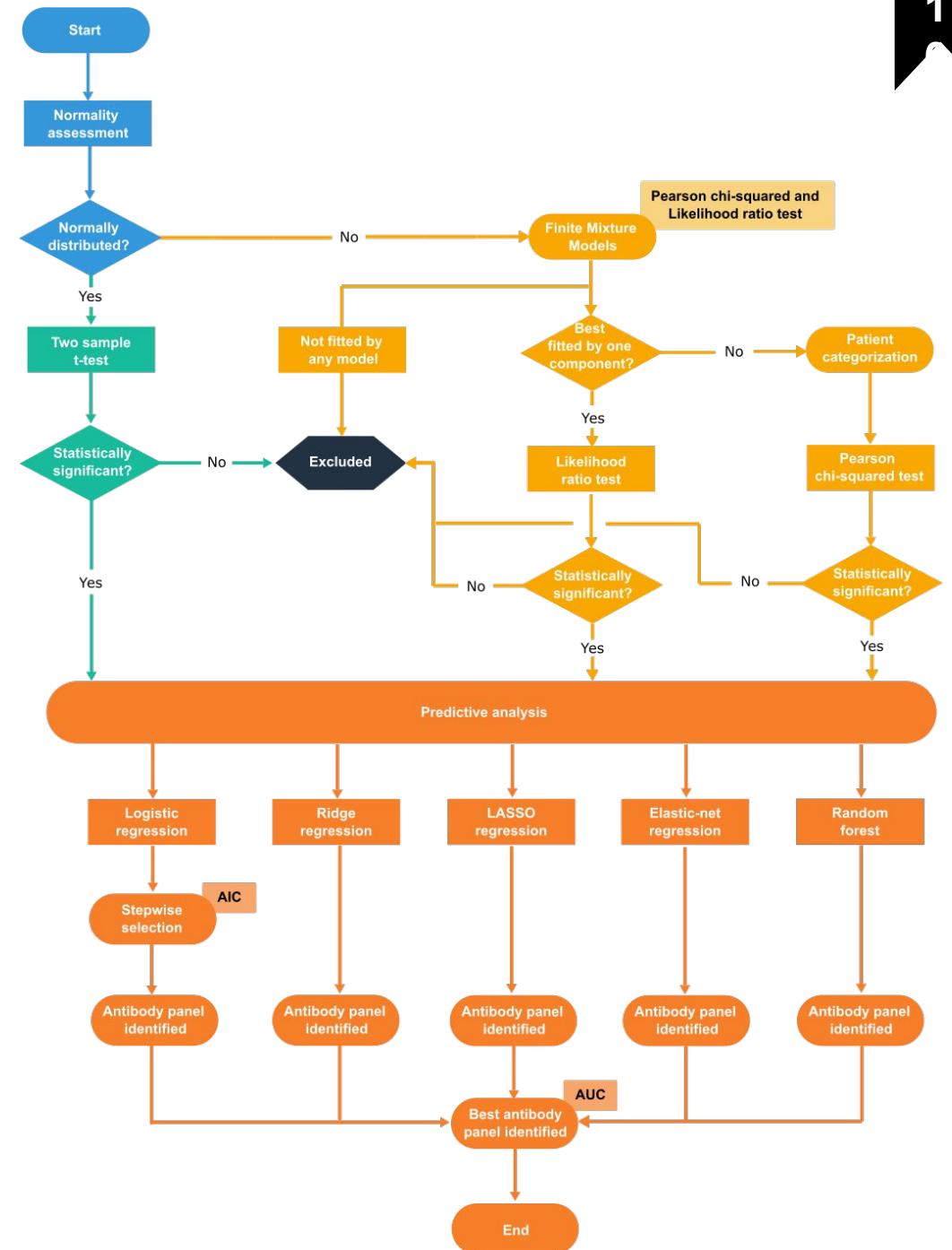
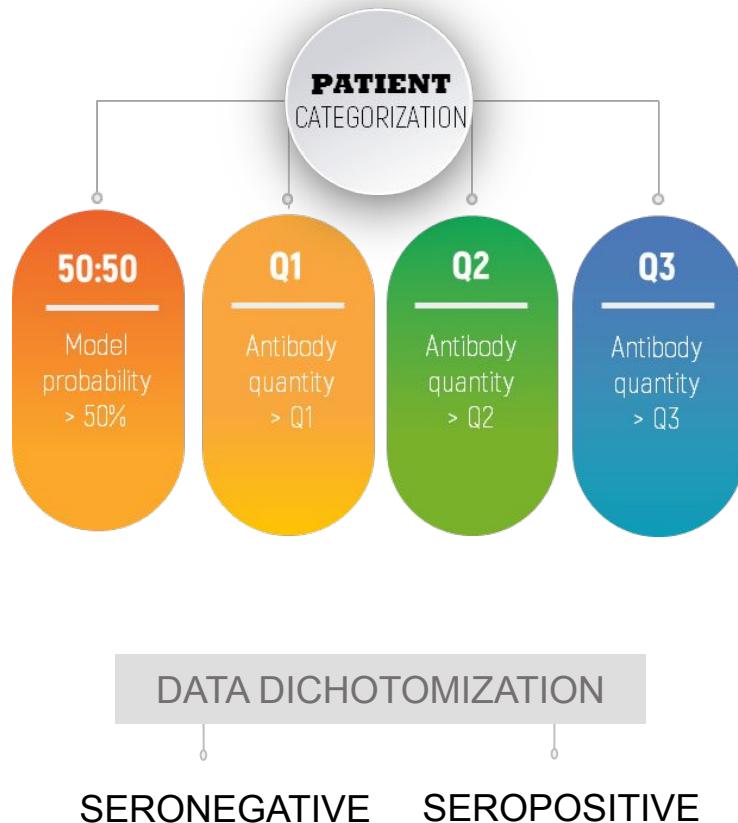


# Solution



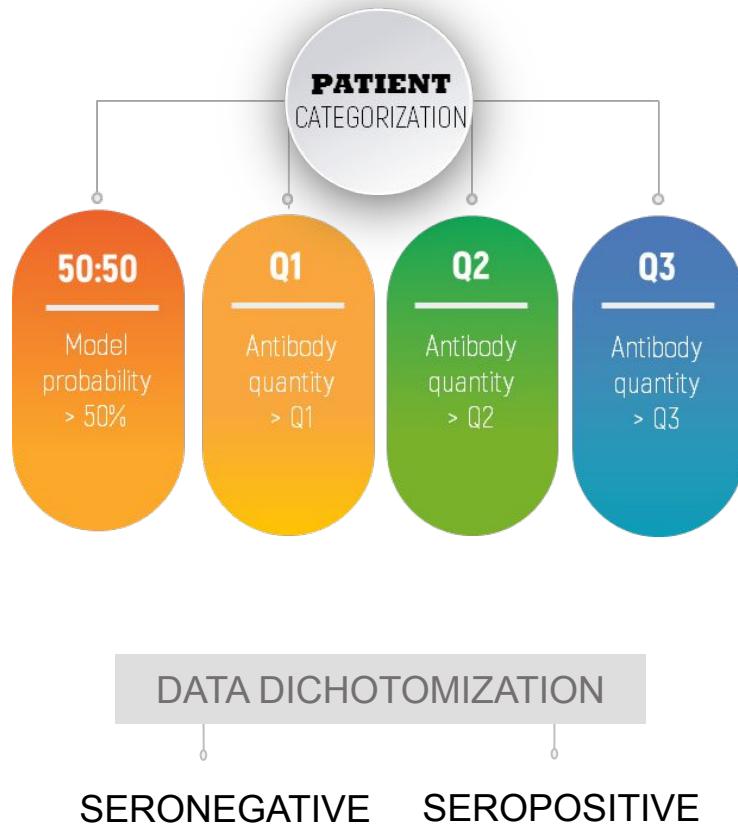


# Solution

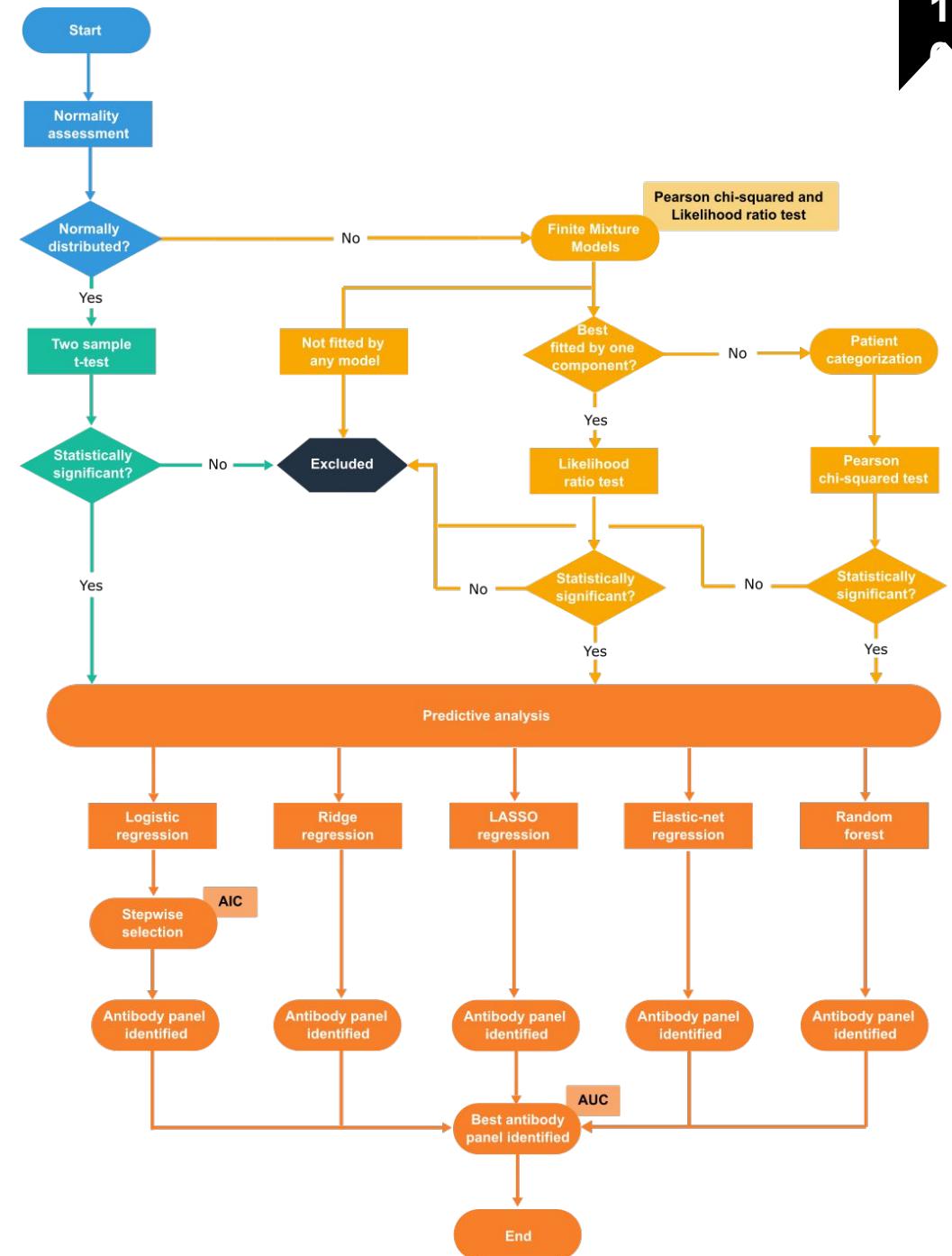




# Solution



Chi-squared test of Seropositivity against Status

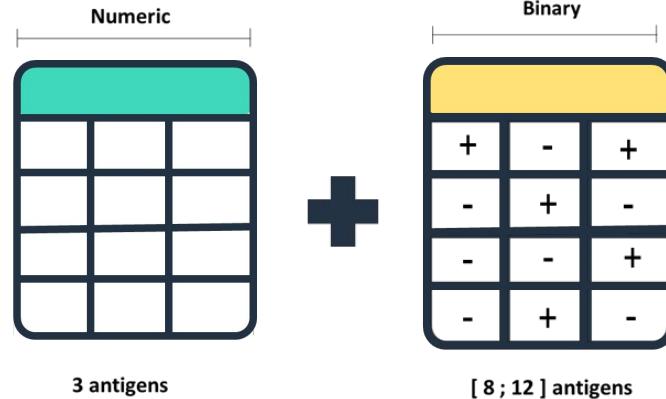




# Solution



## Predictor variables



## Response variable

Binary	
<b>Status</b>	
Protected	
Susceptible	
Susceptible	
Protected	

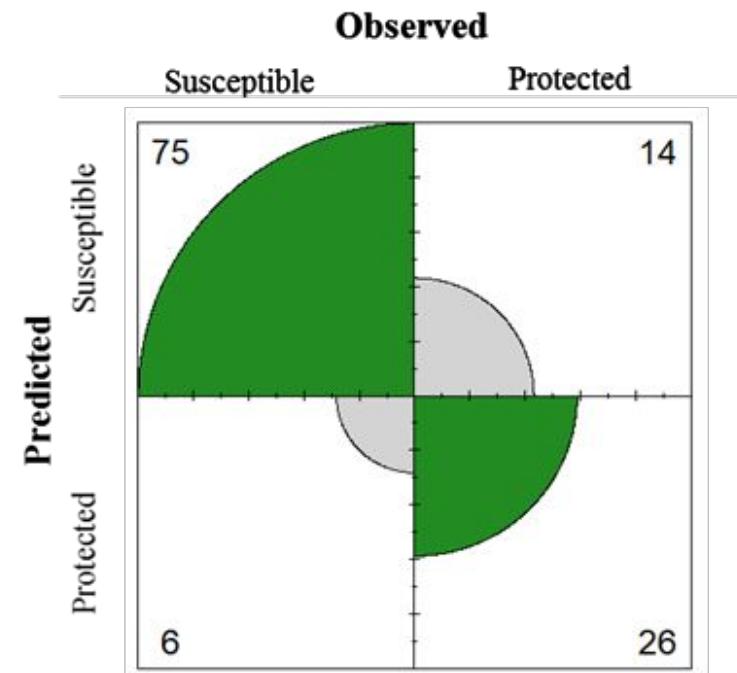
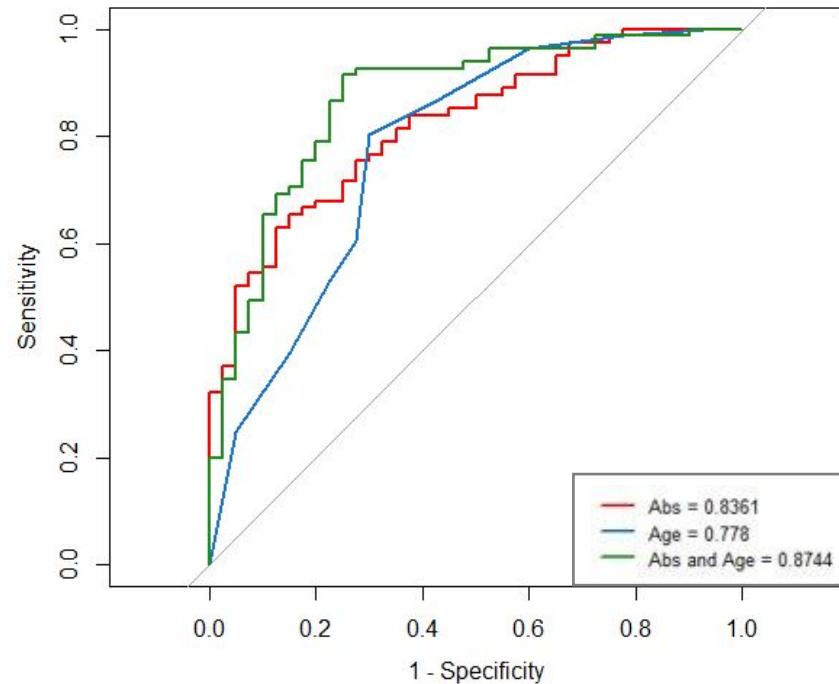




## Results



### Stepwise Logistic regression

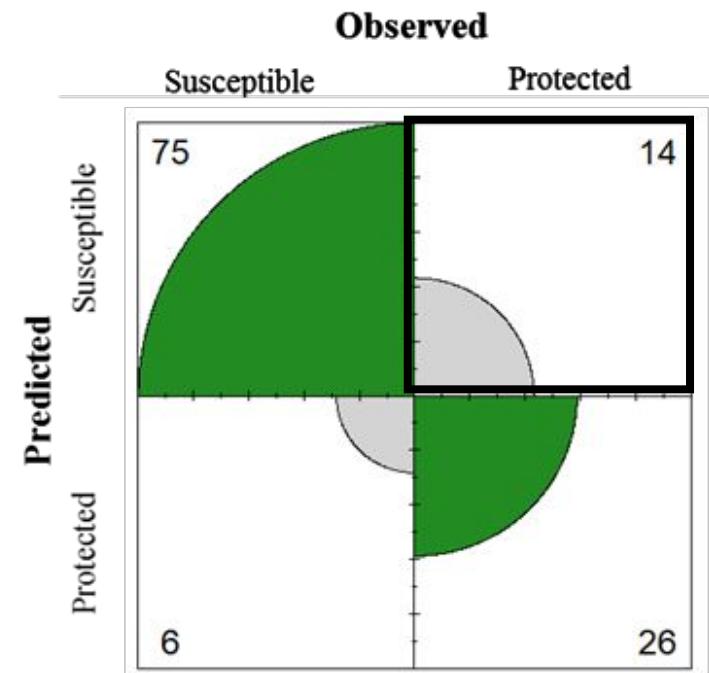
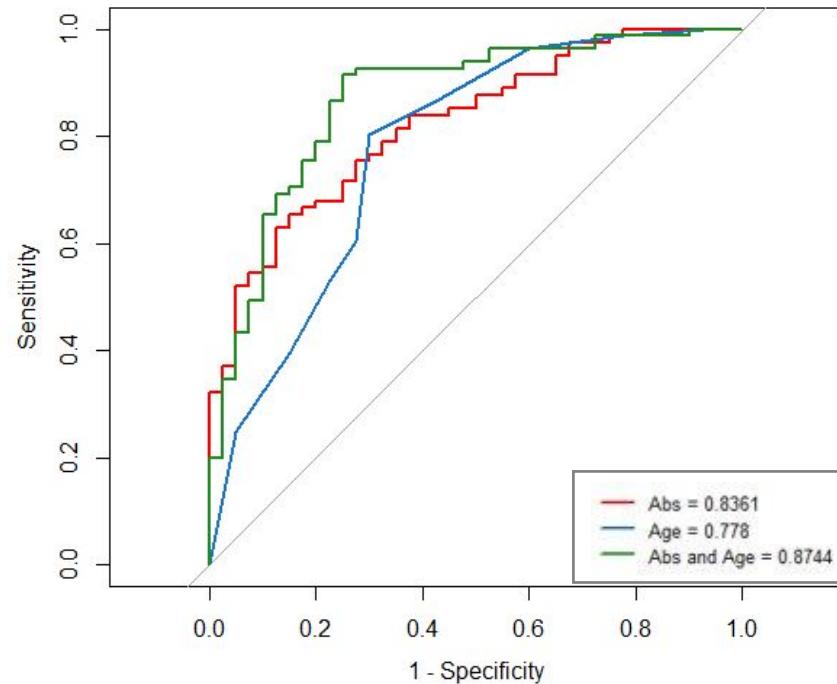




## Results



### Stepwise Logistic regression



Unable to correctly classify  
the protected individuals



## Results

10 fold CV



## Stepwise Logistic regression

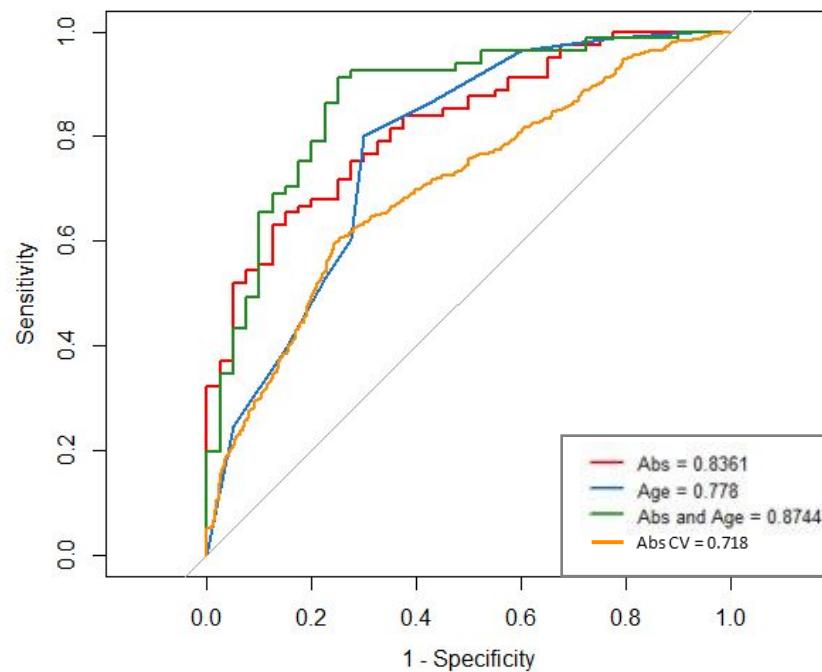


## Results

10 fold CV



## Stepwise Logistic regression



Low sample size ( $n=121$ )

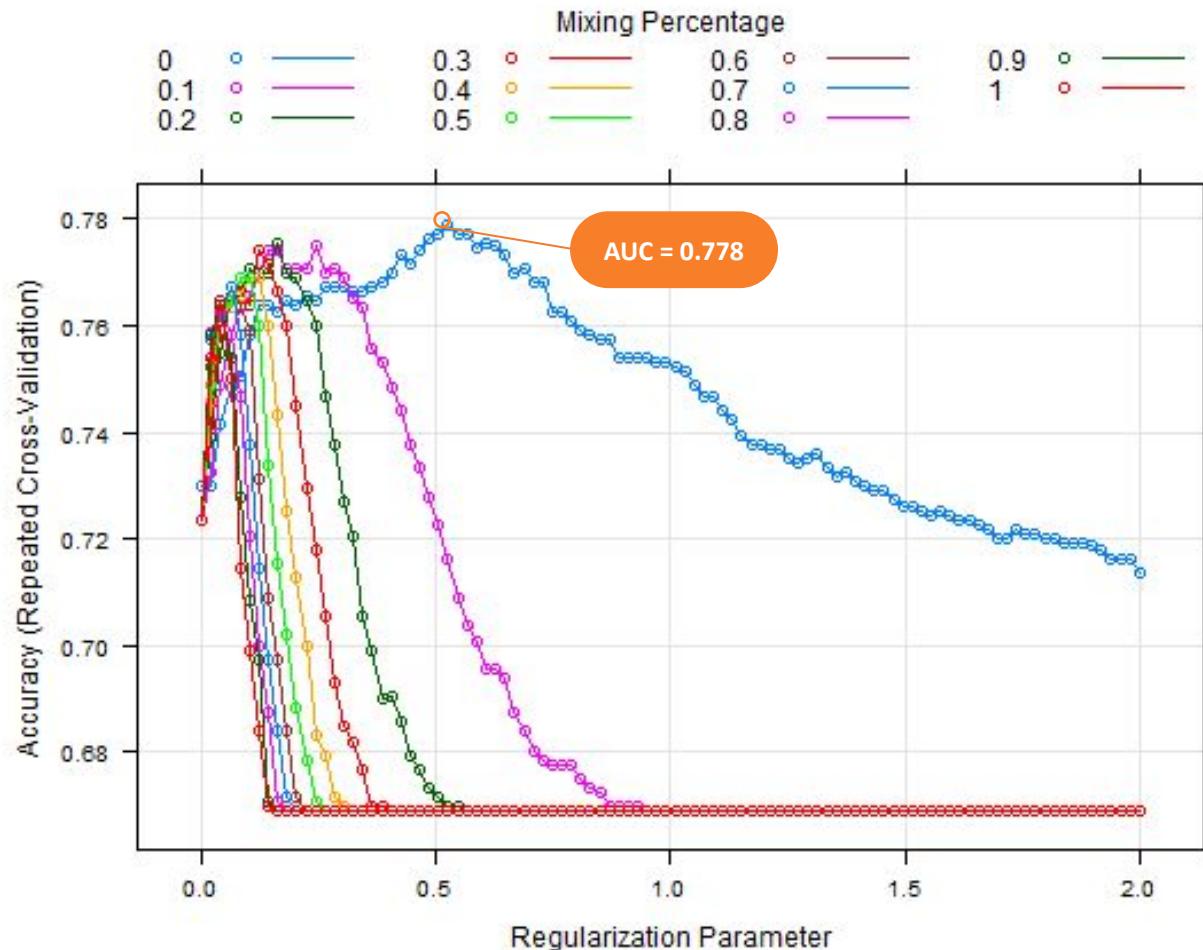


## Results

10 fold CV



## Elastic-net regression



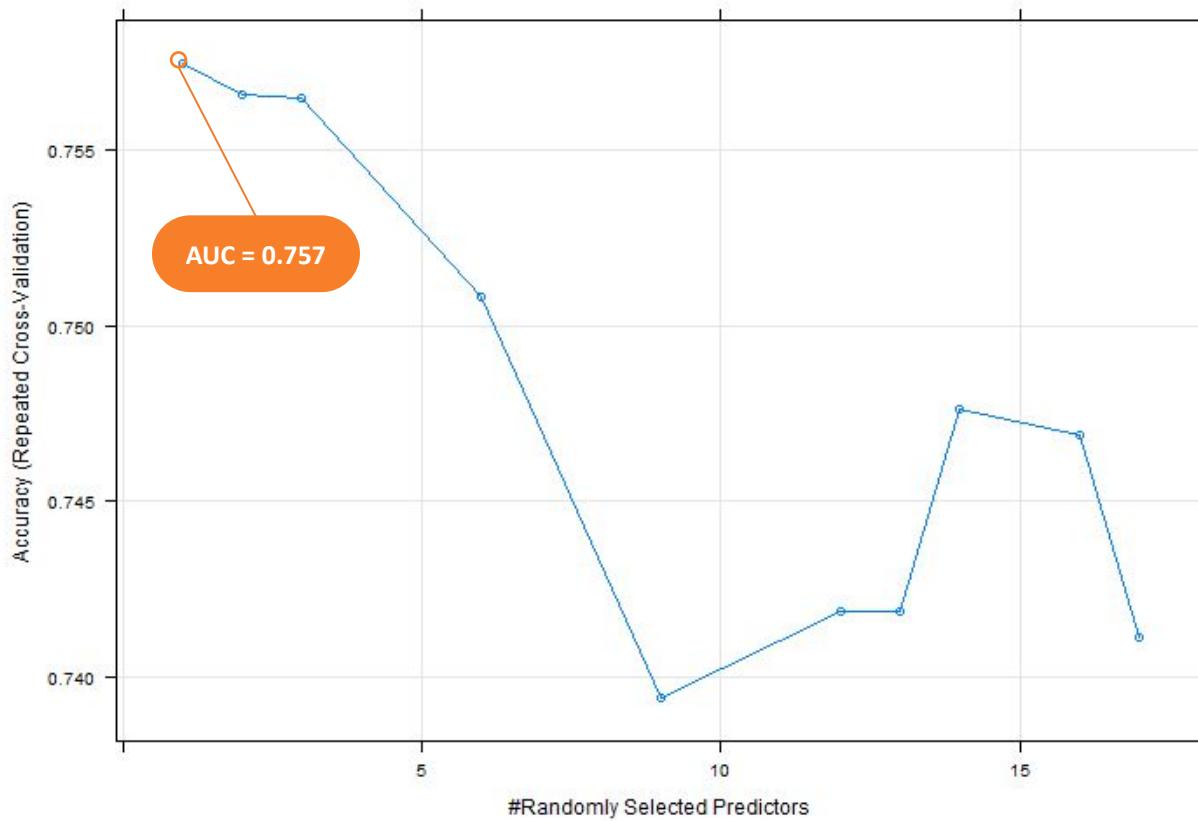


## Results

10 fold CV

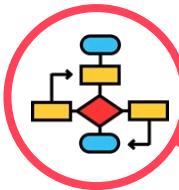


## Random Forest



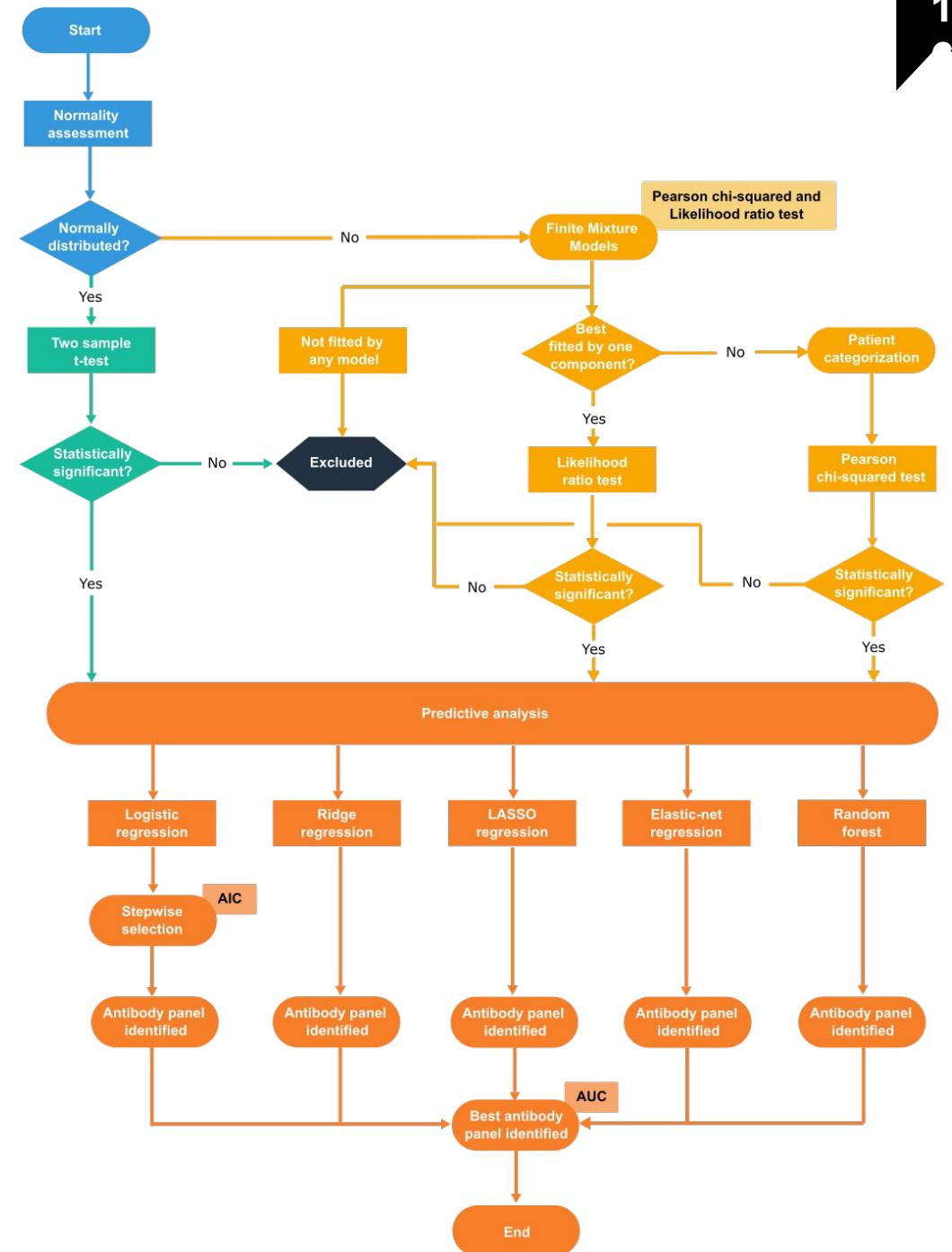


# Solution



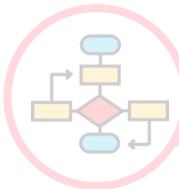
## Parametric pipeline

**Normal distribution and flexible mixture models for seropositivity determination**





## Solution



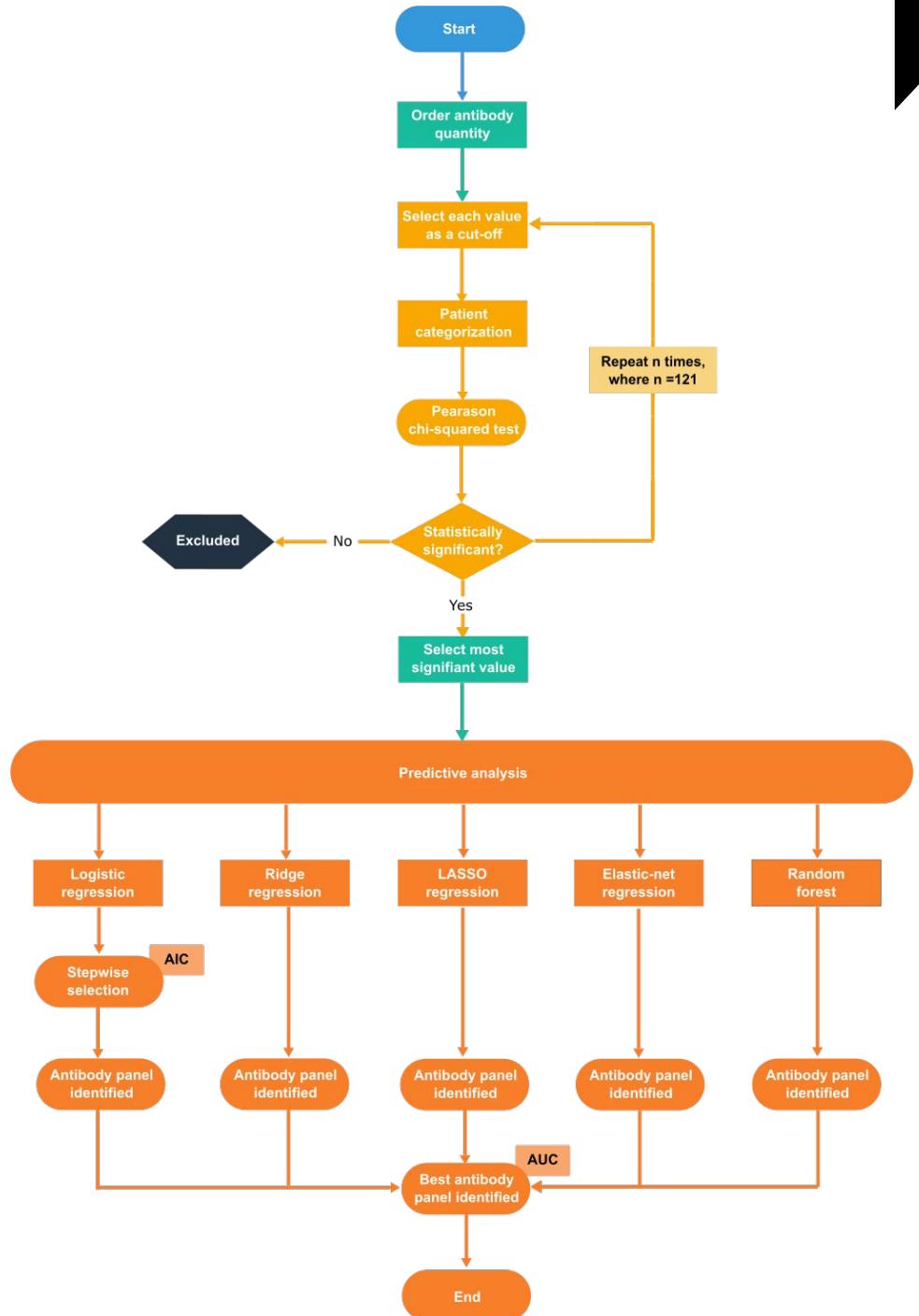
### Parametric pipeline

Normal distribution and flexible mixture models for seropositivity determination



### Pragmatic pipeline

Cut-off that maximized the distinction between susceptible and protected individuals

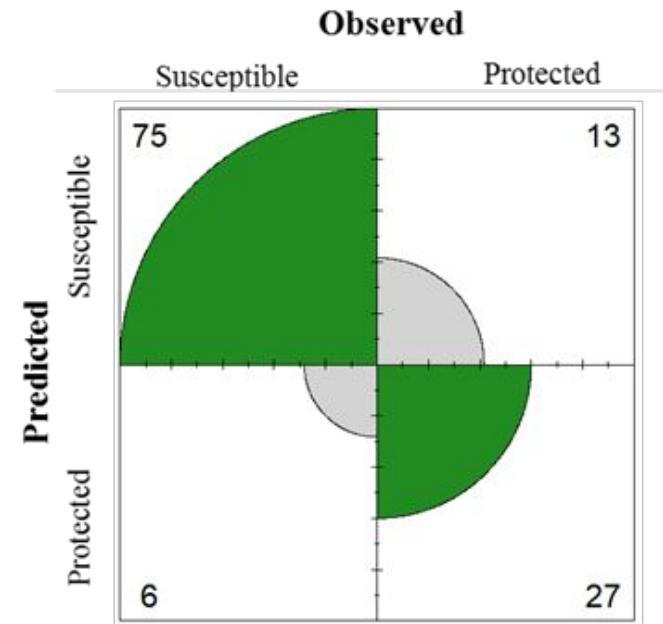
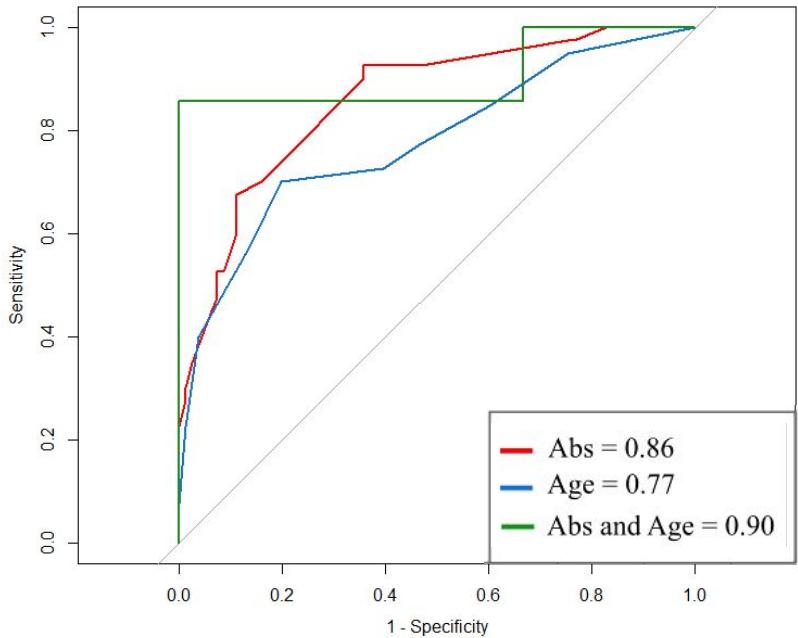




## Results



### Stepwise Logistic regression



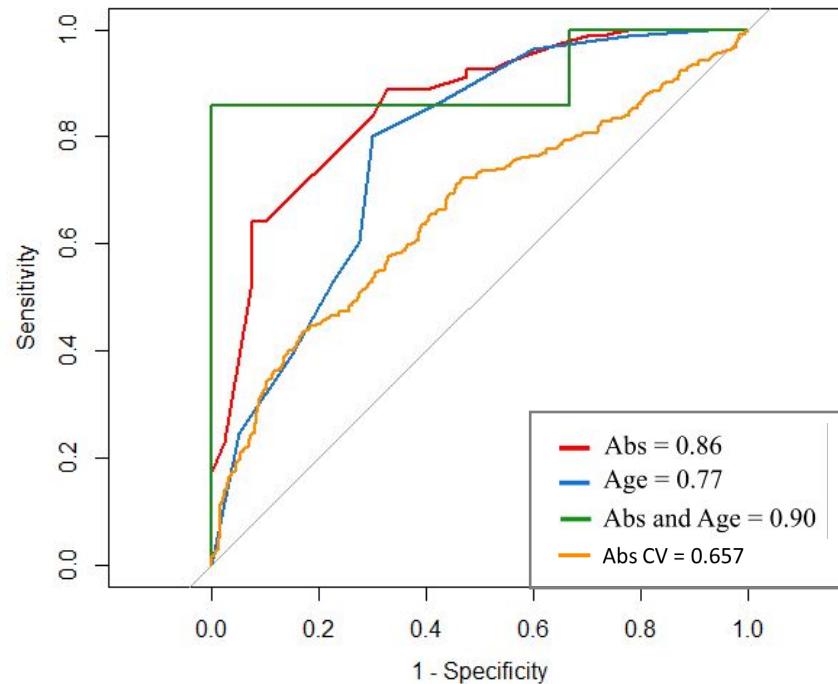


## Results

10 fold CV



## Stepwise Logistic regression



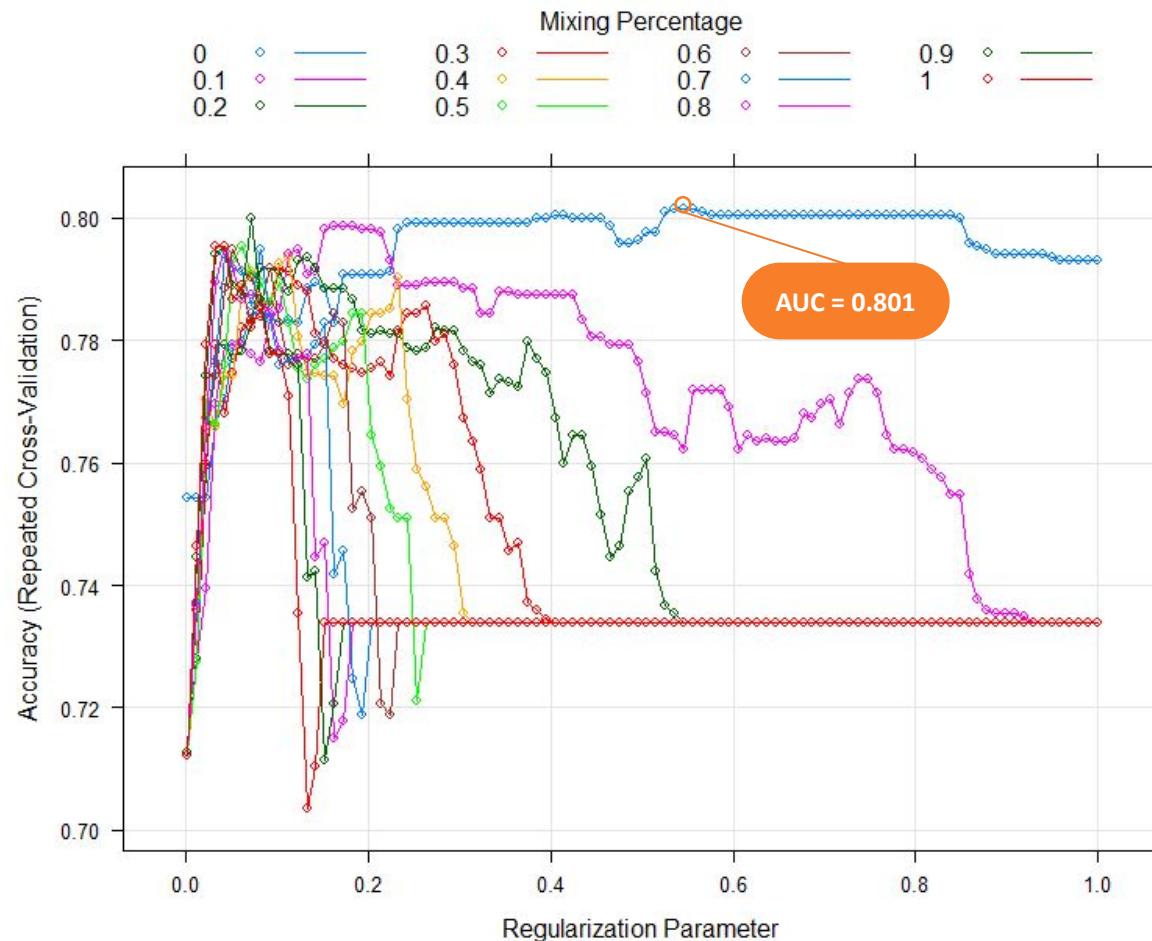


## Results

10 fold CV



## Elastic-net regression



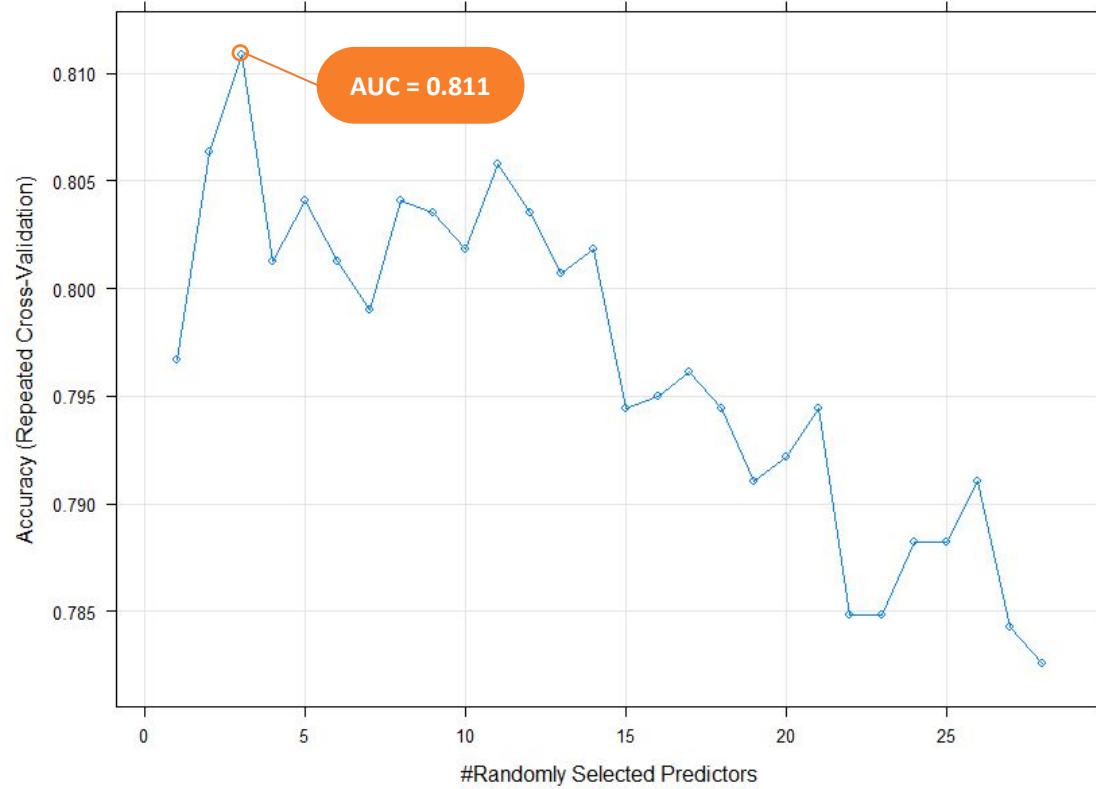


## Results

10 fold CV



### Random Forest





## *Setbacks*



Scale to higher dimension dataset



## Setbacks



### Scale to higher dimension dataset



IgG antibodies against **2320**  
**Plasmodium** **Falciparum**  
 antigens



Antibody quantity  
 information for **186 Malian**  
**individuals** (aged 2-25



years)  
 The **response variable**

"Status" was binary

(Protected vs. Susceptible)

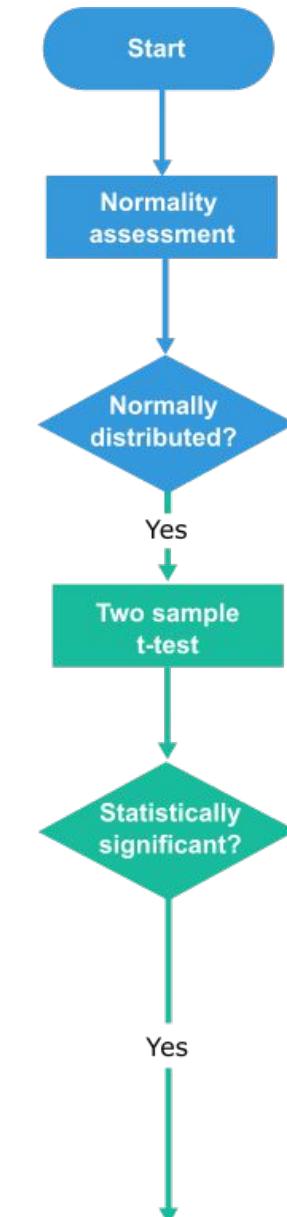
MSP	MSP	AMA	Status
2	7	1	Protected
			Susceptible
			Susceptible
			Protected

**2320 antigens**

**121 children**



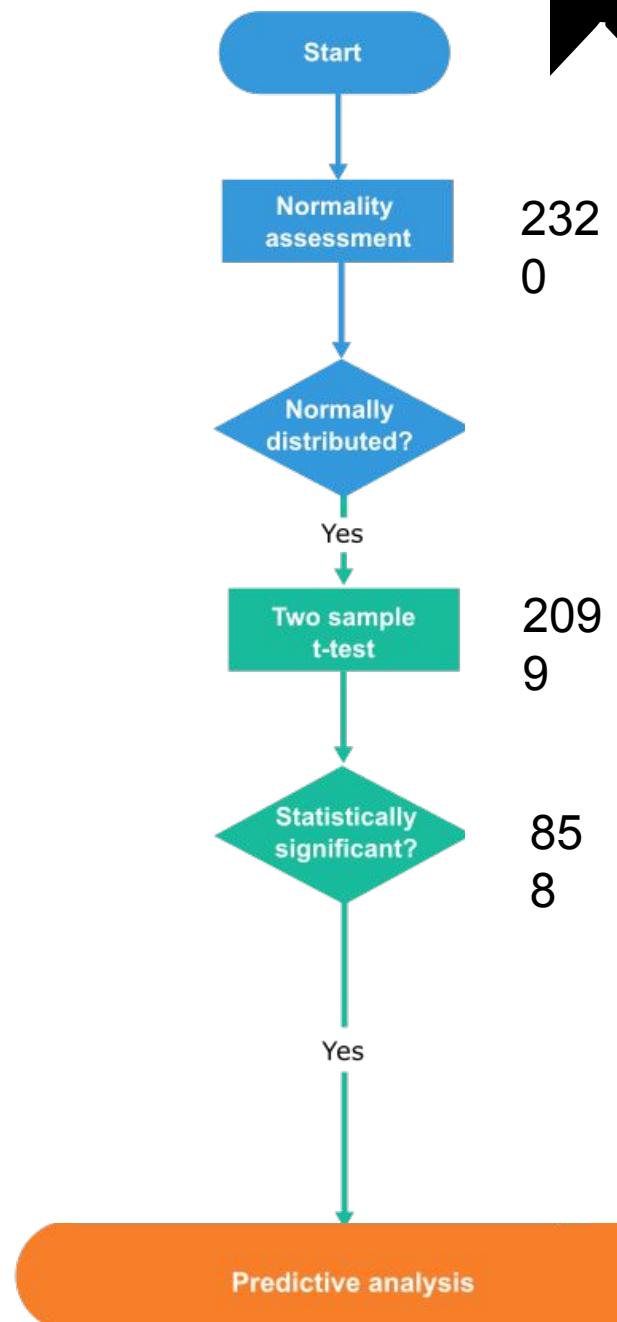
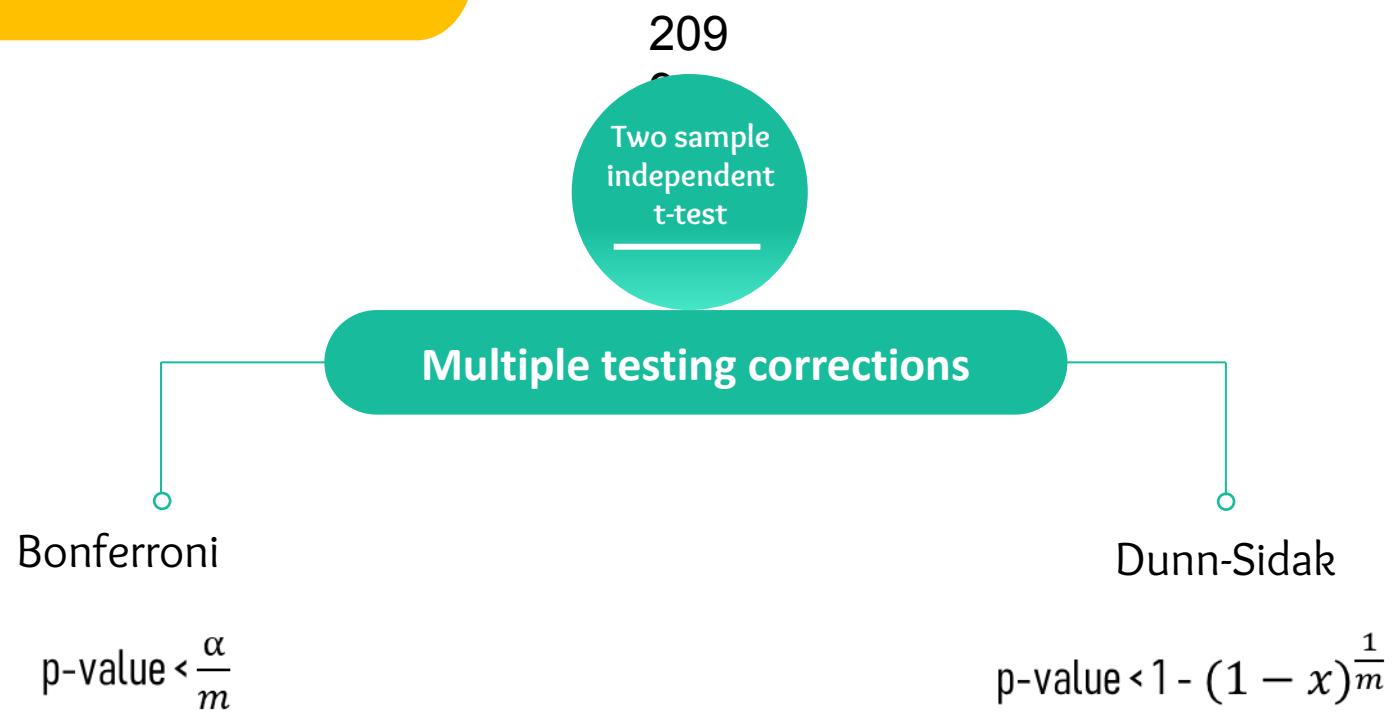
## Setbacks



Predictive analysis

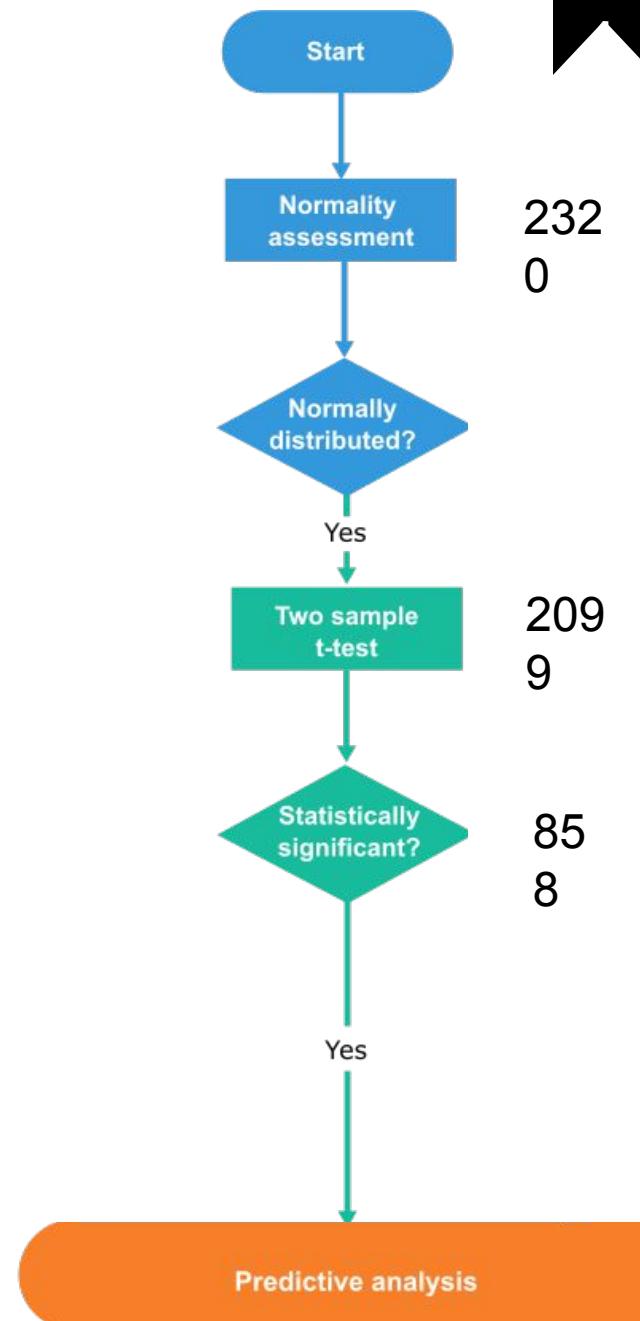
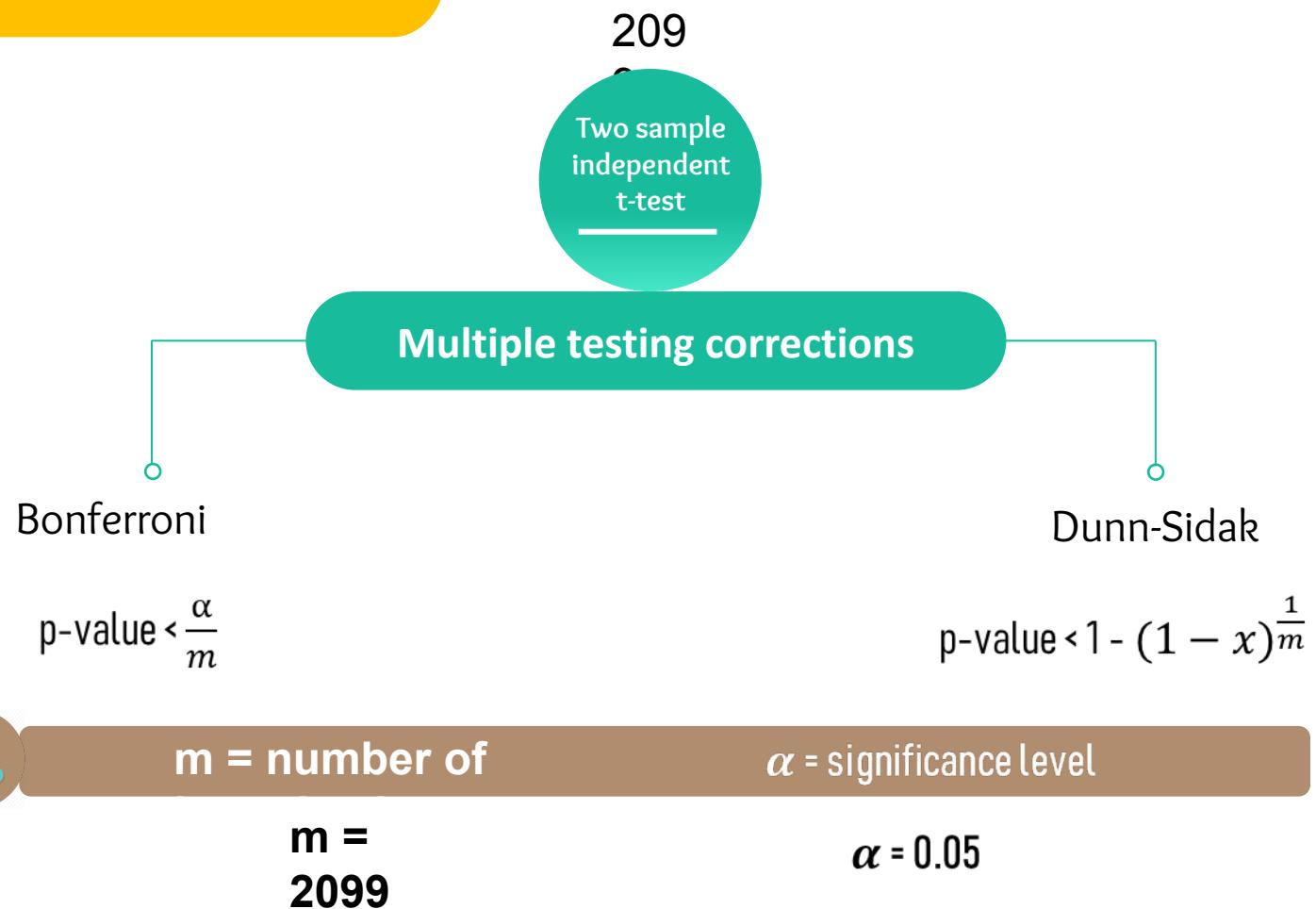


## Solution





## Solution





## Solution

209

Two sample  
independent  
t-test

### Multiple testing corrections

Bonferroni

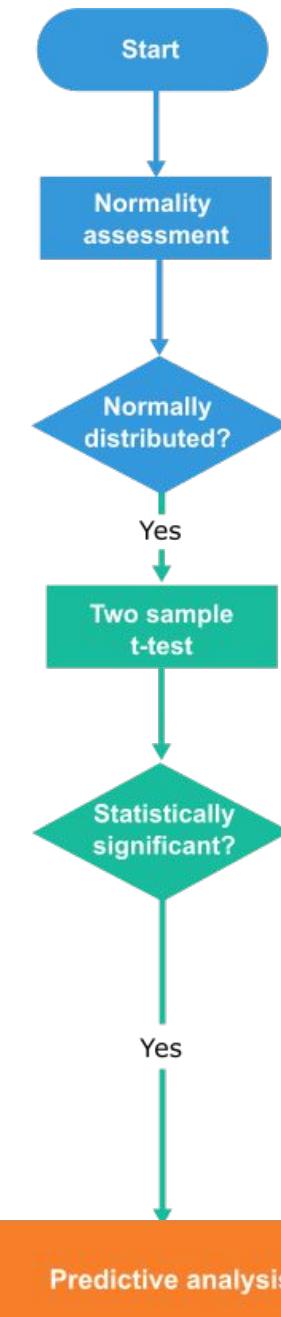
$$p\text{-value} < \frac{\alpha}{m}$$

Dunn-Sidak

$$p\text{-value} < 1 - (1 - x)^{\frac{1}{m}}$$

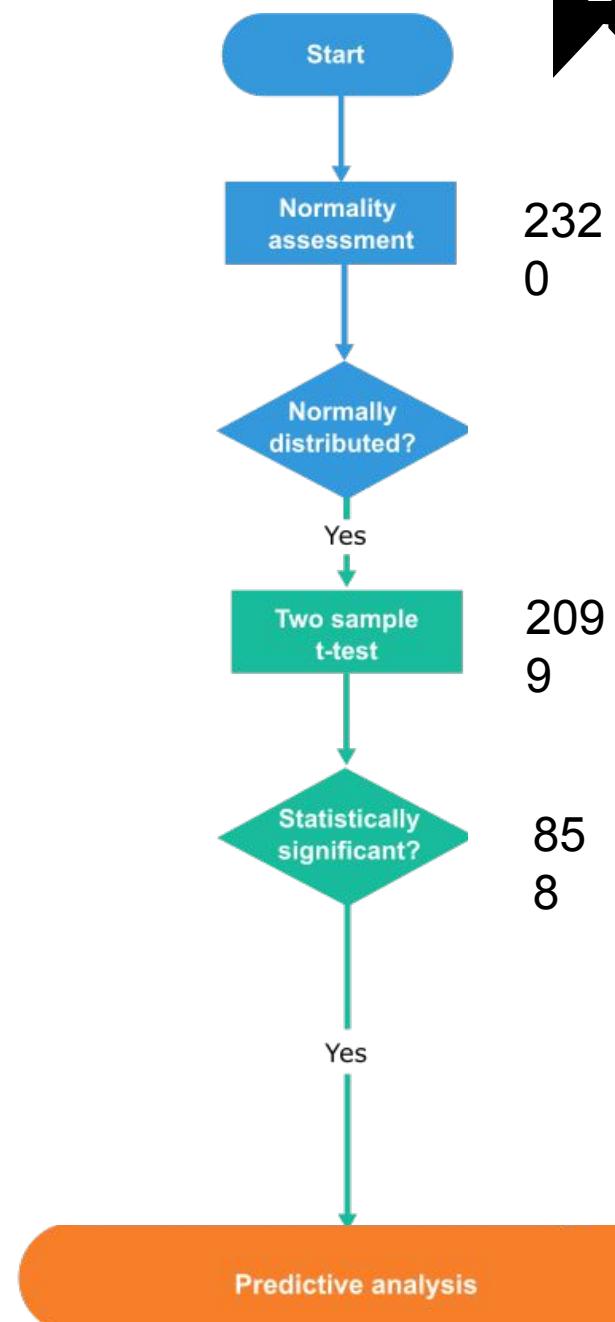
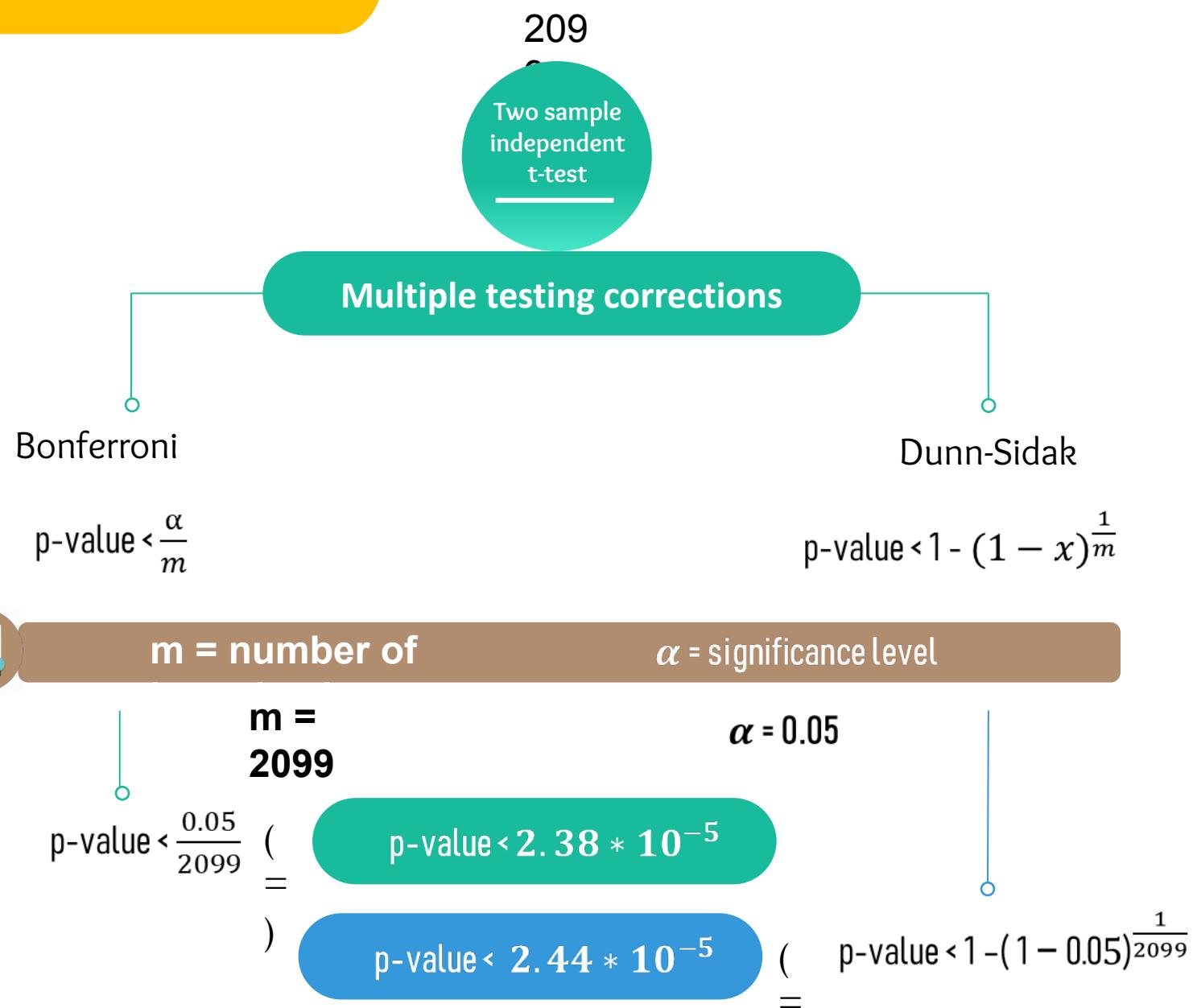
**m = number of** **$\alpha$  = significance level****m =  
2099** **$\alpha = 0.05$** 

$$p\text{-value} < \frac{0.05}{2099} \quad (= \quad p\text{-value} < 2.38 * 10^{-5} \quad )$$

232  
0209  
985  
8

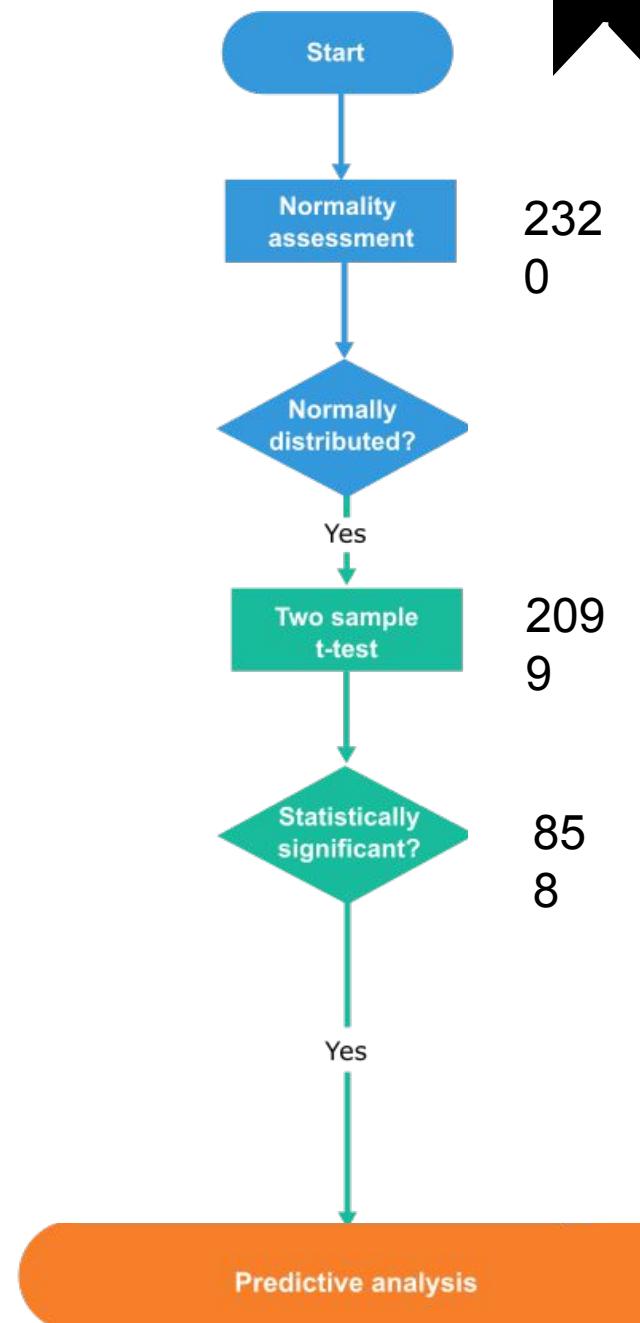
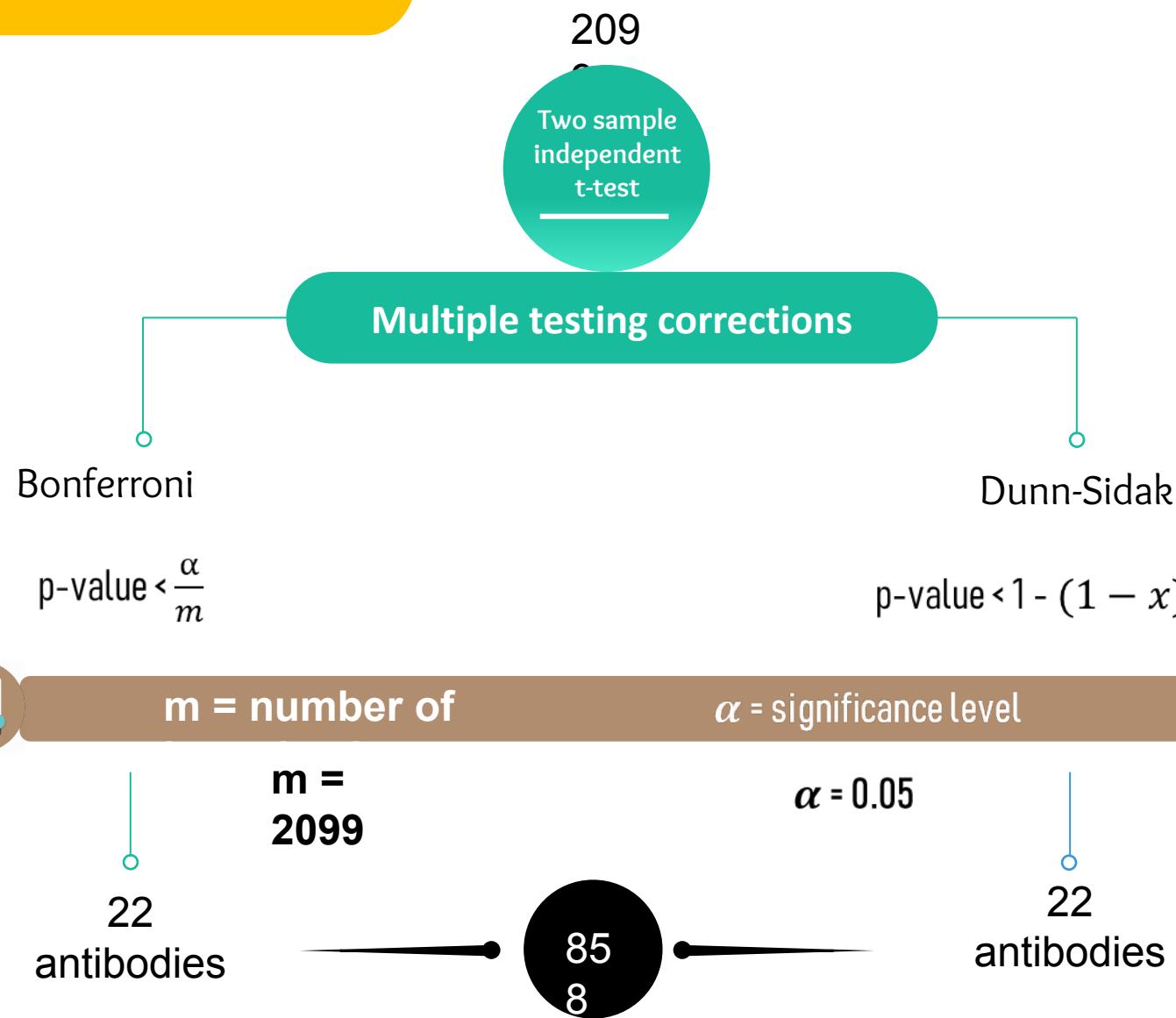


## Solution





## Solution

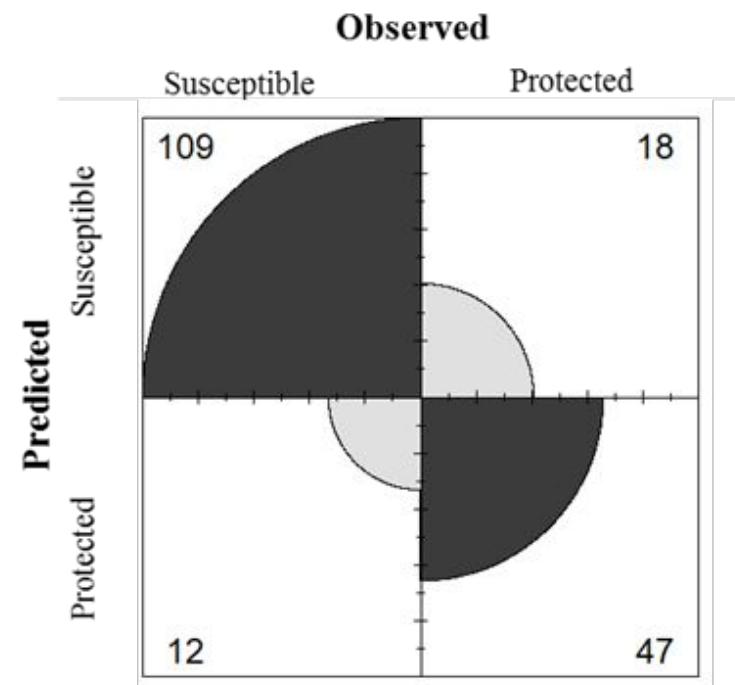
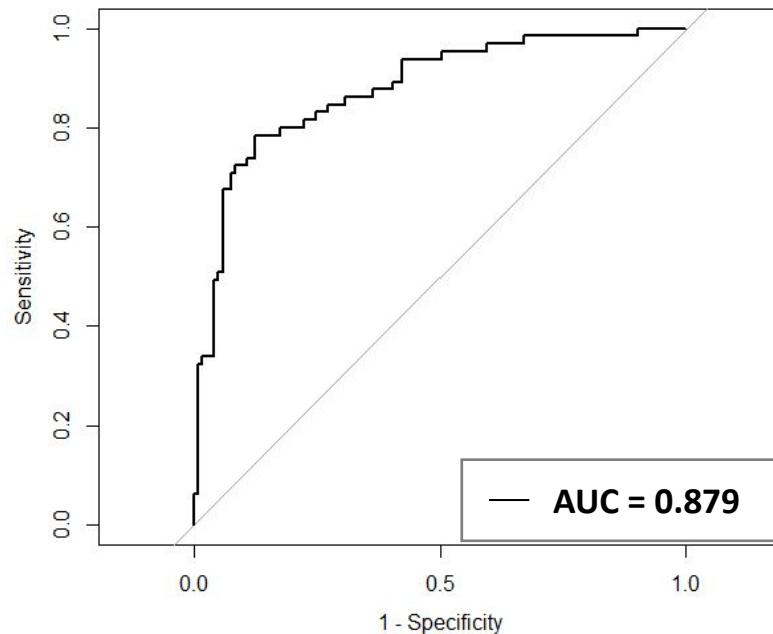




## Results



### Stepwise Logistic regression



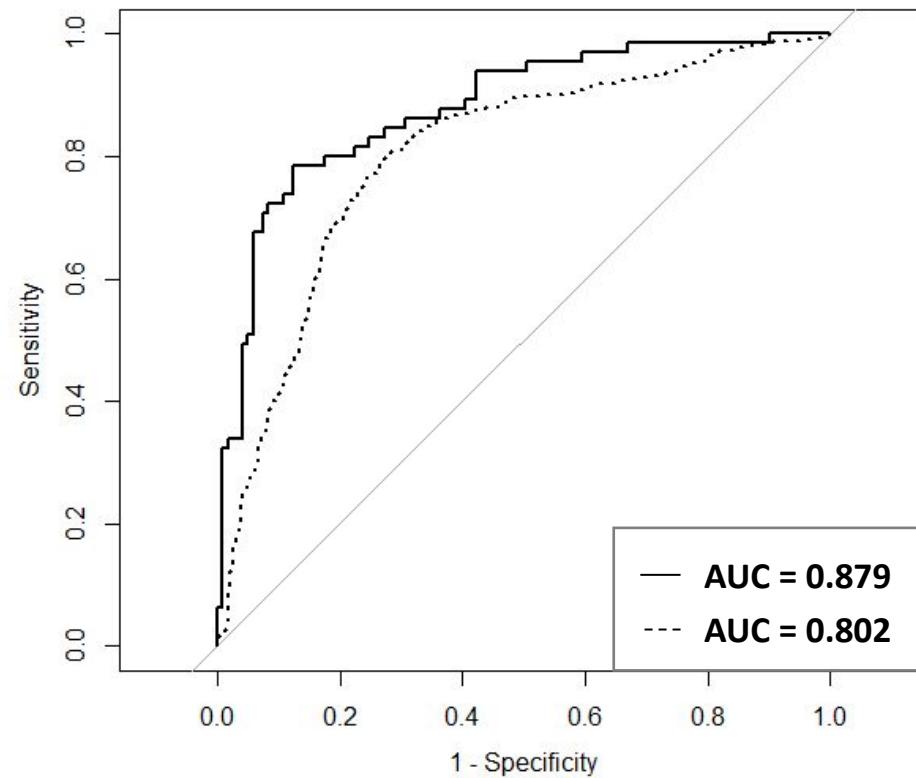


## Results

10 fold CV



### Stepwise Logistic regression



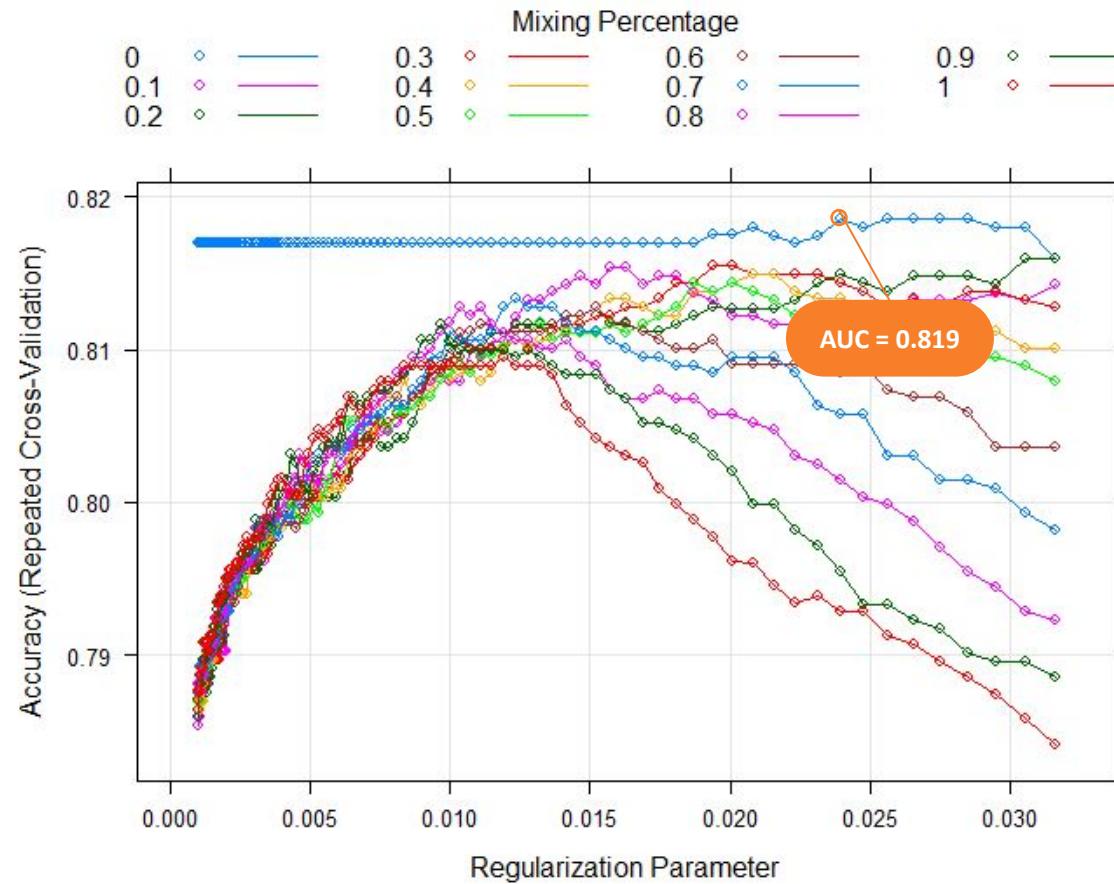


## Results

10 fold CV



### Elastic-net regression



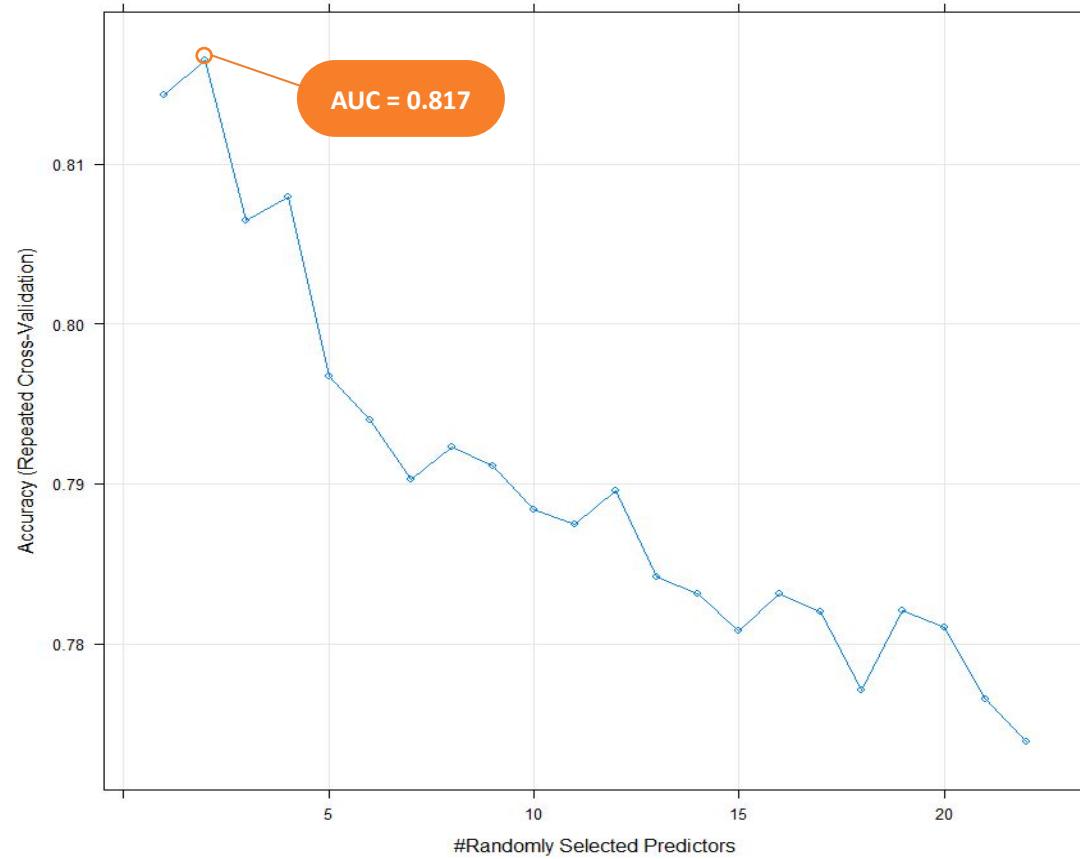


## Results

10 fold CV



### Random Forest



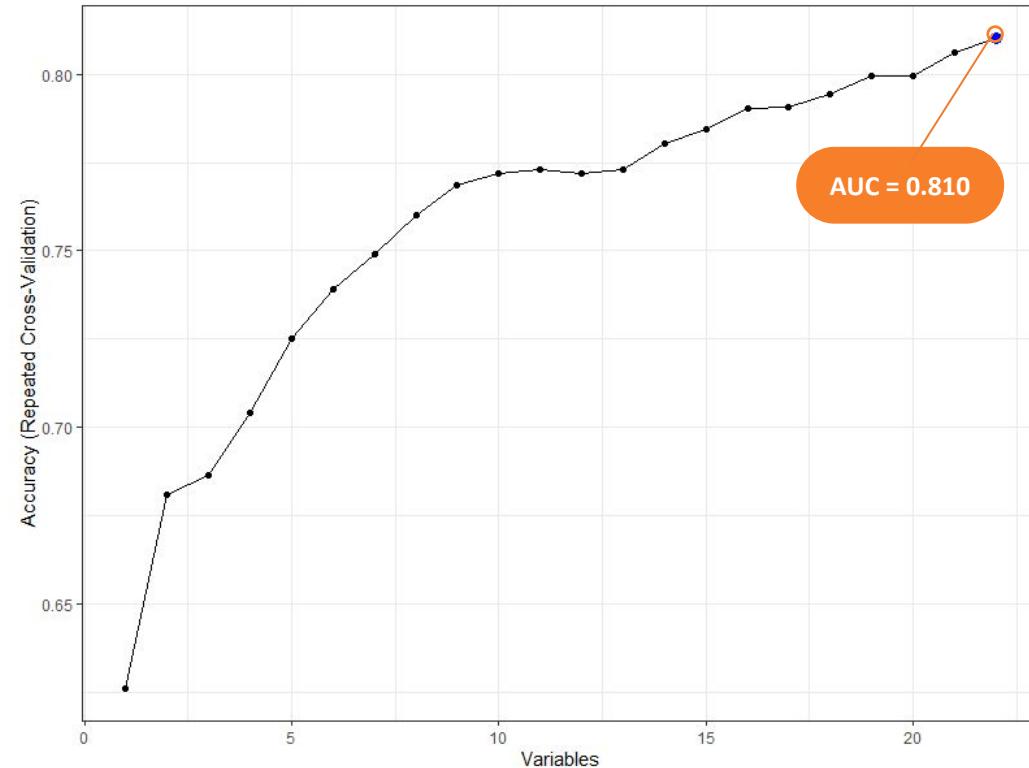


## Results

10 fold CV



### Recursive Feature Elimination



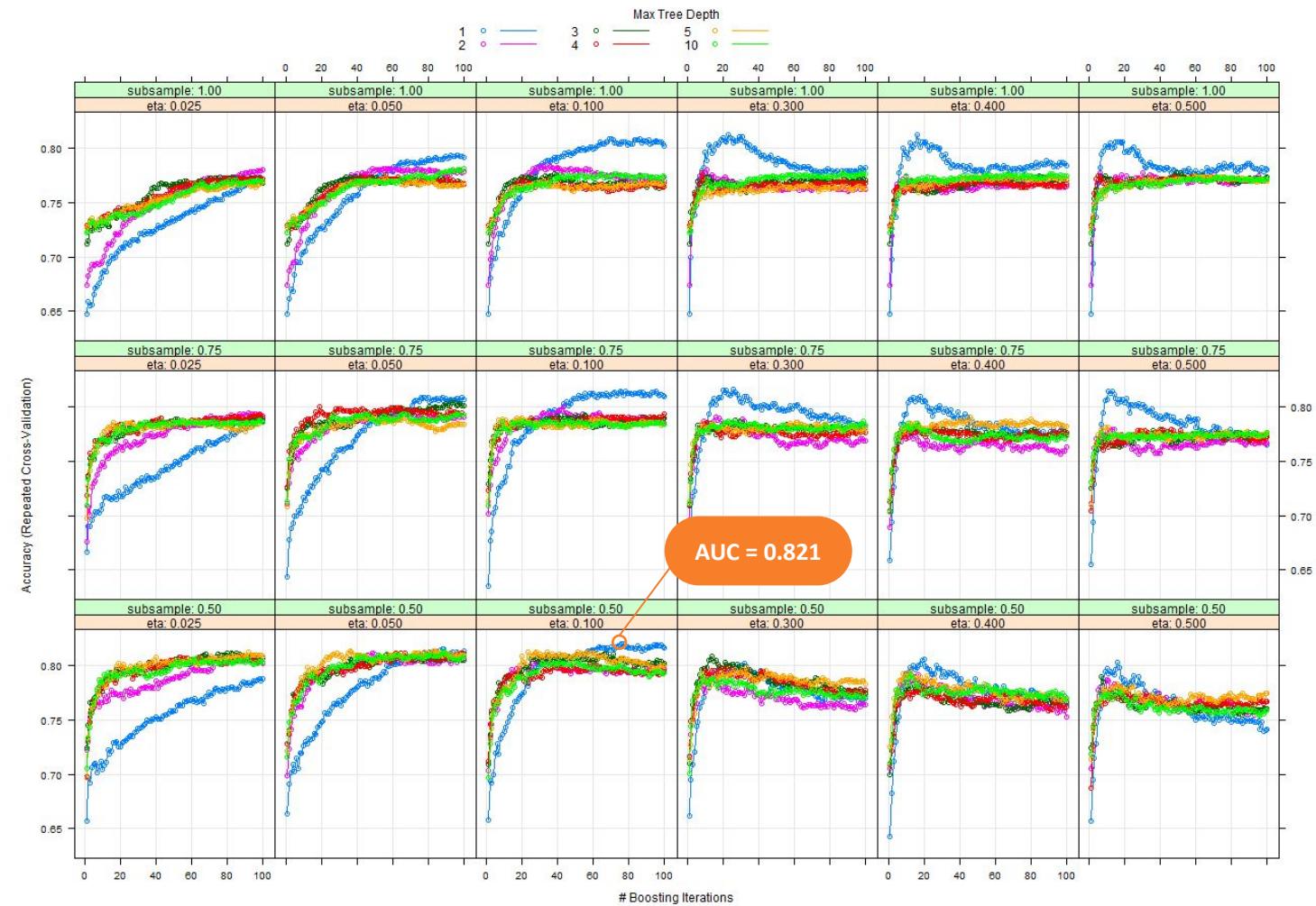


# Results

10 fold CV

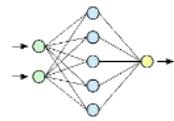


## Stochastic gradient boosting

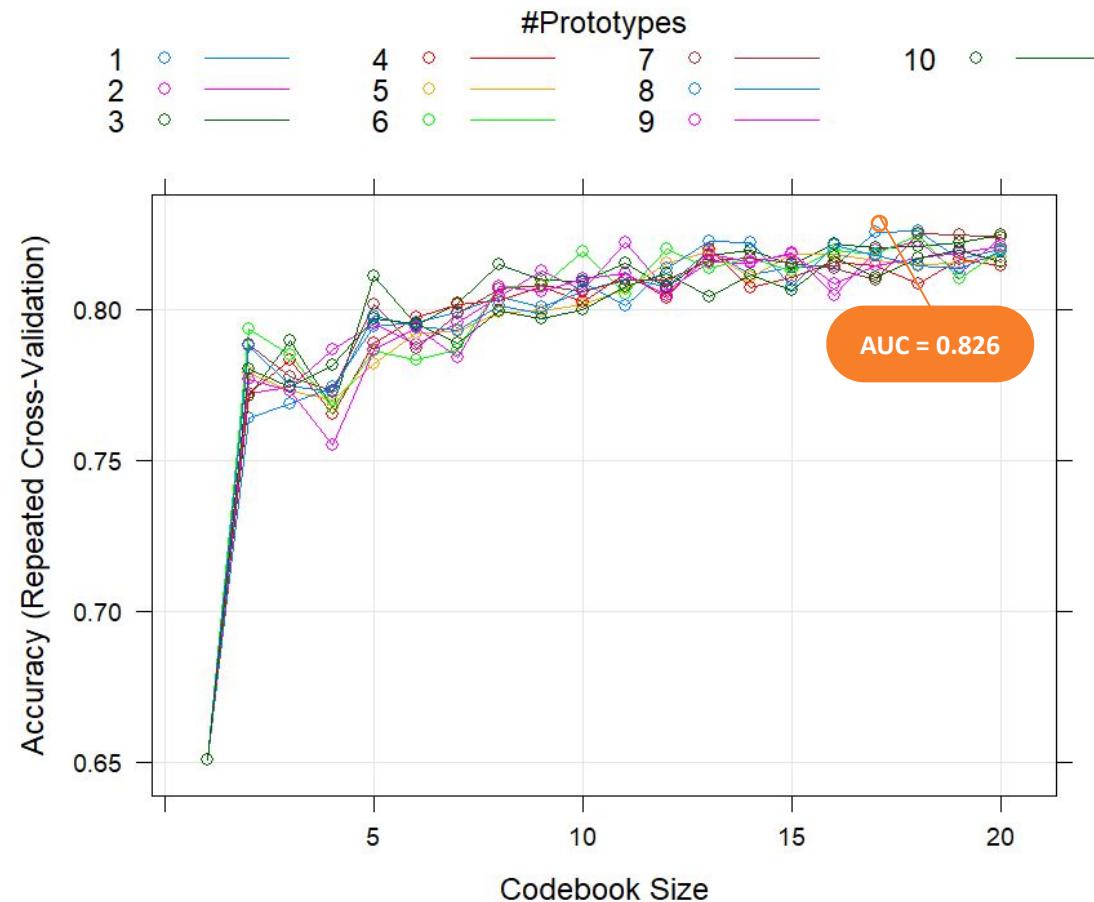




## Results



### Learning Vector Quantization





## *Solution*

What if instead of relying solely on the corrected p-value of the T-test we also rely on the **power** of the result?



# Results

R

```
library(pwr)
```

Power calculations for two samples (different sizes) t-tests of means

## Description

Compute power of tests or determine parameters to obtain target power (similar to as power.t.test).

## Usage

```
pwr.t2n.test(n1 = NULL, n2 = NULL, d = NULL, sig.level = 0.05, power = NULL,
  alternative = c("two.sided",
    "less", "greater"))
```

## Arguments

<code>n1</code>	Number of observations in the first sample
<code>n2</code>	Number of observations in the second sample
<code>d</code>	Effect size
<code>sig.level</code>	Significance level (Type I error probability)
<code>power</code>	Power of test (1 minus Type II error probability)
<code>alternative</code>	a character string specifying the alternative hypothesis, must be one of "two.sided" (default), "greater" or "less"



## Arguments

**n1 = number of susceptibles (121)**

**n2= number of protected (65)**

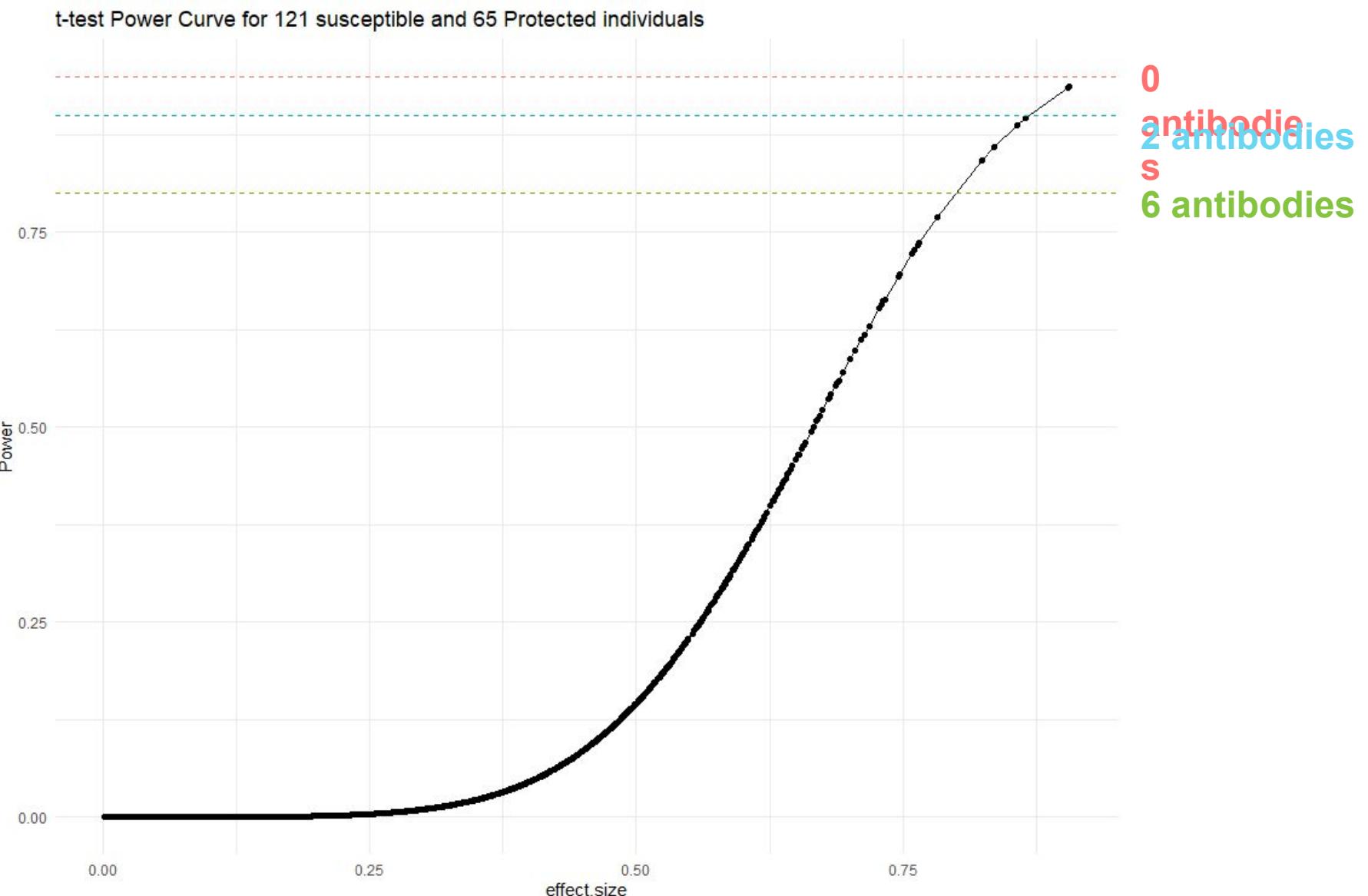
$$d = \frac{\text{mean}(x_1) - \text{mean}(x_2)}{\sqrt{\frac{(n_1-1).sd_1^2 + (n_2-1).sd_2^2}{n_1+n_2-2}}}$$

**sig.level =  $2.38 \times 10^{-5}$**

**alternative = "two-sided"**



# Results

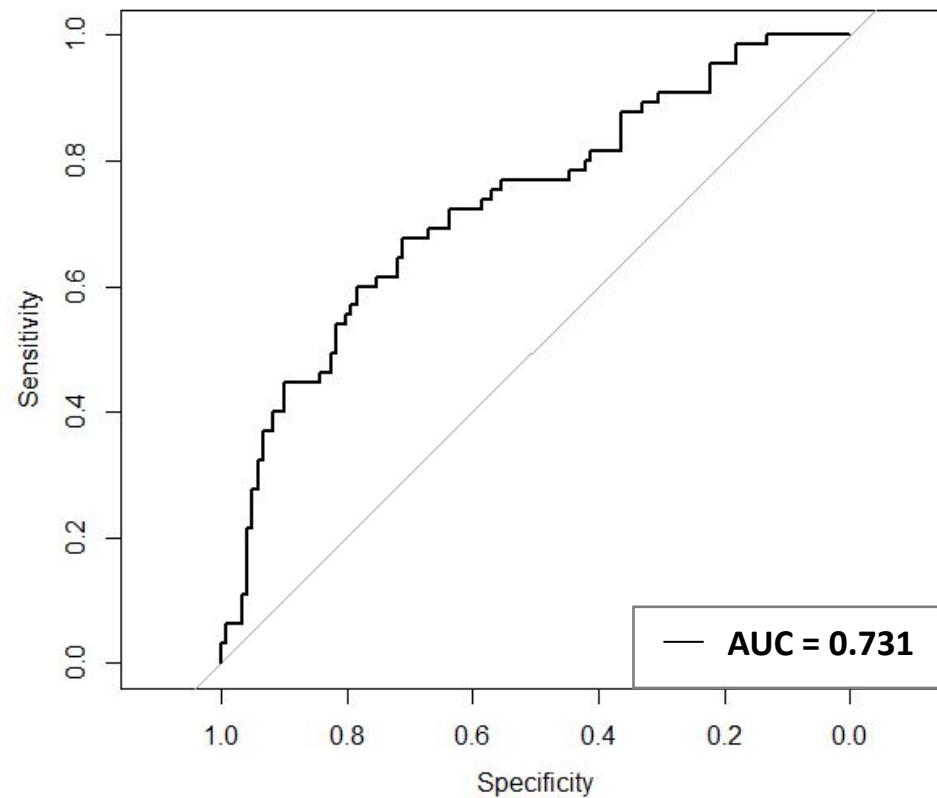




## Results



### Logistic regression





## *Setbacks*



What about the other antibodies that did not follow a normal distribution ( $p=221$ ) ?





## *Discussion*



Our pipeline needs some tweaking to be more suitable for higher dimension datasets



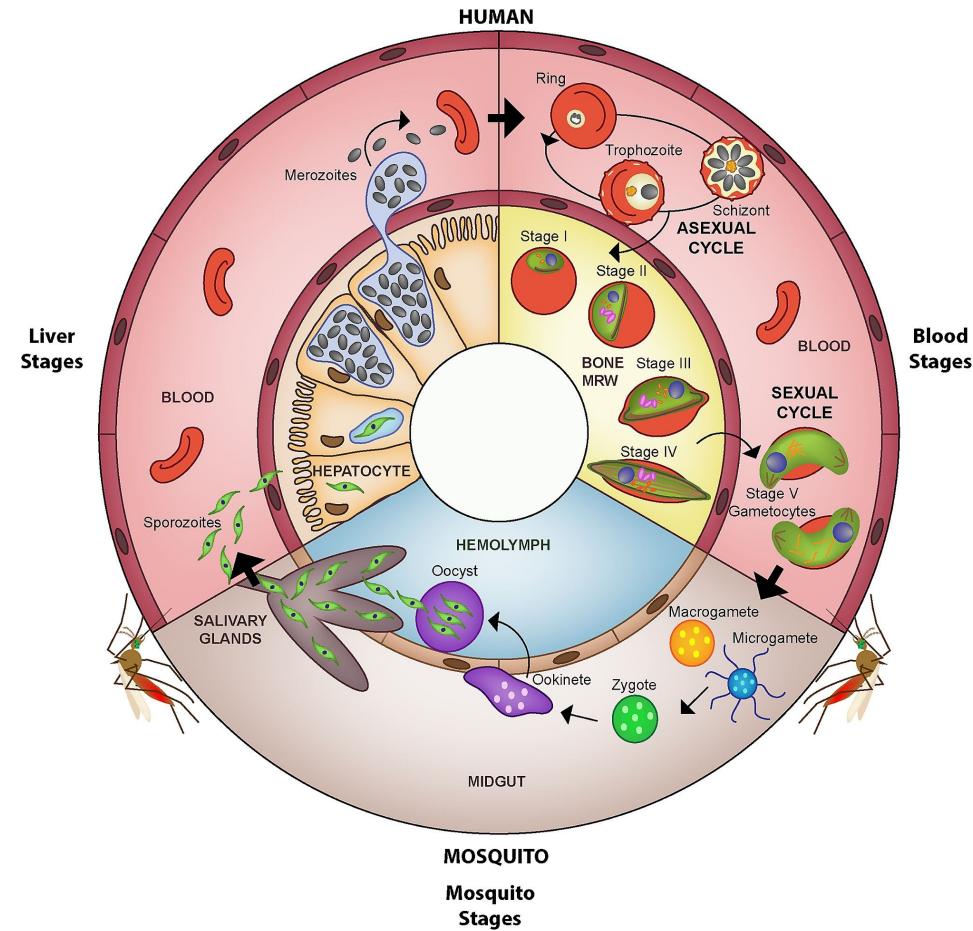
## Discussion



Our pipeline needs some tweaking to be more suitable for higher dimension datasets



High dimension datasets difficult the identification of important antibodies. An alternative could be to analyse antibodies in a more biological oriented approach (e.g. MSP)





## Discussion



Our pipeline needs some tweaking to be more suitable for higher dimension datasets



High dimension datasets difficult the identification of important antibodies. An alternative could be to analyse antibodies in a more biological oriented approach (e.g. MSP)



Pipelines such as these may help to increase reproducibility among studies



## Discussion



Our pipeline needs some tweaking to be more suitable for higher dimension datasets



High dimension datasets difficult the identification of important antibodies. An alternative could be to analyse antibodies in a more biological oriented approach (e.g. MSP)



Pipelines such as these may help to increase reproducibility among studies



Since antibody data transcends the field of malaria these pipelines may provide promising tools to be applied in other settings