# Welcome to CS-304

## A Data Mining Exploration

# Objectives

- Exploring the facets of machine learning

- Data scientist check list

- Practical applications of converse inductive integrals in the context of epsilon

# A Review

$$S_i^{(t)} = \{x_p : \|x_p - m_i^{(t)}\|^2 \leq \|x_p - m_j^{(t)}\|^2 \; \forall \; 1 \leq j \leq k\},$$

$$m_i^{(t+1)} = \frac{1}{|S_i^{(t)}|} \sum_{x_j \in S_i^{(t)}} x_j$$

?????

# Some Background

# Introduction to Machine Learning

# Introduction to ~~Machine~~ ~~Learning~~

# Introduction to Data Mining

# Depth versus Breadth

# Types of problems with machine learning answers

# Supervised versus Unsupervised

# @bryanl

# @thunderboltlabs

# Required Knowledge

# Math

# Papers

# Persistence

# Regression

# Linear regression
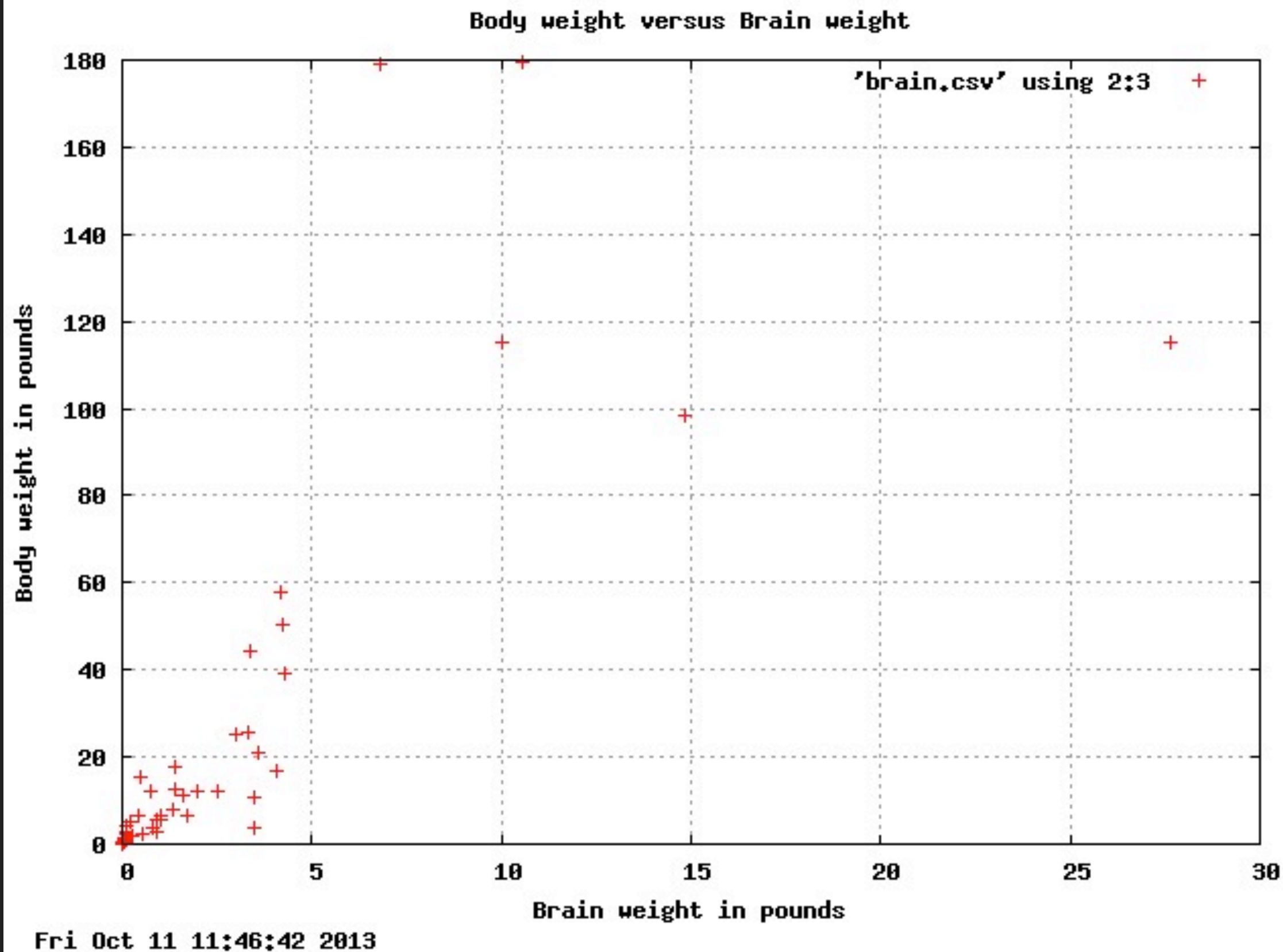
$$y = \alpha + \beta x,$$

$$\text{Find } \min_{\alpha, \beta} Q(\alpha, \beta), \text{ where } Q(\alpha, \beta) = \sum_{i=1}^{n} \hat{\varepsilon}_i^2 = \sum_{i=1}^{n} (y_i - \alpha - \beta x_i)^2$$

$$y = mx + b$$
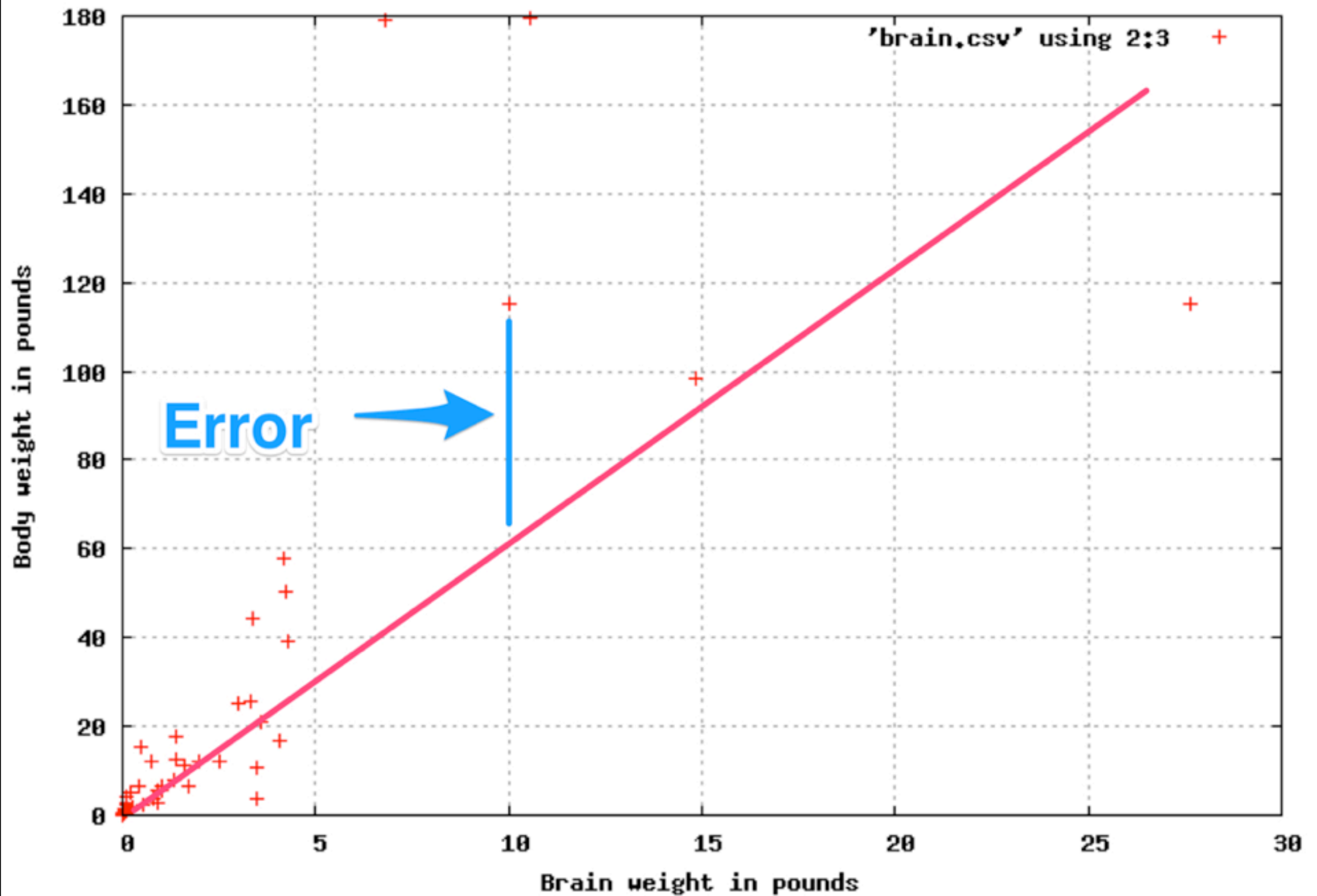
$$y = mx + \alpha$$

$$y = \beta\chi + \alpha$$

```
3,1.350,8.100
4,465.000,423.000
5,36.330,119.500
6,27.660,115.000
7,14.830,98.200
8,1.040,5.500
9,4.190,58.000
10,0.425,6.400
11,0.101,4.000
12,0.920,5.700
13,1.000,6.600
14,0.005,0.140
15,0.060,1.000
16,3.500,10.800
17,2.000,12.300
18,1.700,6.300
```

Body weight versus Brain weight

'brain.csv' using 2:3

Body weight in pounds

Brain weight in pounds

Fri Oct 11 11:46:42 2013

Sunday, October 13, 13

Body weight versus Brain weight

'brain.csv' using 2:3

Error

Body weight in pounds

Brain weight in pounds

Fri Oct 11 11:46:42 2013

Sunday, October 13, 13

$$error=(y_i-\alpha-\beta x_i)$$

$$Q=\sum(error)^2$$

# Classification

# How do we classify?

# Binary classification

# Linear classification

# Support Vector Machines

# Decision trees

# Clustering

# Jaccard Coefficient

# Group documents

# Detect plaguirism

# K-Means Clustering

# Survey of the Ruby Landscape

# AI4R

# SciRuby

# JRuby and Mahout

# Rails on Ruby

# Fast Math

# Easy Plotting

# Integrated Environment

# Let me dance!

# Moving Forward

# Want to learn more?

# Linear Algebra

# Calculus

# Coursera ML

# Wikipedia

# Now, if you want to get serious

# 1. Find a dataset

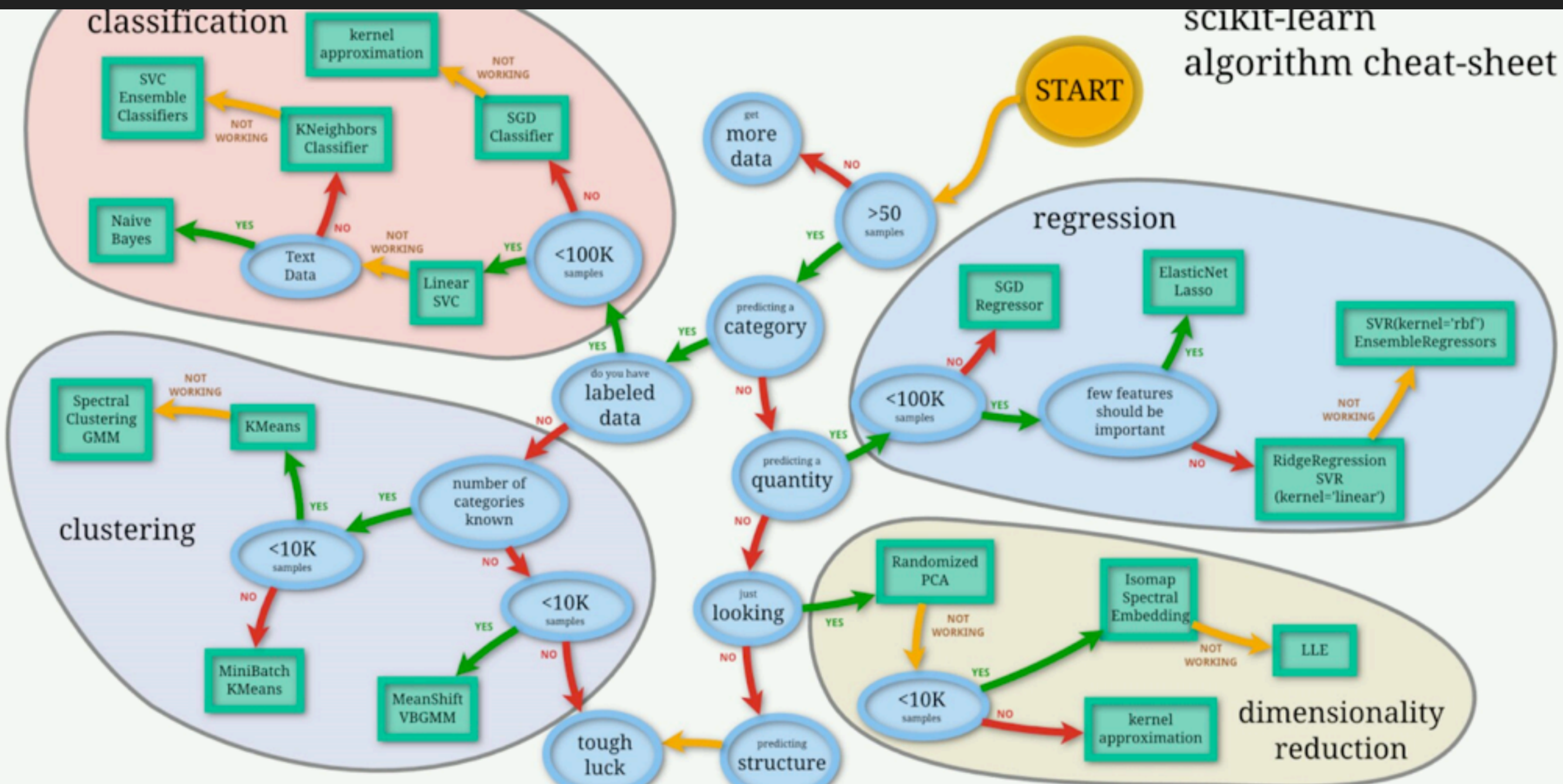# 2. Find another language

# 3. ...

# 4. Profit?

# We haven't event scratched the surface

scikit-learn algorithm cheat-sheet

# BigML

# Dundas

# Kaggle

# Python Land

# Mahout

# Shark with Spark