# DATA SCIENCE
## 11 WEEK PART TIME COURSE

## Week 3 – Regression
## Monday 4th April

1. Motivation
2. Supervised Vs Unsupervised learning
3. What is Linear Regression?
4. How do Run a Linear Regression Model?
5. Lab
6. Discussion / Review / Homework

# WHAT ARE THE GOALS OF STATISTICAL LEARNING?

# scikit-learn algorithm cheat-sheet

**START**

## classification
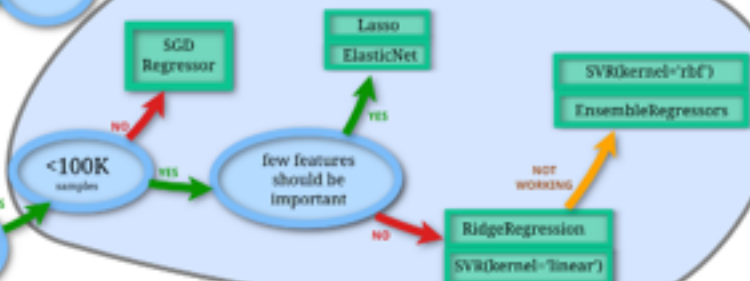
- kernel approximation
- SVC
- Ensemble Classifiers
- KNeighbors Classifier
- SGD Classifier
- Naive Bayes
- Text Data
- Linear SVC
- <100K samples

get more data

>50 samples — NO → get more data

>50 samples — YES → predicting a category

predicting a category — YES → do you have labeled data

do you have labeled data — YES → <100K samples (classification)

<100K samples — YES → Text Data
<100K samples — NO → SGD Classifier

Text Data — YES → Naive Bayes
Text Data — NO → KNeighbors Classifier

KNeighbors Classifier — NOT WORKING → SVC / Ensemble Classifiers
Linear SVC — NOT WORKING → KNeighbors Classifier
SGD Classifier — NOT WORKING → kernel approximation

## regression

- SGD Regressor
- Lasso / ElasticNet
- SVR(kernel='rbf') / EnsembleRegressors
- few features should be important
- RidgeRegression / SVR(kernel='linear')
- <100K samples

predicting a category — NO → predicting a quantity

predicting a quantity — YES → <100K samples (regression)

<100K samples — NO → SGD Regressor
<100K samples — YES → few features should be important

few features should be important — YES → Lasso / ElasticNet
few features should be important — NO → RidgeRegression / SVR(kernel='linear')

RidgeRegression / SVR(kernel='linear') — NOT WORKING → SVR(kernel='rbf') / EnsembleRegressors

## clustering

- Spectral Clustering / GMM
- KMeans
- number of categories known
- <10K samples
- MiniBatch KMeans
- MeanShift / VBGMM

do you have labeled data — NO → number of categories known

number of categories known — YES → <10K samples
number of categories known — NO → <10K samples (MeanShift)

<10K samples — YES → KMeans
<10K samples — NO → MiniBatch KMeans

KMeans — NOT WORKING → Spectral Clustering / GMM

<10K samples — YES → MeanShift / VBGMM
<10K samples — NO → tough luck

## dimensionality reduction

- Randomized PCA
- Isomap / Spectral Embedding
- LLE
- <10K samples
- kernel approximation

predicting a quantity — NO → just looking

just looking — YES → Randomized PCA
just looking — NO → predicting structure

Randomized PCA — NOT WORKING → <10K samples

<10K samples — YES → Isomap / Spectral Embedding
<10K samples — NO → kernel approximation

Isomap / Spectral Embedding — NOT WORKING → LLE

predicting structure → tough luck

**tough luck**

Back

scikit learn

We want to predict some value, let's call it **y**, based on some observed data we have, let's call that **x**.

We will use statistical learning to estimate a function that approximates **y** based on the input, **x**.

**y** is also called; label, dependent variable, target

**x** is also called; predictor, independent variable, features

We want to predict the price of a house, let's call it y, based on some observed data we have about the area, number of bedrooms, size of the house, and if it has a pool or not.

The area, number of bedrooms, size of the house, and if it has a pool or not would be our x variables (sometimes you might see this denoted as X)

What we want is y = f(X), a way to describe the house price based on observed data

If the y variable is numeric then we have a regression problem - we are trying to predict a continuous number

If the y variable is a category (for example trying to predict a type of flower) the we have a classification problem - we are trying to classify what group that y belongs to.

We want to find some underlying structure or patterns in the data but in this case we don't have any labeled data.

So for example, if we have a large group of customers but would like to separate them into groups (or clusters) to better target them.

# WHAT IS LINEAR REGRESSION?

We want to model a linear relationship (think straight line) between our target variable y and our input variable x.

$$y = X\beta + \epsilon$$

$$y = X\beta + \epsilon$$

‣ y = target variable

‣ X = input variable

‣ β = coefficients

‣ ∊ = error term

Note, one of our input variables can be 1 so we have an intercept parameter

‣ Linear relationship in the parameters, β, we can transform the actual values of the inputs if we want

‣ Variance of the error term, $\epsilon$, is constant. This means there is no systematic pattern in the values of X and the variance of $\epsilon$

‣ The mean of $\epsilon = 0$

‣ $\epsilon$ has a normal distribution

‣ No perfect (or near perfect) co-linearity between any of the input variables. Otherwise the fitting procedure will break.

# HOW TO RUN LINEAR REGRESSION?

$$SS_{res} = \sum_{i=1}^{n}(y_i - f(x_i))^2$$

Basically, what we are trying to do is minimise the Residual Sum of Squares. This is the Sum of the squared difference between our observed value and the value from the model

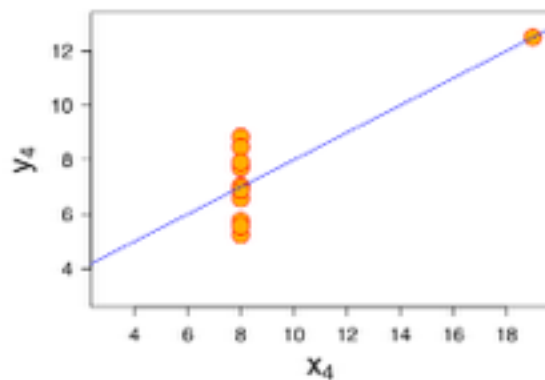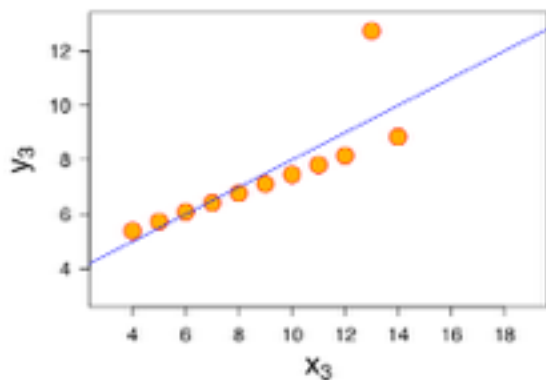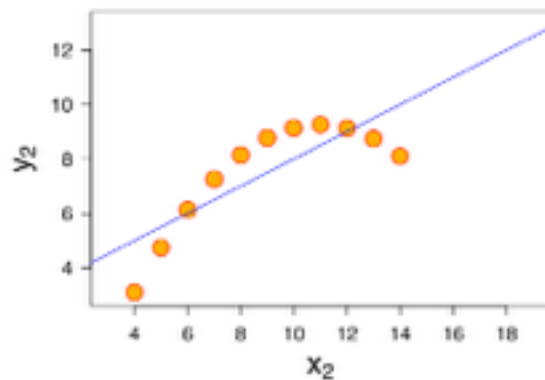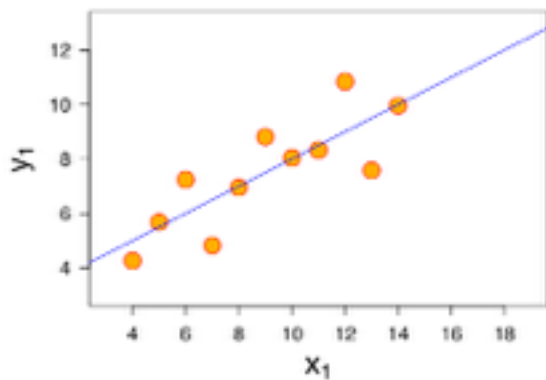$$SS_{res} = \sum_{i=1}^{n} (y_i - f(x_i))^2$$

Basically, what we are trying to do is minimise the Residual Sum of Squares. This is the Sum of the squared difference between our observed value and the value from the model

Oak Diameter vs. Age

‣ Make sure you visualise your data and check the actual model fit !!!

‣ The fitting a model to the four datasets in the table on the right produce the same fit statistics, model coefficients and standard error

‣ See anything wrong?

**Anscombe's quartet**

|  | I |  | II |  | III |  | IV |
|---|---|---|---|---|---|---|---|
| x | y | x | y | x | y | x | y |
| 10.0 | 8.04 | 10.0 | 9.14 | 10.0 | 7.46 | 8.0 | 6.58 |
| 8.0 | 6.95 | 8.0 | 8.14 | 8.0 | 6.77 | 8.0 | 5.76 |
| 13.0 | 7.58 | 13.0 | 8.74 | 13.0 | 12.74 | 8.0 | 7.71 |
| 9.0 | 8.81 | 9.0 | 8.77 | 9.0 | 7.11 | 8.0 | 8.84 |
| 11.0 | 8.33 | 11.0 | 9.26 | 11.0 | 7.81 | 8.0 | 8.47 |
| 14.0 | 9.96 | 14.0 | 8.10 | 14.0 | 8.84 | 8.0 | 7.04 |
| 6.0 | 7.24 | 6.0 | 6.13 | 6.0 | 6.08 | 8.0 | 5.25 |
| 4.0 | 4.26 | 4.0 | 3.10 | 4.0 | 5.39 | 19.0 | 12.50 |
| 12.0 | 10.84 | 12.0 | 9.13 | 12.0 | 8.15 | 8.0 | 5.56 |
| 7.0 | 4.82 | 7.0 | 7.26 | 7.0 | 6.42 | 8.0 | 7.91 |
| 5.0 | 5.68 | 5.0 | 4.74 | 5.0 | 5.73 | 8.0 | 6.89 |

# LAB

git remote -v

git remote add upstream https://github.com/ihansel/SYD_DAT_3.git

git remote -v

git fetch upstream

git checkout master

git merge upstream/master

# DISCUSSION TIME

- ‣ Review of last class
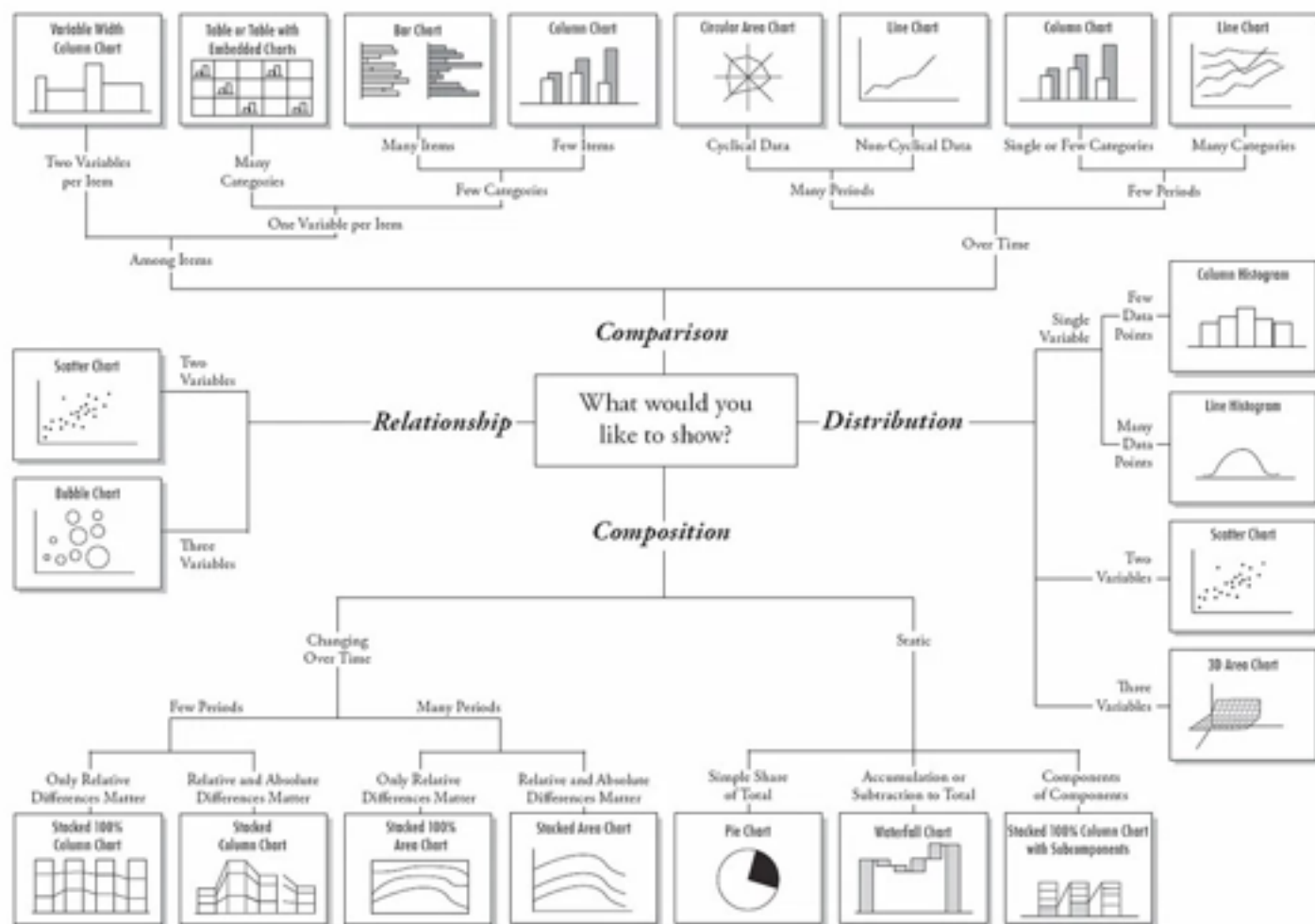- ‣ Further Reading for Regression
- ‣ Check in with homework/course project

Week 2 Monday 14th

☑ Understand goals of Data Viz.

☑ Visualise a dataset

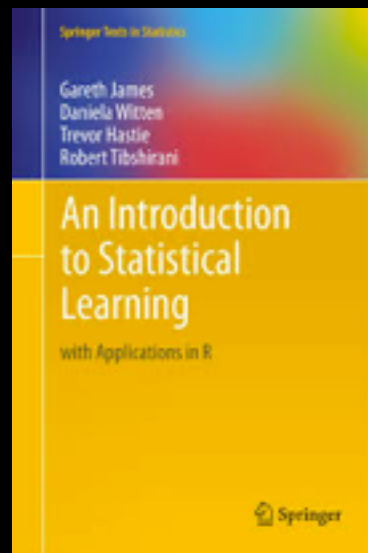☑ Understand 3 different graph types

☑ Examples & Sources to Review

# Chart Suggestions—A Thought-Starter



**Variable Width Column Chart** — Two Variables per Item

**Table or Table with Embedded Charts** — Many Categories

**Bar Chart** — Many Items

**Column Chart** — Few Items

Few Categories

One Variable per Item

Among Items

**Circular Area Chart** — Cyclical Data

**Line Chart** — Non-Cyclical Data

Many Periods

**Column Chart** — Single or Few Categories

**Line Chart** — Many Categories

Few Periods

Over Time

**Comparison**

**Scatter Chart** — Two Variables

**Bubble Chart** — Three Variables

**Relationship**

## What would you like to show?

**Distribution**

**Composition**

**Column Histogram** — Single Variable, Few Data Points

**Line Histogram** — Many Data Points

**Scatter Chart** — Two Variables

**3D Area Chart** — Three Variables

Changing Over Time

Static

Few Periods — Only Relative Differences Matter: **Stacked 100% Column Chart**; Relative and Absolute Differences Matter: **Stacked Column Chart**

Many Periods — Only Relative Differences Matter: **Stacked 100% Area Chart**; Relative and Absolute Differences Matter: **Stacked Area Chart**

Simple Share of Total: **Pie Chart**

Accumulation or Subtraction to Total: **Waterfall Chart**

Components of Components: **Stacked 100% Column Chart with Subcomponents**

# DISCUSSION TIME

**An Introduction to Statistical Learning**

‣ **Chapter 3 – Linear Regression**

‣ **Chapter 6 – Linear Model Selection and Regularization**

# DISCUSSION TIME

**Homework/Course Project**

‣ **How's it going ?**