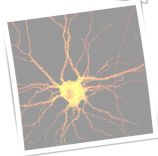
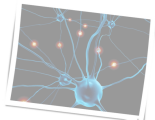
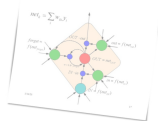
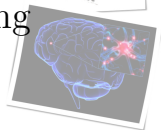


A quick introduction to machine learning

Spyros Samothrakis
Senior Lecturer, IADS
University of Essex
MiSoC

June 22, 2022



WELCOME/COURSE CONTENTS

- ▶ What will this course cover?
 - ▶ Day 1: An intro to machine learning (ML)
 - ▶ Day 1: ML labs
 - ▶ Day 2: An intro to causal inference
 - ▶ Day 2: ML and causal inference labs
- ▶ Textbooks?
 - ▶ Mitchell, T. M. (1997). Machine learning.¹
 - ▶ Bishop, C. M. (2006). Pattern recognition and machine learning. springer.²
 - ▶ Wasserman, L. (2013). All of statistics: a concise course in statistical inference. Springer Science & Business Media.³

¹<http://www.cs.cmu.edu/~tom/mlbook.html>

²<https://www.microsoft.com/en-us/research/publication/pattern-recognition-machine-learning/>

³<http://www.stat.cmu.edu/~larry/all-of-statistics/index.html>

BETTER SCIENCE THROUGH DATA

Hey, Tony, Stewart Tansley, and Kristin M. Tolle. “Jim Gray on eScience: a transformed scientific method.” (2009).⁴

- ▶ Thousand years ago: empirical branch
 - ▶ You observed stuff and you wrote down about it
- ▶ Last few hundred years: theoretical branch
 - ▶ Equations of gravity, equations of electromagnetism
- ▶ Last few decades: computational branch
 - ▶ Modelling at the micro level, observing at the macro level
- ▶ Today: data exploration
 - ▶ Let machines create models using vast amounts of data

⁴<http://language.log.ldc.upenn.edu/myl/JimGrayOnE-Science.pdf>

BETTER BUSINESS THROUGH DATA

- ▶ There was a report by Mckinsey

Manyika, J., Chui, M., Brown, B., Bughin, J., Dobbs, R., Roxburgh, C., & Hung Byers, A. (2011). Big data: The next frontier for innovation, competition, and productivity. McKinsey Global Institute.⁵

- ▶ Urges everyone to monetise “Big Data”
- ▶ Use the data provided within your organisation to gain insights
- ▶ Has some numbers as to how much this is worth
- ▶ Proposes a number of methods, most of them associated with machine learning and databases

⁵<http://www.mckinsey.com/business-functions/digital-mckinsey/our-insights/big-data-the-next-frontier-for-innovation>

WHY IS IT POPULAR NOW?

- ▶ **Algorithms + data + tools**
- ▶ Breiman, L. (2001). Statistical modeling: The two cultures (with comments and a rejoinder by the author). Statistical science, 16(3), 199-231.⁶
- ▶ Anderson, P. W. (1972). More is different. Science, 177(4047), 393-396.⁷
- ▶ Pedregosa, et.al. (2011). Scikit-learn: Machine learning in Python. the Journal of machine Learning research, 12, 2825-2830.⁸

⁶http://projecteuclid.org/download/pdf_1/euclid.ss/1009213726%20

⁷https://www.tkm.kit.edu/downloads/TKM1_2011_more_is_different_PWA.pdf

⁸<https://www.jmlr.org/papers/volume12/pedregosa11a/pedregosa11a.pdf>

⁸<https://www.jmlr.org/papers/volume12/pedregosa11a/pedregosa11a.pdf>

<https://www.jmlr.org/papers/volume12/pedregosa11a/pedregosa11a.pdf>

SO THIS COURSE COVERS TOOLS

- ▶ ML theory
 - ▶ *Supervised learning Regression Classification*
 - ▶ Understanding basic modelling
 - ▶ Confirming your model is sane
 - ▶ Tuning your model
 - ▶ **All within a very applied setting**
- ▶ Tools
 - ▶ Numpy
 - ▶ Scikit-learn

WHAT IS SUPERVISED LEARNING?

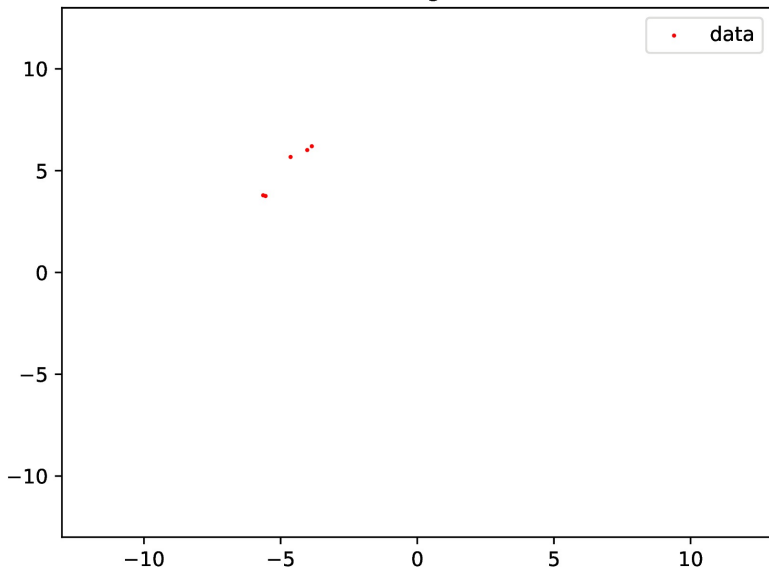
- ▶ Imagine someone gives you a group of smokers
 - ▶ And asks the question – what is their life expectancy?
- ▶ **Completely made up imaginary data**

SOME ABSTRACTION

- ▶ We are given inputs $x_0, x_1 \dots x_n$ and we are looking to predict y
- ▶ Let's plot!

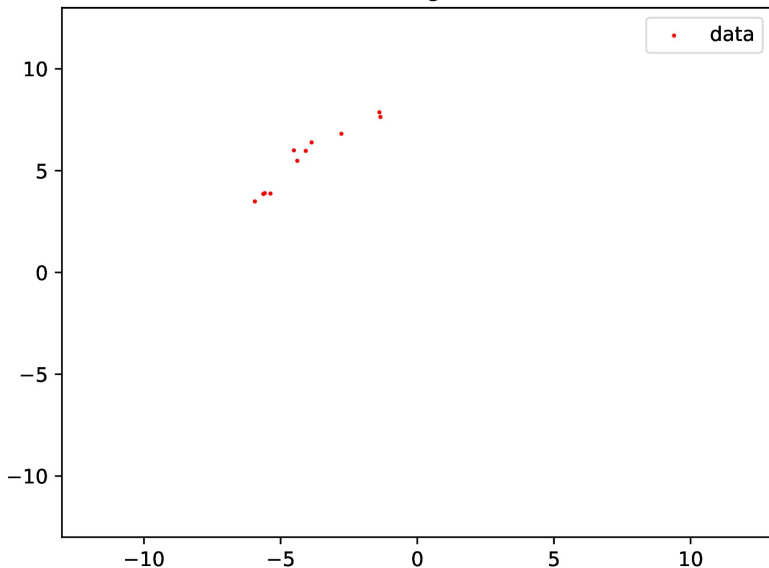
REGRESSION - LINK THE DOTS (1)

training set



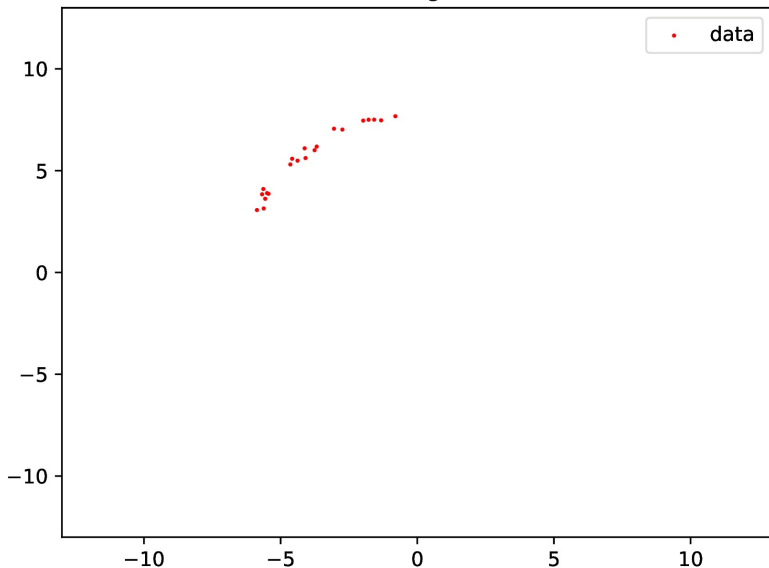
REGRESSION - LINK THE DOTS (2)

training set



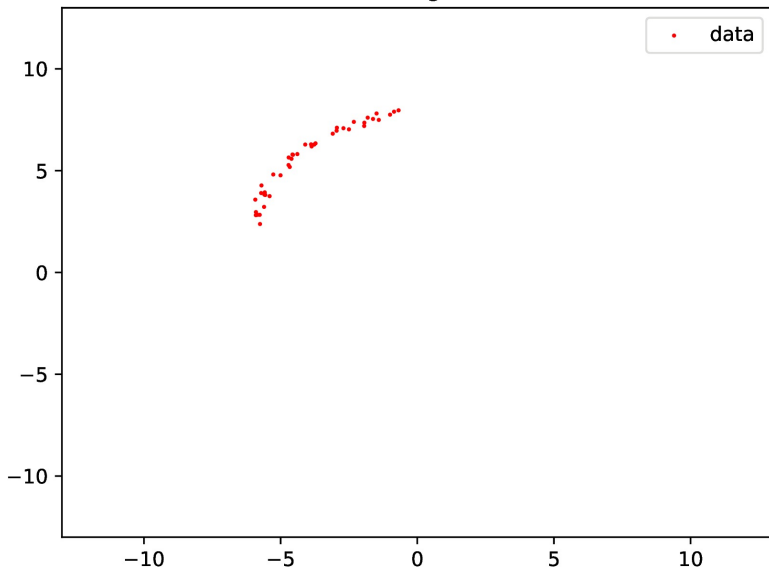
REGRESSION - LINK THE DOTS (3)

training set



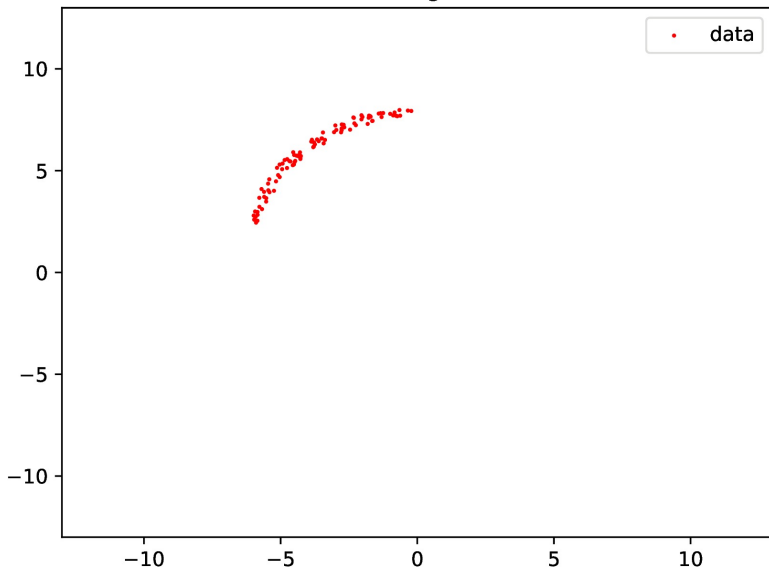
REGRESSION - LINK THE DOTS (4)

training set



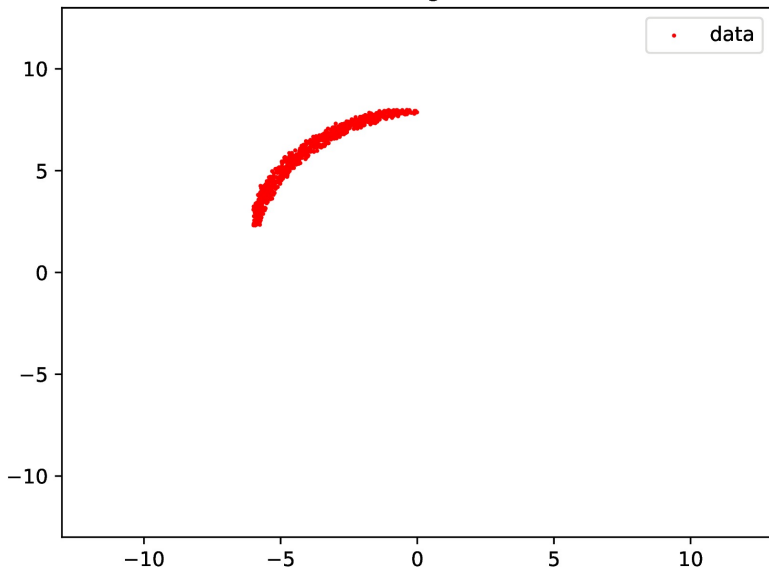
REGRESSION - LINK THE DOTS (5)

training set



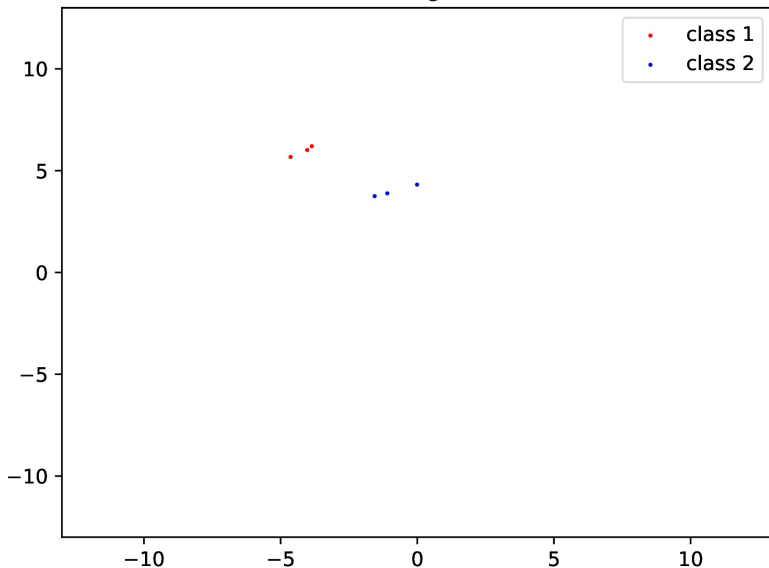
REGRESSION - LINK THE DOTS (6)

training set



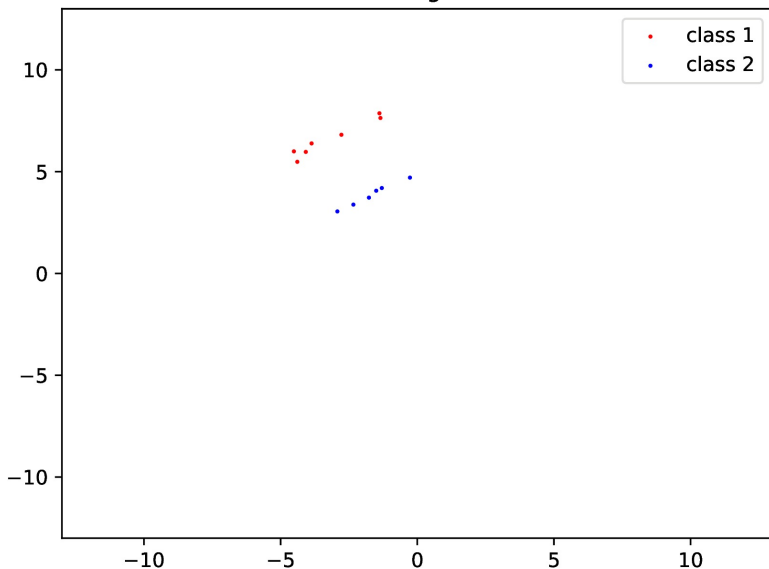
CLASSIFICATION - DRAW A BOUNDARY (1)

training set



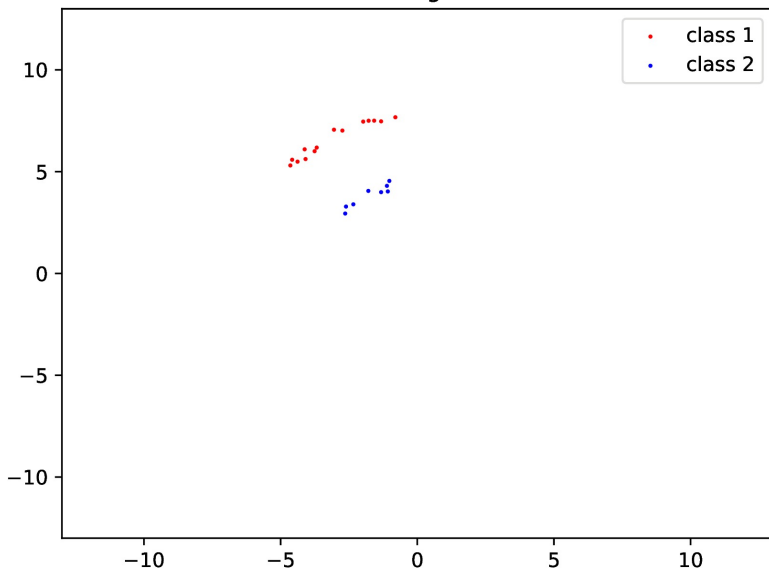
CLASSIFICATION - DRAW A BOUNDARY (2)

training set



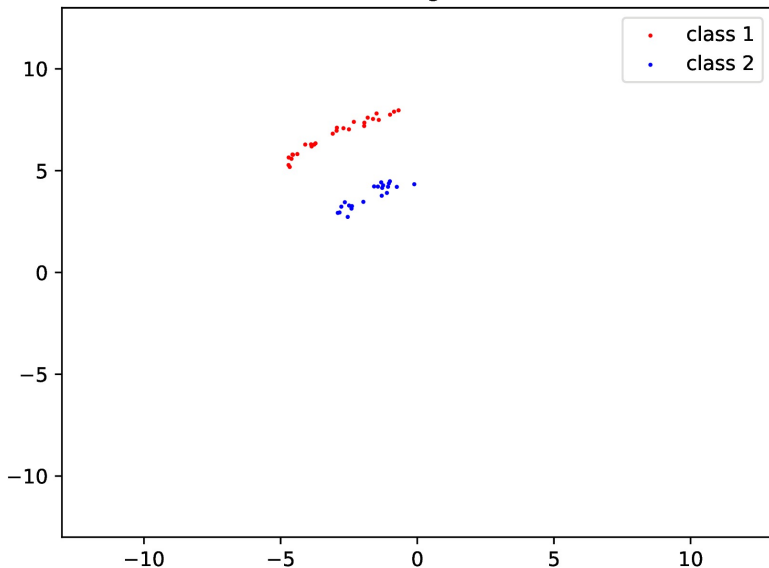
CLASSIFICATION - DRAW A BOUNDARY (3)

training set



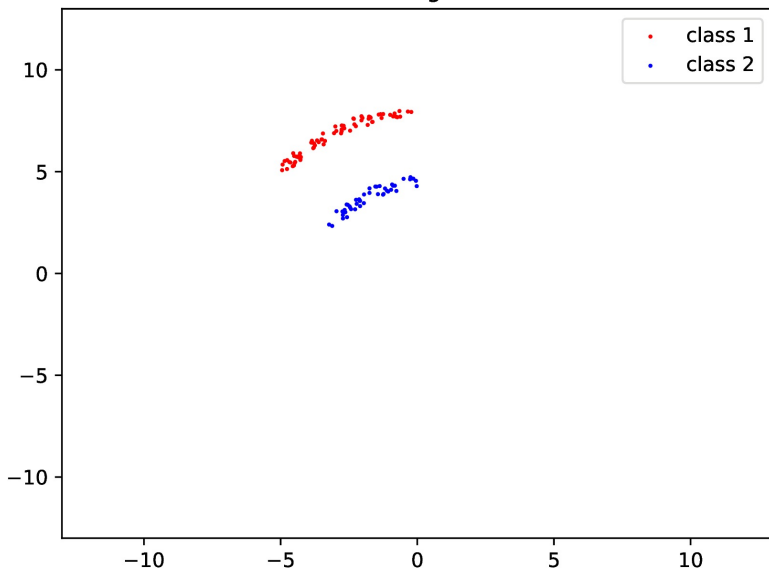
CLASSIFICATION - DRAW A BOUNDARY (4)

training set



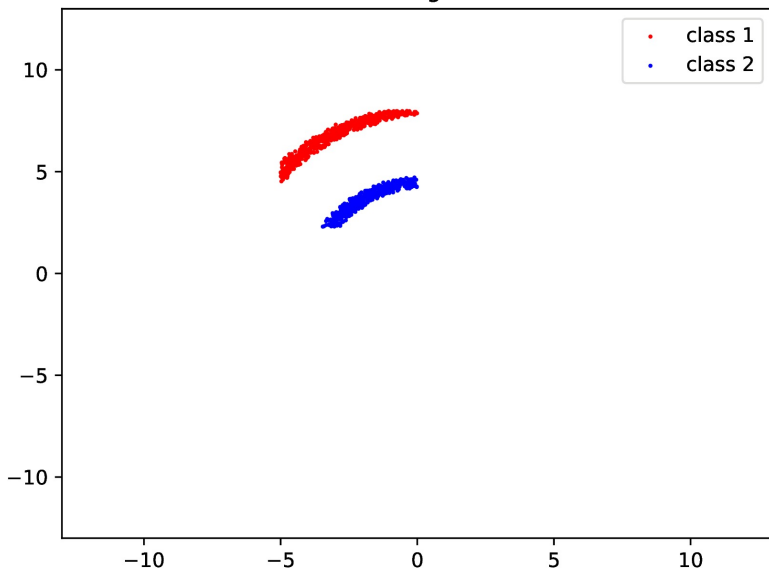
CLASSIFICATION - DRAW A BOUNDARY (5)

training set

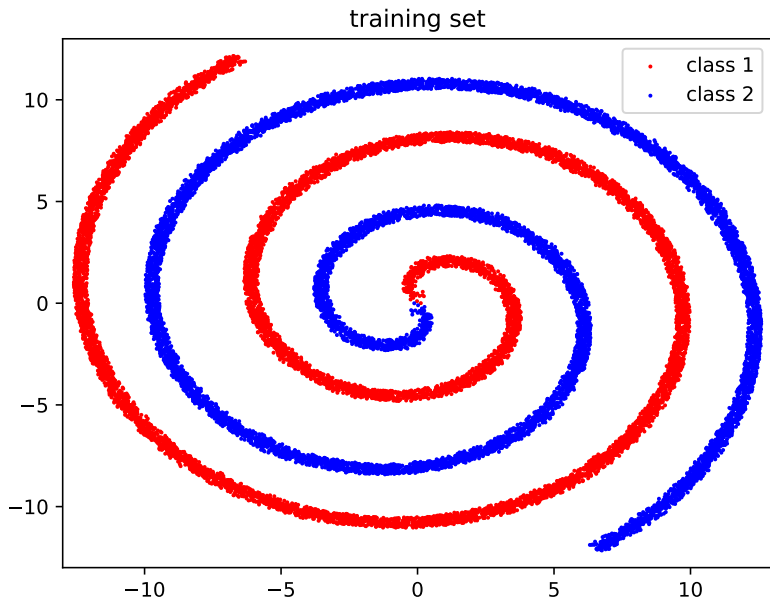


CLASSIFICATION - DRAW A BOUNDARY (6)

training set



FULL DATA



INTUITION

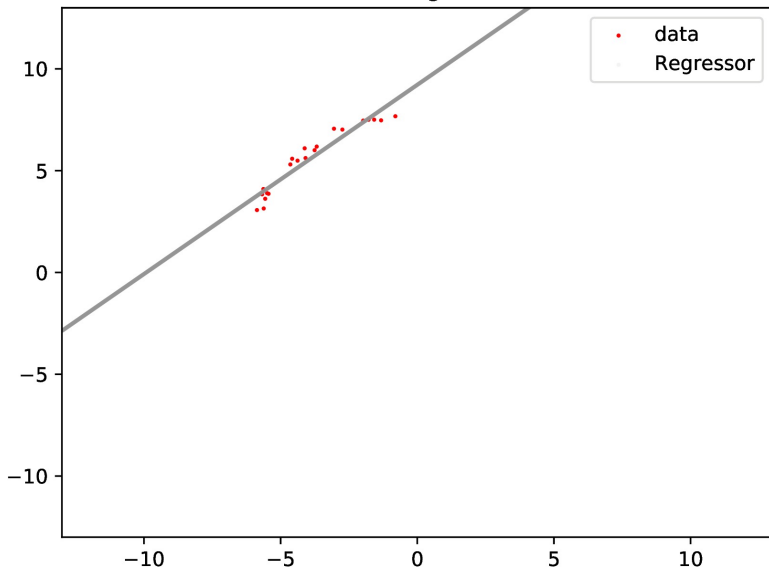
- ▶ That's it - we are given data, and we need to come up with an algorithm to join it up – but in high dimensions
 - ▶ Can be binary, categorical, real-valued
- ▶ How well a function joins the data is called the “loss”
- ▶ Very low loss is not good, it might not generalise that well to unseen data points – you can learn to memorise data instances

LINEAR REGRESSION

- ▶ Linear and logistic regression
 - ▶ Logistic regression does classification
- ▶ You just assume everything is a line

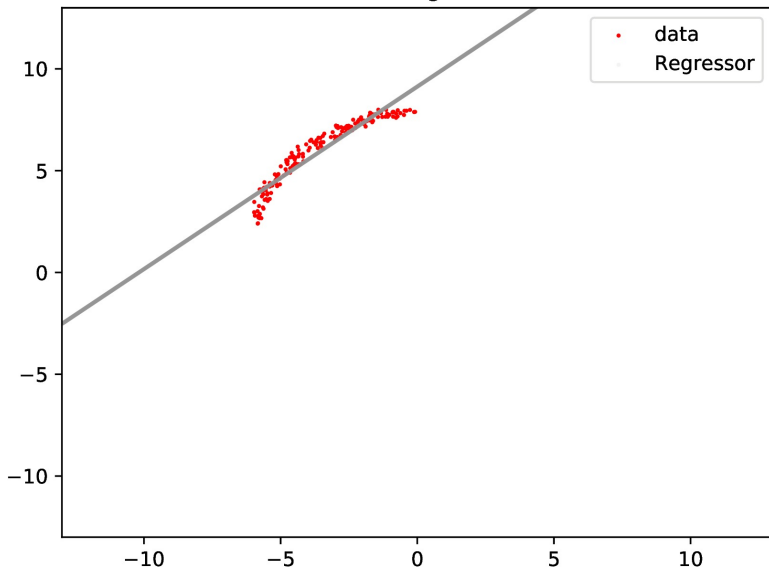
EXAMPLE (LINEAR REGRESSION)

training set



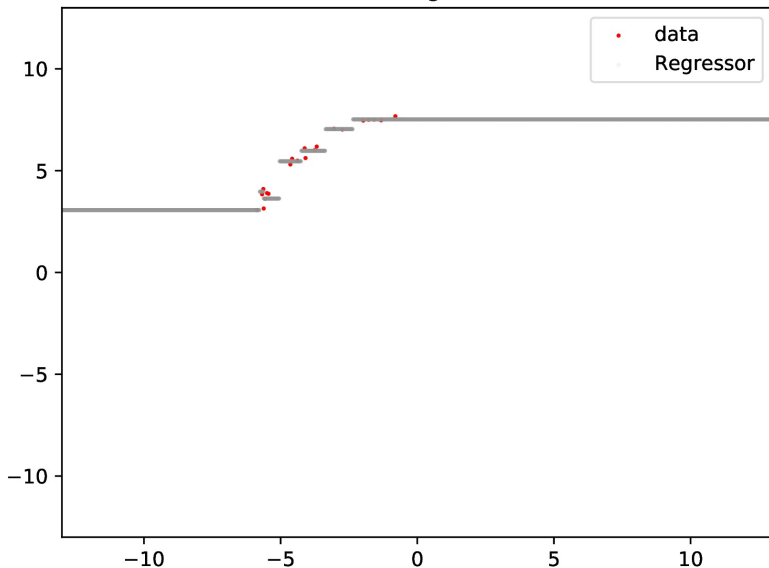
EXAMPLE (LINEAR REGRESSION)

training set



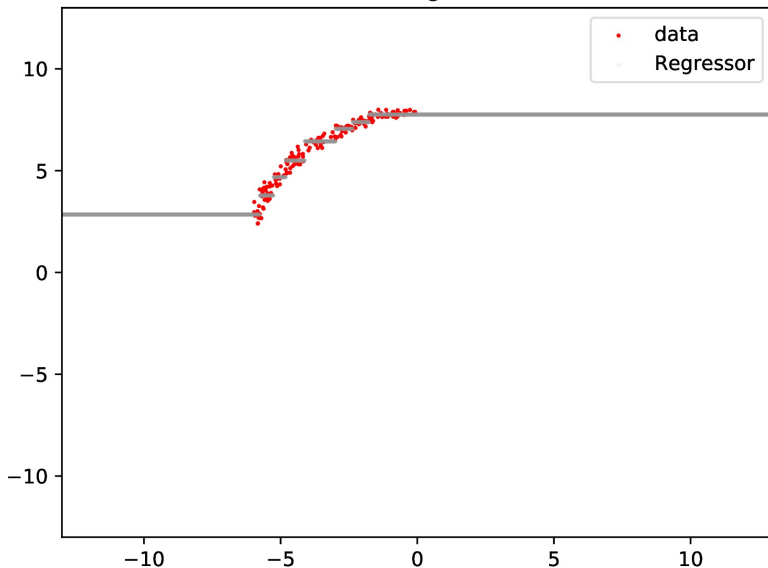
EXAMPLE (DECISION TREE)

training set

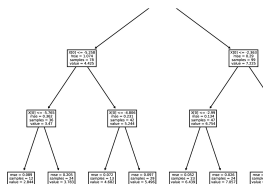


EXAMPLE (DECISION TREE)

training set

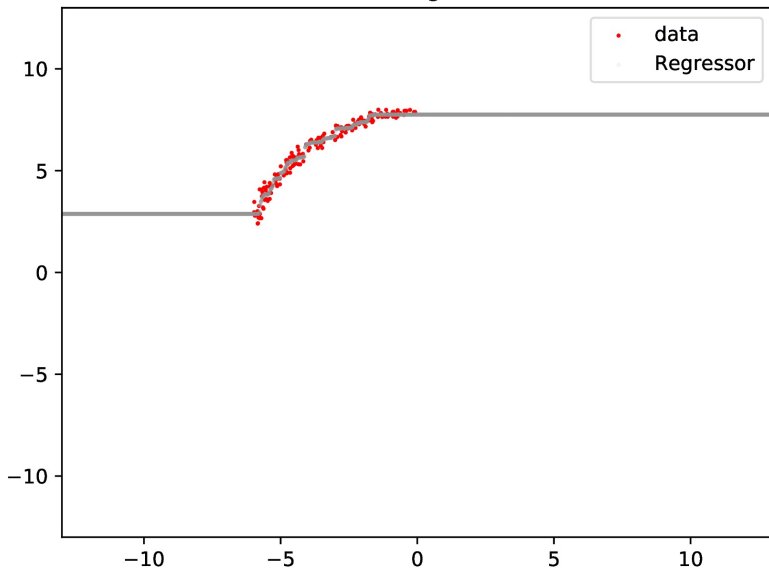


EXAMPLE (DECISION TREE — INTERNAL)



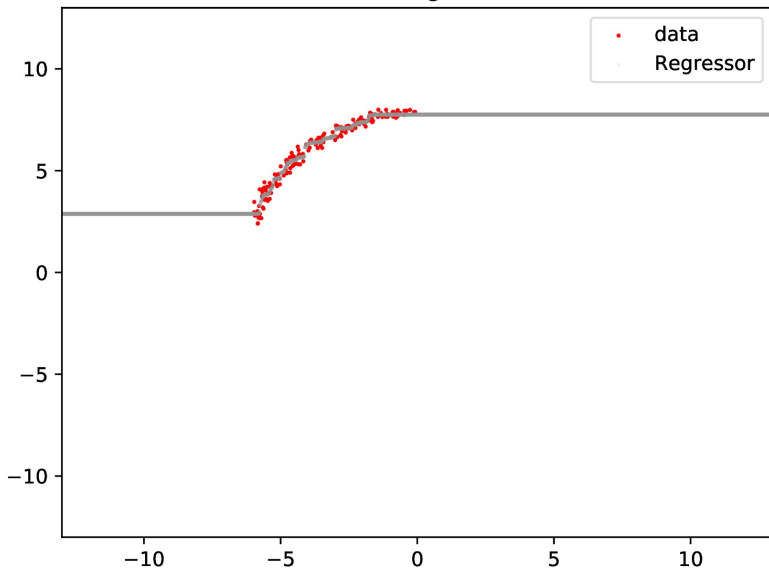
EXAMPLE (RANDOM FOREST)

training set



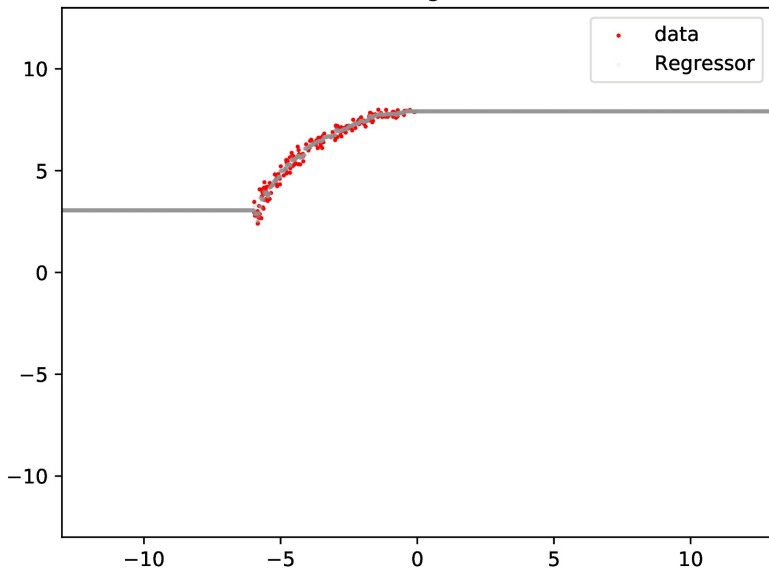
EXAMPLE (RANDOM FOREST)

training set



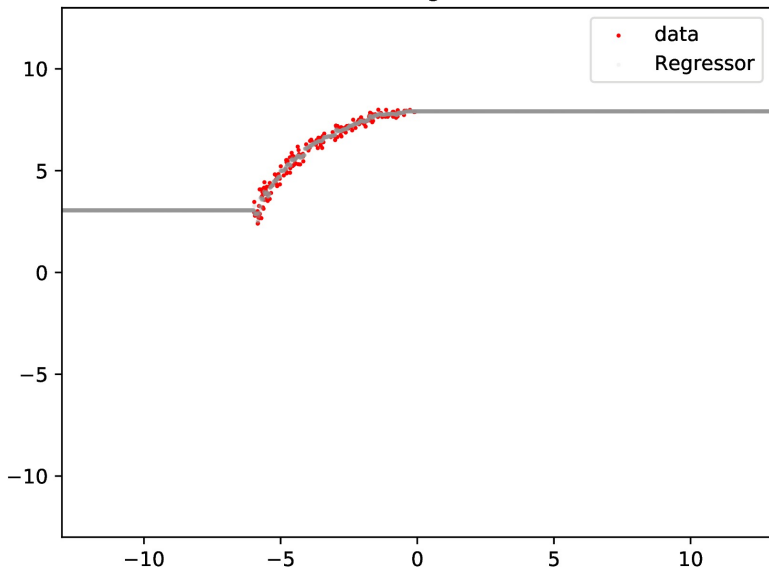
EXAMPLE (RANDOM FOREST)

training set

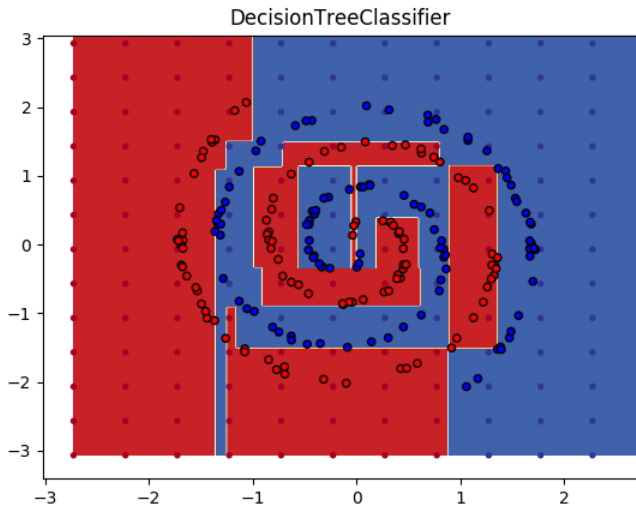


EXAMPLE (GRADIENT BOOSTING)

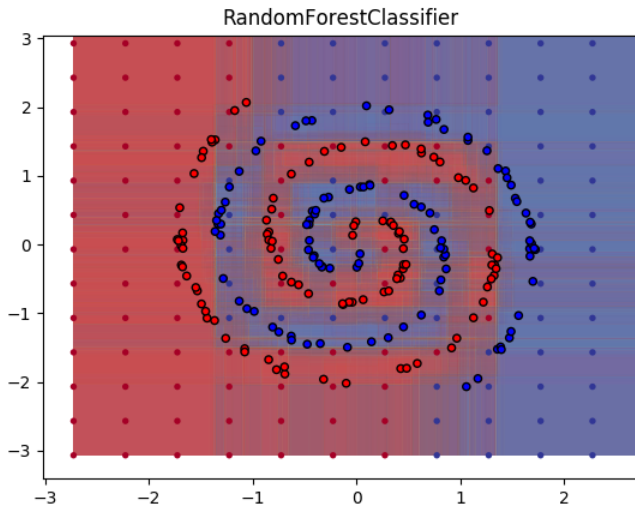
training set



CLASSIFICATION (DECISION TREES)



CLASSIFICATION (RANDOM FORESTS)



UNTIL NOW

BUT HOW DO WE KNOW THIS WILL GENERALISE WELL?

- ▶ Train/Validation/Test split
- ▶ Cross validation

HYPERTPARAMETERS

- ▶ How many trees?
- ▶ Tree depth?
- ▶ l2?