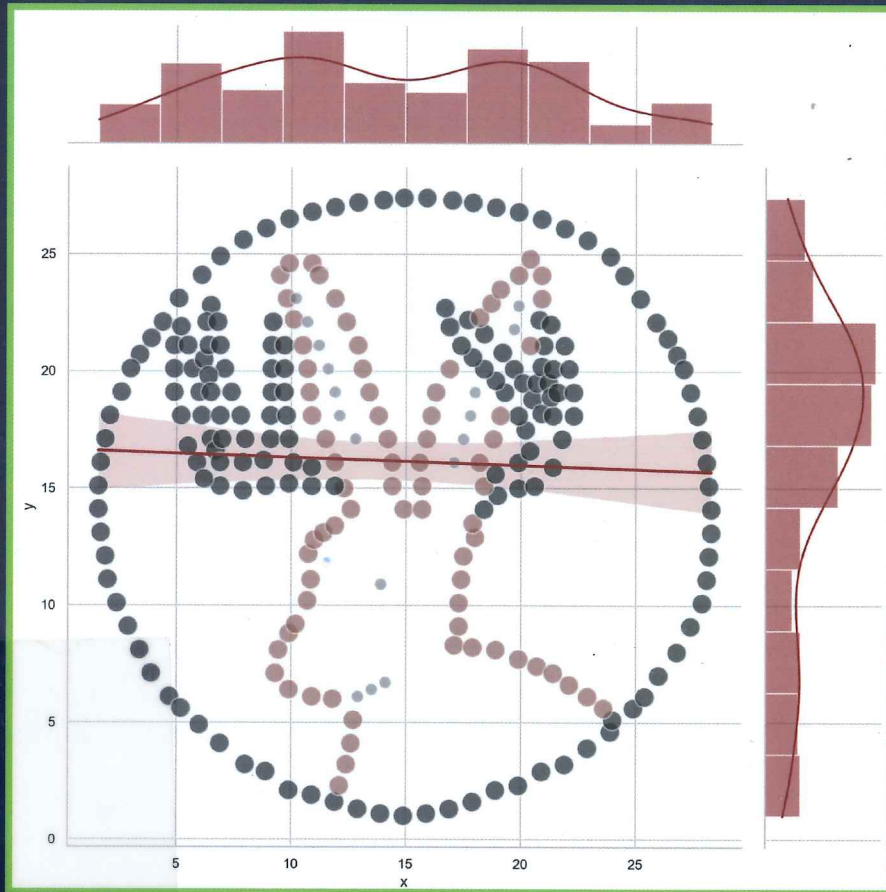# STATISTICS AND DATA VISUALISATION WITH PYTHON



## JESÚS ROGEL-SALAZAR

CRC Press
Taylor & Francis Group

A CHAPMAN & HALL BOOK

# Statistics and Data Visualisation with Python

Jesús Rogel-Salazar

# Contents