

MSc Data Mining

Topic 04 : Regression

Part 01 : Overview

Dr Bernard Butler and Dr Kieran Murphy

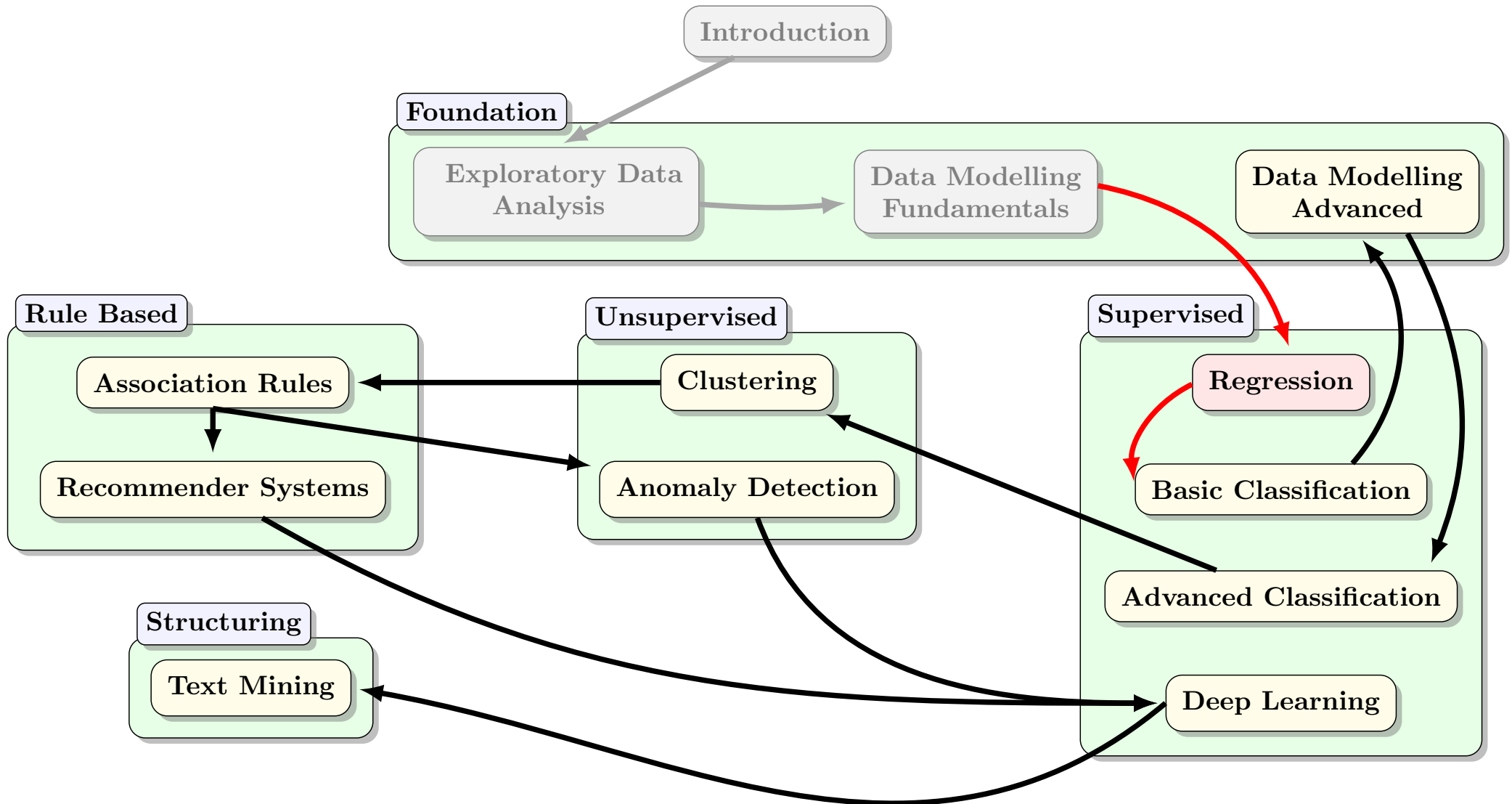
Department of Computing and Mathematics, WIT.
(bernard.butler@setu.ie; kmurphy@wit.ie)

Spring Semester, 2023

Outline

- Regression as a means of minimising sum of the squared errors
- Regression assumptions - what they mean, how they can be used for validation and model building
- Case studies from Advertising, Diamond sales, Credit Balance prediction

Data Mining (Week 4)



Overview — Summary

1. Introduction
2. Linear regression assumptions
3. Reviewing regression results
4. Case Study 1: Generated
5. Case Study 2: Diamonds
6. Case Study 3: Advertising
7. Case Study 4: Credit Balances
8. Diagnostics and Plots - how to fix problems
9. Resources

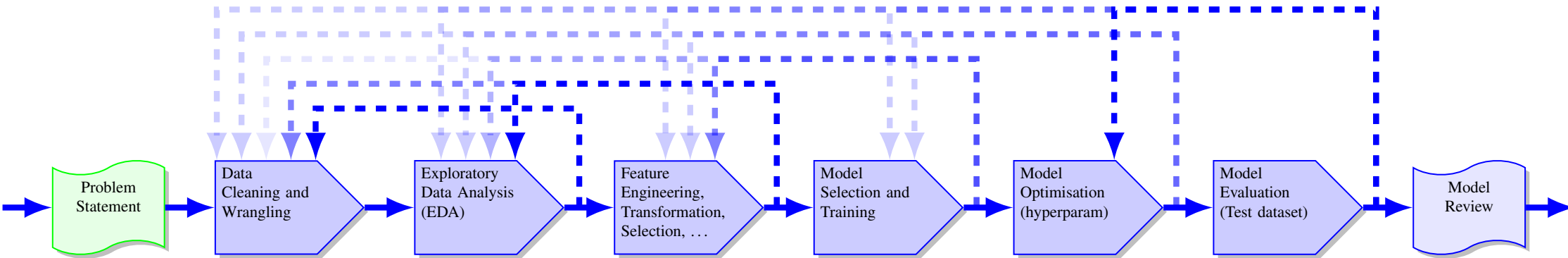
This Week's Aim

This week's aim is to give an overview of linear regression: fitting linear models to data, to **predict a numeric value**.

- High level view of regression: where it came from, what it attempts to do.
- Examine some extensions to the simplest case of linear regression.
- Consider how to check that the regression was successful, and make some improvements if necessary
- To provide context we will use the following datasets:
 - Generated data (various)
 - Diamond dataset: predicting diamond prices given their weights
 - Advertising dataset: predicting widgets sold based on spending in different advertising channels
 - Credit dataset: predicting credit balance using income, status, etc.

If you need **to predict a categorical value** (such as whether a passenger survived the sinking of the Titanic), that is classification, not regression.

Data Mining Workflow: where Regression Modelling fits



- Data pre-processing (Data Cleaning + EDA + (Basic) Feature Engineering) is the foundation of Data Mining.
- Train-test split, feature engineering and reducing errors (covered by Kieran last week) play key roles in Regression.
- Data Models are needed for prediction: learning from training data, so that predictions can be made on new data.
- Regression models have common characteristics and assumptions - we consider them today.
- Regression is the process of building, improving, validating and visualising such models.

But first - where did the term come from?

“Regression to the Mean”

- Sir Francis Galton (1822-1911): polymath, genius, eugenicist
- Found that children of tall parents “regressed” towards average height of population
 - Let δh be the height difference between their parents and the general population.
 - Then the height difference of the children versus their peers is predicted to be $\frac{2}{3}\delta h$.
 - So: the height of children can be predicted using a simple formula (mathematical expression).
- Techniques used, and underlying statistics, were applicable to other problems, so the name stuck.
- Regression has been expanded to handle larger and more complex data.
- Galton “fitted” the formula by eye, but Legendre and Gauss previously (early 1800s) invented *least squares* for predicting orbits in astronomy.
- Used to require hours of tedious calculation, but not any more!

Simple Linear Regression: Background

- With small data sets, calculations can be done by hand, but they are tedious and error-prone.
- The goal is simple: Given a (training) set of (x, y) data where y is assumed to have a linear relationship with x
 - Find the line that is the “best fit” to that data
 - Use the specification of that line to *predict* y for the (test) x values
- Note that the “linear relationship” of y upon x is just one of the underlying assumptions
- In practice, the data does not have an exact linear relationship, but it should be “close enough”—but what does that mean?
- In terms of Week 3’s **ML models taxonomy**: regression is **geometric** and **parametric**
- In terms of Week 3’s **Components of a Machine Learning Problem**
 - **Representation** is based on (fitting) hyperplanes to point clouds
 - **Evaluation** is based on a MSE metric, with assumption checks to help identify the best model family
 - **Optimization** is one-step (no search needed) because we have a constraint on the errors we allow
- Hyperparameter tuning: which features to include (e.g., polynomial terms).

Review: Linear combinations (scalar product)

Definition 1

Given two vectors **a** and **b**, each with n elements, the *linear combination* (c) of **a** and **b** is

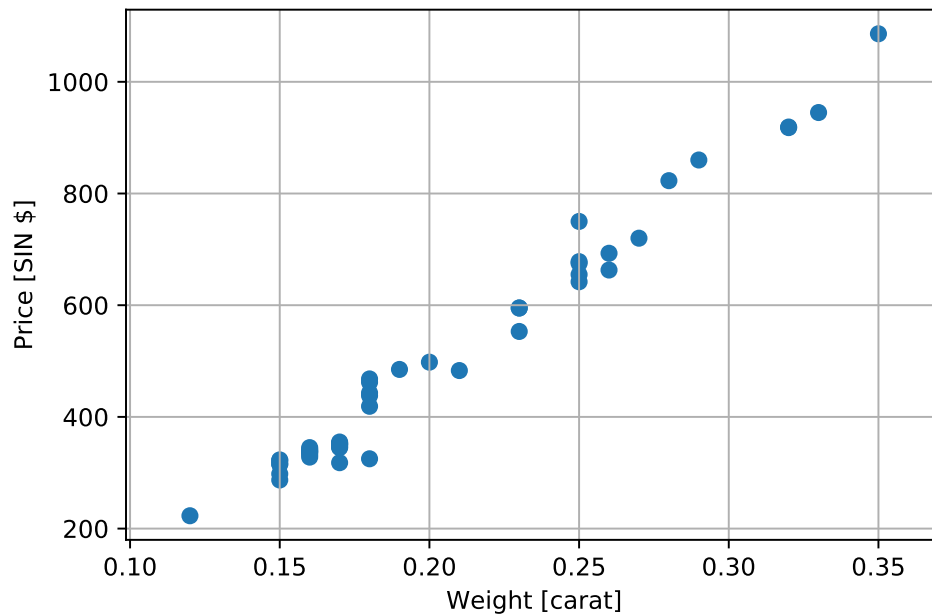
$$c \equiv a_1b_1 + a_2b_2 + \dots + a_nb_n = \sum_{i=1}^n a_ib_i \equiv |\mathbf{a}||\mathbf{b}| \cos(\mathbf{a}, \mathbf{b})$$

Remarks

- The linear combination of 2 vectors is a scalar, which can be seen as “mixing” two vectors.
- Matrix-vector multiplication $A\mathbf{x}$ can be seen as the linear combination of each row in the matrix A with the (column) vector \mathbf{x} .
- Matrix-matrix multiplication AB can be seen as the linear combination of each row in the matrix A with each column in the matrix B .
- Two vectors **a** and **b** can have a scalar product that is zero even if neither $|a|$ nor $|b|$ is zero, but only if $\cos(\mathbf{a}, \mathbf{b}) = 0$, i.e., the **a** and **b** vectors are perpendicular to each other.
- **Linear combinations are used for prediction from linear (regression) models.**

Motivating example: Diamond data

Relation between diamonds' price and weight



Diamond Prices by Weight

- Given the data on the left, can we use it to predict the price of a diamond that weighs 0.22 carat?
- NB - we have not seen a diamond with that weight before in the data
- Can you think of at least 3 other factors that might affect the price?
- Various(!) - some examples: clarity, cut, provenance, part of a set, ...

Simple Linear Regression: Formulation

Definition 2 (Matrix formulation)

- Given data $\{x_i, y_i\}$ where $i = 2, 3, \dots, n$ and β_0, β_1 as the (unknown, but to be determined) *intercept* and *slope* of the regression line for this data.
- For $n = 2$ points with $x_2 \neq x_1$, this can be solved uniquely for β_0, β_1 , using techniques you learnt for your Junior/Inter Cert.
- For $n > 2$ *collinear* points, just pick any two points on the line and solve as before.
- Otherwise you need a more general formulation using matrices, and can use linear algebra to solve for β_0, β_1 .
- General equation is $y_i = \beta_0 + \beta_1 x_i + \epsilon_i = \hat{y}_i + \epsilon_i$ (data = model + error), where \hat{y} is the predicted y for these values of β_0, β_1 .
- Matrix form is $\mathbf{y} = X\beta$. Remember matrix-vector multiplication: inner product of i^{th} row of X times the vector $\beta = 1 \times \beta_0 + x_i \times \beta_1 = \hat{y}_i$.
- However, we don't know β yet, nor do we know \hat{y}_i , so we use y_i as an estimate of \hat{y}_i and solve for all data in the training set.
- So: our task is to solve the *overdetermined* (number of rows exceeds the number of columns) system of equations $\mathbf{y} = X\beta$ for β

Simple Linear Regression: Normal Equations

$$\begin{aligned}\mathbf{y} &\approx X\beta \\ \mathbf{y} &= X\beta + \epsilon \\ X^T \mathbf{y} &= X^T X \beta + X^T \epsilon \\ X^T \mathbf{y} &= X^T X \beta\end{aligned}$$

because we require $X^T \epsilon \equiv 0$ for a least-squares fit*. Swapping sides, we have

$$\begin{aligned}(X^T X) \beta &= X^T \mathbf{y} \\ (X^T X)^{-1} (X^T X) \beta &= (X^T X)^{-1} X^T \mathbf{y}\end{aligned}$$

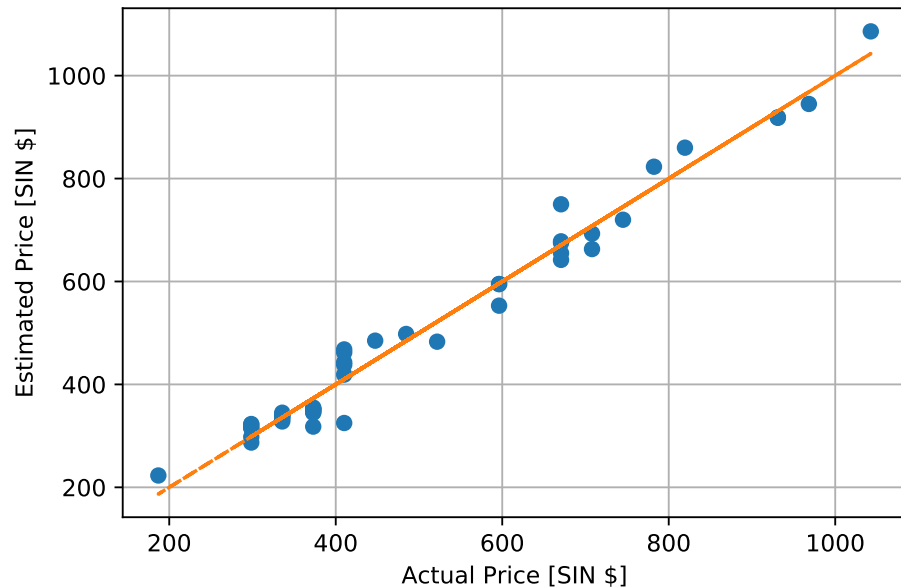
which is equivalent to the *Normal equations*

$$\beta = (X^T X)^{-1} X^T \mathbf{y} \tag{1}$$

Note that everything on the right is a set of operations on the training data.

Simple Linear Regression: Balanced Errors

Relation between estimated and actual diamonds' prices



- More generally: a weighted sum of the errors should be 0.
- Weights should depend on the features.
- The $X^T \epsilon = 0$ criterion works well, so we apply it.

What makes this look like a good fit?

*The fitted line passes through the data centroid and errors pass are **balanced** - cf. see-saw*

Simple Linear Regression: Implementation

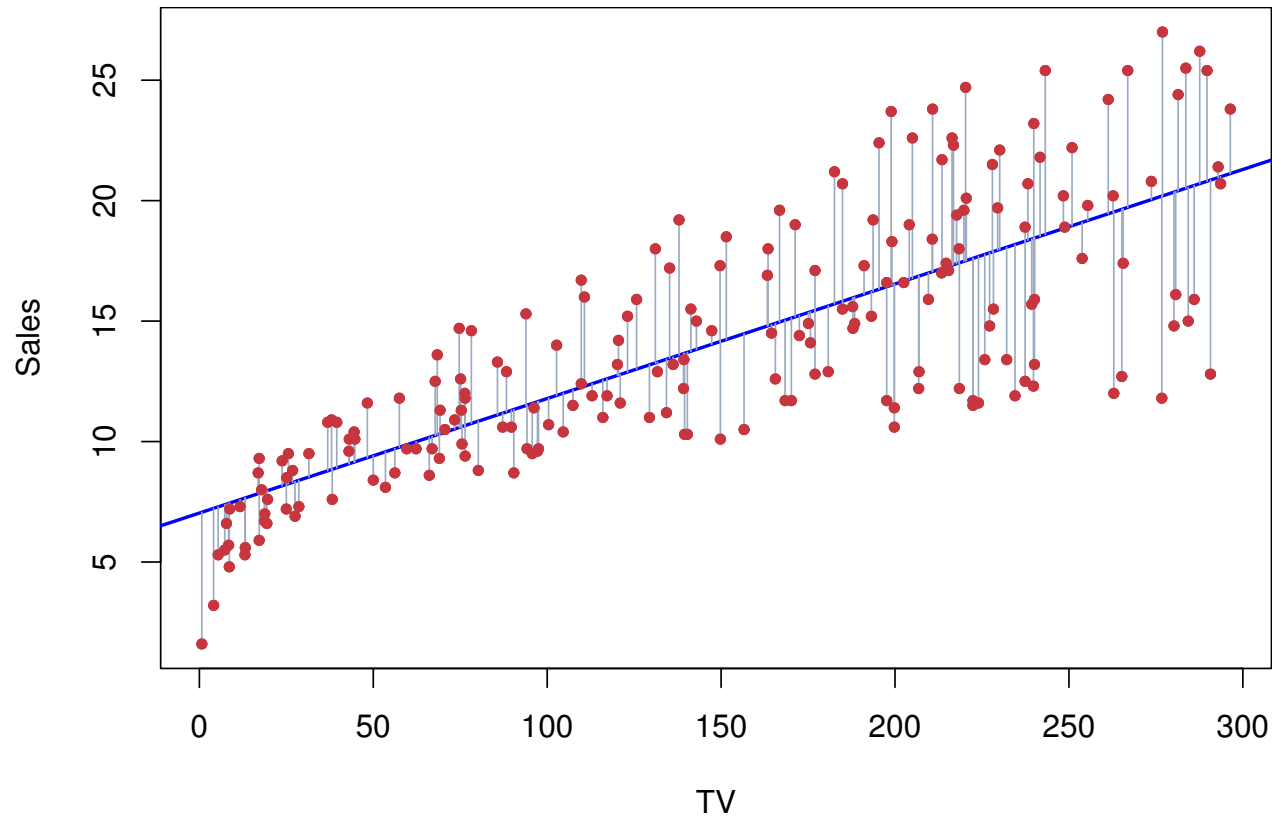
When implemented in software, the Normal equations are not used directly: better (in the sense of being faster and more numerically accurate!) algorithms are used instead, but the results are equivalent in exact arithmetic (remember: digital computers perform finite-precision arithmetic and so cannot be exact).

One option is to use statsmodels: consistent with R (separate model specification), excellent diagnostics as standard

Another option is to use sklearn: consistent with other sklearn algorithms, more controls

Remember: after *learning* the β parameters, it is then possible to predict \hat{y} for “new” x values, using separate *prediction* function calls.

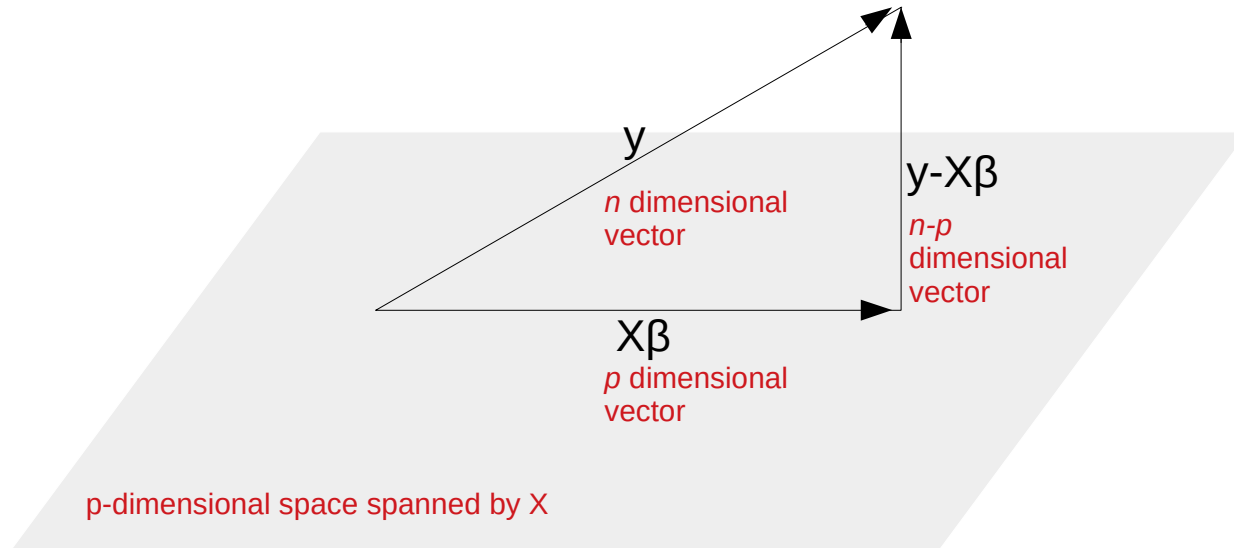
SLR: Residual Plot for the model



Source: ISLR, Fig 3.1: Advertising data with the model “ $\text{Sales} \sim \text{TV}$ ”.

Note the vertical distance between the red dots (data points) y and the corresponding \hat{y} on the regression line, which is termed the *error* ϵ .

Geometrical interpretation of regression



- The X matrix spans the $p \times p$ space represented by the grey plane.
- But \mathbf{y} has $n > p$ dimensions and so is represented by a point that lies outside the grey plane.
- When \mathbf{y} is projected onto the nearest point in the X space,
 - The projected point is $\hat{\mathbf{y}}$.
 - The residuals (errors) ϵ are $\mathbf{y} - X\beta \equiv \mathbf{y} - \hat{\mathbf{y}}$.

This decomposition of n data dimensions (observations) into p model parameters and n residuals with rank $n - p$ is helpful when interpreting regression diagnostics.

OLS and Linear Regression

Definition 3 (BLUE)

According to the Gauss-Markov theorem, *Ordinary Least Squares* (OLS), which uses the Normal equations to minimise the sum of the squares of the errors ($\|\epsilon\|_2 \equiv \sqrt{e_1^2 + e_2^2 + e_3^2 + \dots + e_n^2}$), is the *Best, Linear, Unbiased, Estimator* of that model that can be derived from the training data, provided some reasonably loose assumptions hold.

When we discuss Bias, Variance and Irreducible Error, it is clear that low bias is not enough. OLS might be BLUE but that does not guarantee low variance, because overfitting can still be a problem.

In practice, the assumptions required for OLS to be appropriate can be stated in terms of properties of the residual vector ϵ .

In the rest of this lecture, we will generalise from Simple to Multiple Linear Regression, where $\beta = (\beta_0, \beta_1, \beta_2, \dots, \beta_p)$ and $2 \leq p \leq n$, so instead of fitting lines, we fit (hyper)planes to data.

Assumptions required for the linear model to be meaningful

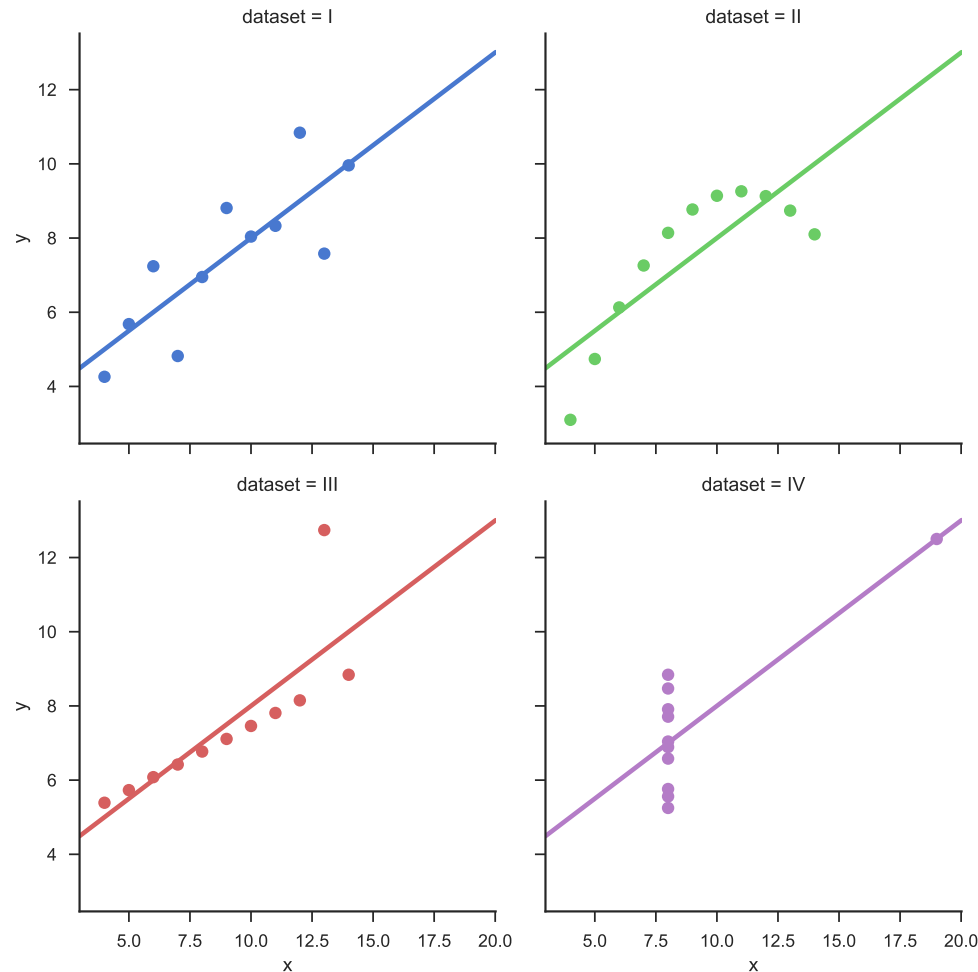
Definition 4 (Linear Regression Assumptions)

- ① The underlying relationship between the predictors and the response is linear in the regression parameters β .
- ② The residual errors ϵ are drawn from a (multivariate) Normal distribution $N(\mu, \sigma^2)$ where $\mu = \mathbf{0}$.
- ③ The predictors are not pairwise collinear, i.e., each pair of predictors β_{j_1} and β_{j_2} (associated with columns $X(:, j_1)$ and $X(:, j_2)$) have low correlation (equivalently, the inner product of $X(:, j_1)$ and $X(:, j_2)$ is far from zero).
- ④ There is no auto-correlation in \mathbf{y} .
- ⑤ The errors are *homoscedastic* (i.e., $\text{Var}(\epsilon)$ is constant over the range of \mathbf{x} or \mathbf{y}).

Remarks

- Because these assumptions depend both on the data and on the model fitted to that data, it is meaningless to say that “Data set A does not satisfy the linear regression assumptions”, because this observation might not apply to all formulations of all models applied to that data.
- Consequently, these assumptions can be used **constructively**, when model building, or **as checks**, when validating models.

Anscombe's quartet (1973)



Francis Anscombe devised 4 data sets to show different forms of misalignment between data and models. Sets I,II,III share the same x values. All 4 sets share approximately the same descriptive statistics (mean and variance), but little else is common to all 4!

Only I appears suited as it stands. The other data sets require some work, particularly IV.

What do you think needs to be done for each data set?

Common Cost Functions in Regression Models

Remember: for *least squares* regression, we are trying to minimise a loss function based on the error ϵ .

Measure	Definition	Purpose
Mean square error (MSE)	$\frac{(p_1 - a_1)^2 + \dots + (p_m - a_m)^2}{m}$	Mathematically tractable but places greater emphasise on observations with large error
Root mean square error (RMSE)	$\sqrt{\frac{(p_1 - a_1)^2 + \dots + (p_m - a_m)^2}{m}}$	Has same units as data
Mean absolute error (MAE)	$\frac{ p_1 - a_1 + \dots + p_m - a_m }{m}$	Does not overemphasise observations with large error (like MSE)
Relative square error (RSE)	$\frac{(p_1 - a_1)^2 + \dots + (p_m - a_m)^2}{(p_1 - \bar{a})^2 + \dots + (p_m - \bar{a})^2}$	Relative metric compares the error in the predictions with errors in the simplest model possible (a model just always predicting the average value of y)
Root Relative square error (RRSE)	$\sqrt{\frac{(p_1 - a_1)^2 + \dots + (p_m - a_m)^2}{(p_1 - \bar{a})^2 + \dots + (p_m - \bar{a})^2}}$	
Relative absolute error (RAE)	$\frac{ p_1 - a_1 + \dots + p_m - a_m }{ p_1 - \bar{a} + \dots + p_m - \bar{a} }$	

where a_j is the actual value, p_j is the predicted value, m is the number of observations, and \bar{a} represents the mean of the a_j .

Choices of Vector norms

Definition 5 (Manhattan norm)

$\ell_1(\dots) = \|\dots\|_1$ is the *Manhattan* norm (length) of a vector. Let $\mathbf{x} = (x_1, x_2, \dots, x_m)$. Then $\ell_1(\dots) = \|\dots\|_1 = |x_1| + |x_2| + \dots + |x_m|$ is the *Manhattan* distance of \mathbf{x} from the origin. Think of having to *walk* from one junction in Manhattan to another, the distance is the difference in the Street numbers plus the difference in the Avenue numbers.

Definition 6 (Euclidean norm)

$\ell_2(\dots) = \|\dots\|_2$ is the *Euclidean* norm (length) of a vector. Let $\mathbf{x} = (x_1, x_2, \dots, x_m)$. Then $\ell_2(\dots) = \|\dots\|_2 = \sqrt{x_1^2 + x_2^2 + \dots + x_m^2}$ is the *Euclidean* distance of \mathbf{x} from the origin. Think of being able to *fly* over all the buildings using the shortest route (think: Pythagoras theorem!) from one junction in Manhattan to another.

The Euclidean norm is very common, but the Manhattan norm is gaining popularity, because it is robust to outliers and computers are becoming powerful enough to handle it. However we generally use the Euclidean norm in this module.

Sidebar: Distance Measures for numeric data

Definition 7 (Minkowski p -norm)

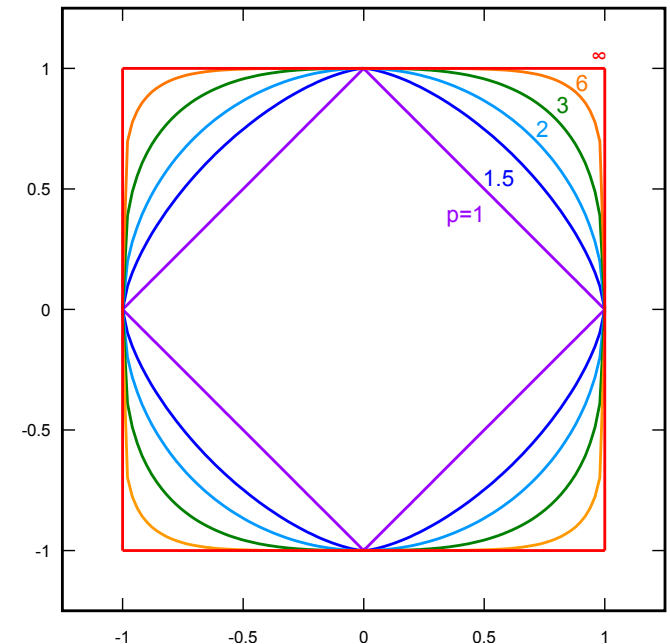
For a real number $1 \leq p < \infty$, the p -norm of \mathbf{x} is defined by

$$\|\mathbf{x}\|_p \equiv (|x_1|^p + |x_2|^p + \dots + |x_n|^p)^{\frac{1}{p}}.$$

The limiting case of $p = \infty$ is defined as

$$\|\mathbf{x}\|_\infty \equiv \max\{|x_1|, |x_2|, \dots, |x_n|\}.$$

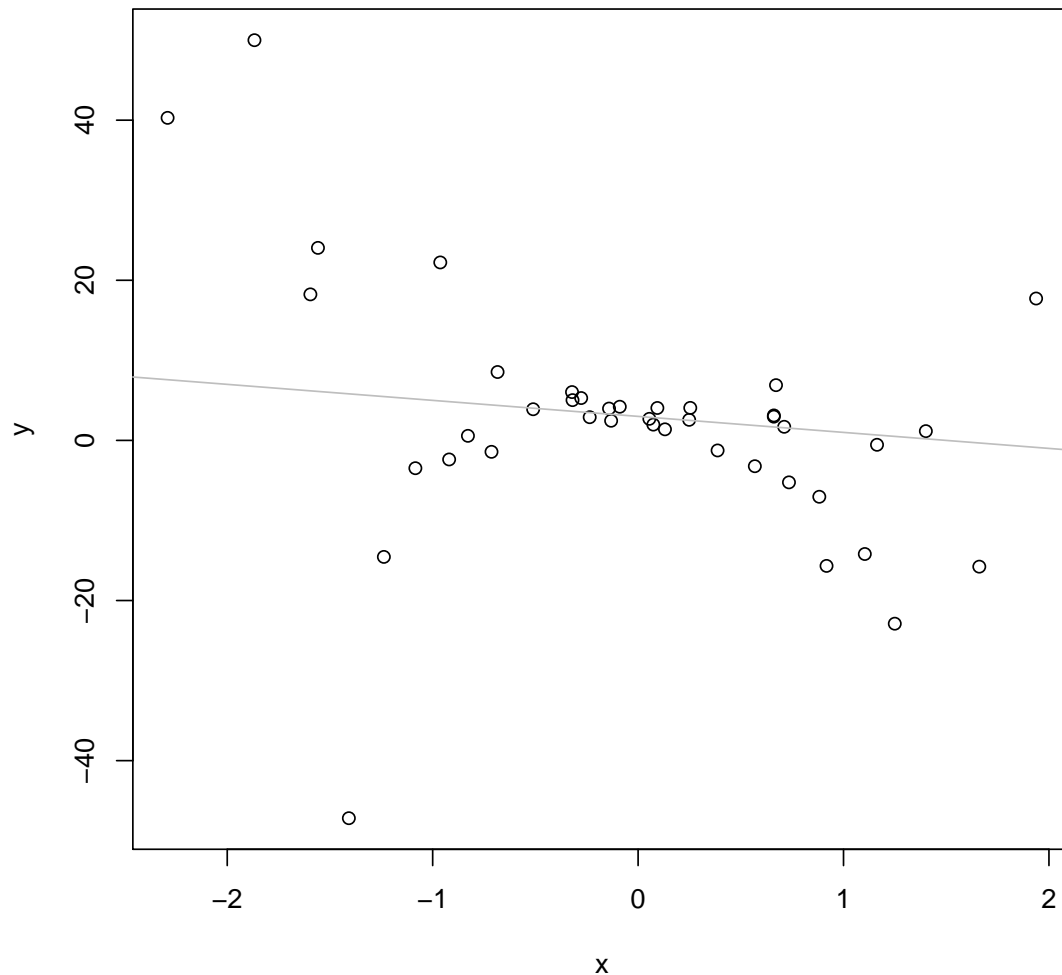
See the visualisation of the “unit balls” alongside, for $p = 1, 1.5, 2, 3, 6, \infty$.



Source: wikipedia

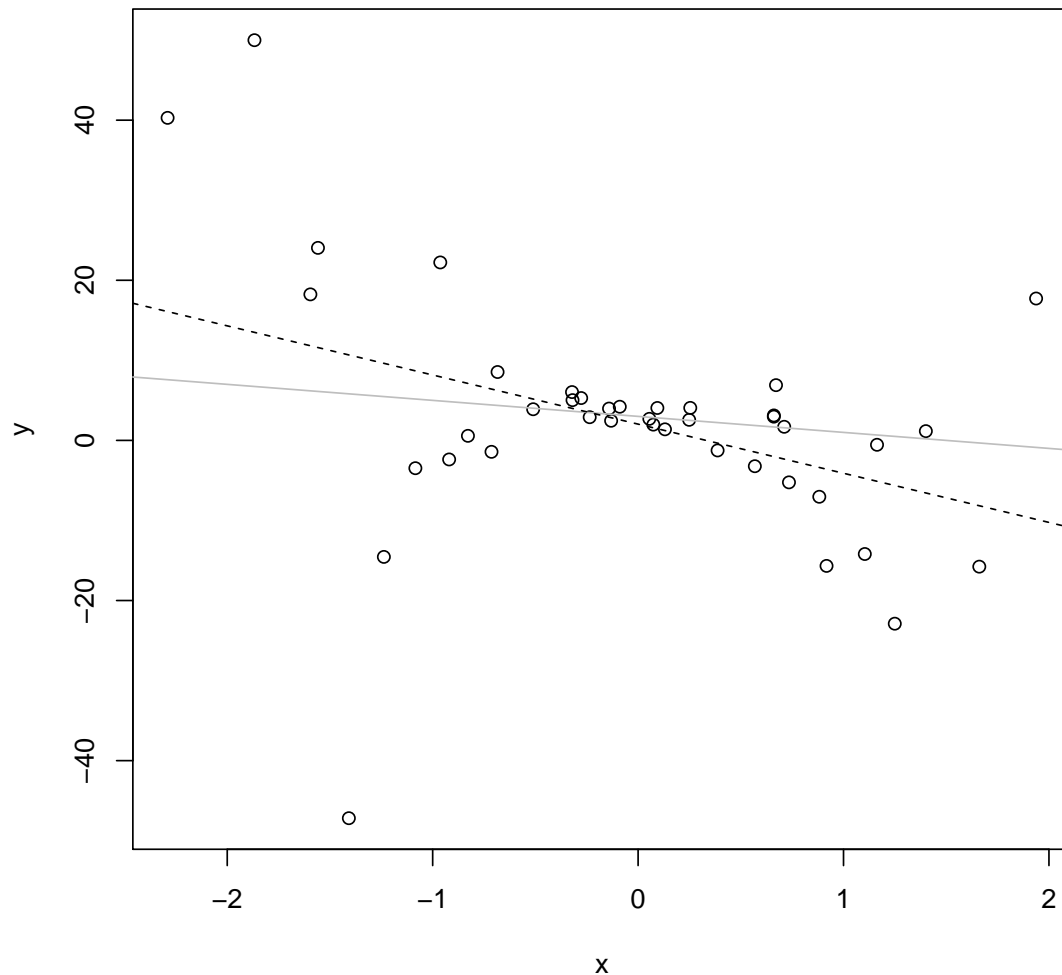
The most common norms are when $p = 1, 2$, or, ∞ . Choice of p depends on the application scenario. Can you think of when you would use each?

Case Study 1: Heteroscedasticity - Step 1



I generated 41 x, y points based on $y = 3 - 2x$, but with added errors that increase away from $x = 0$. The plot shows the line with $\beta = (3, -2)$ in grey.

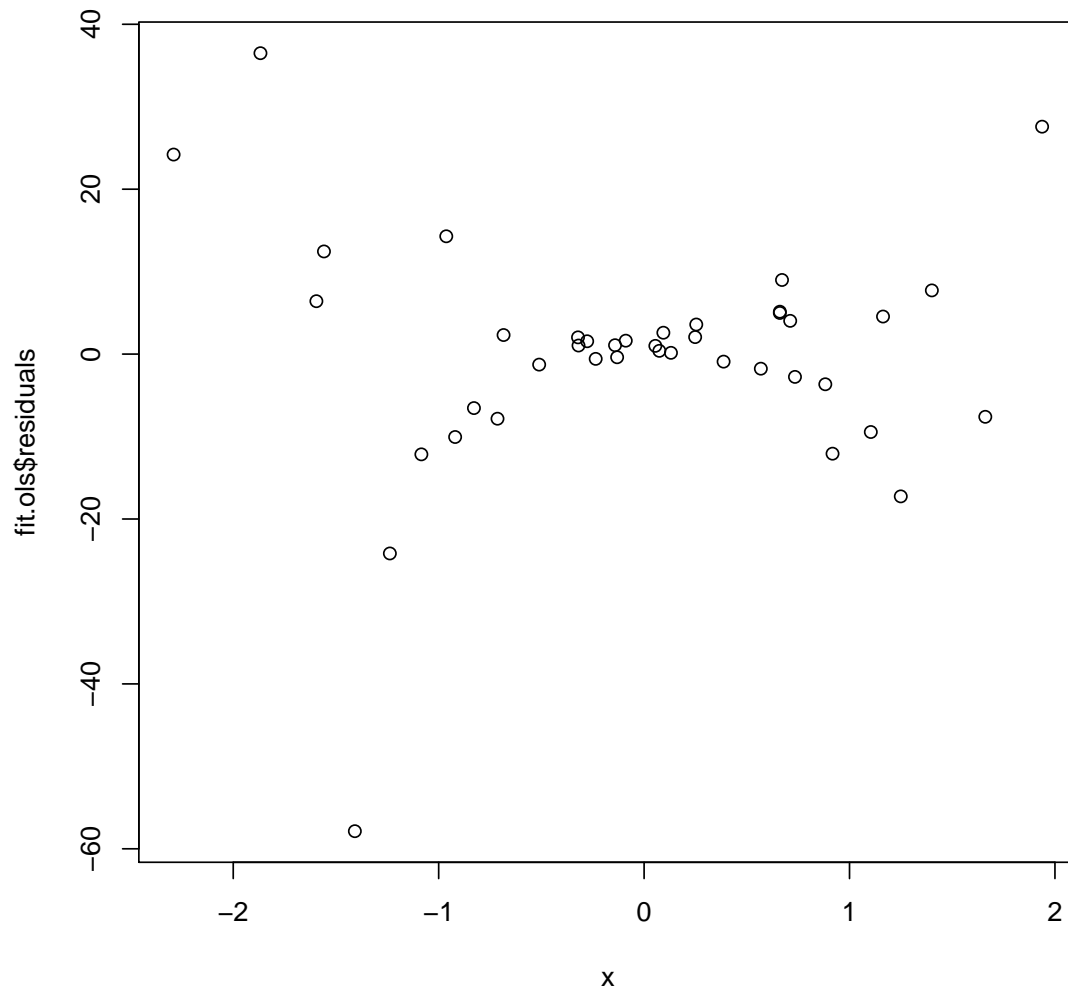
Case Study 1: Heteroscedasticity - Step 2



In this plot I added the OLS fit as a dashed line. Note that the parameters of the fit are quite different:

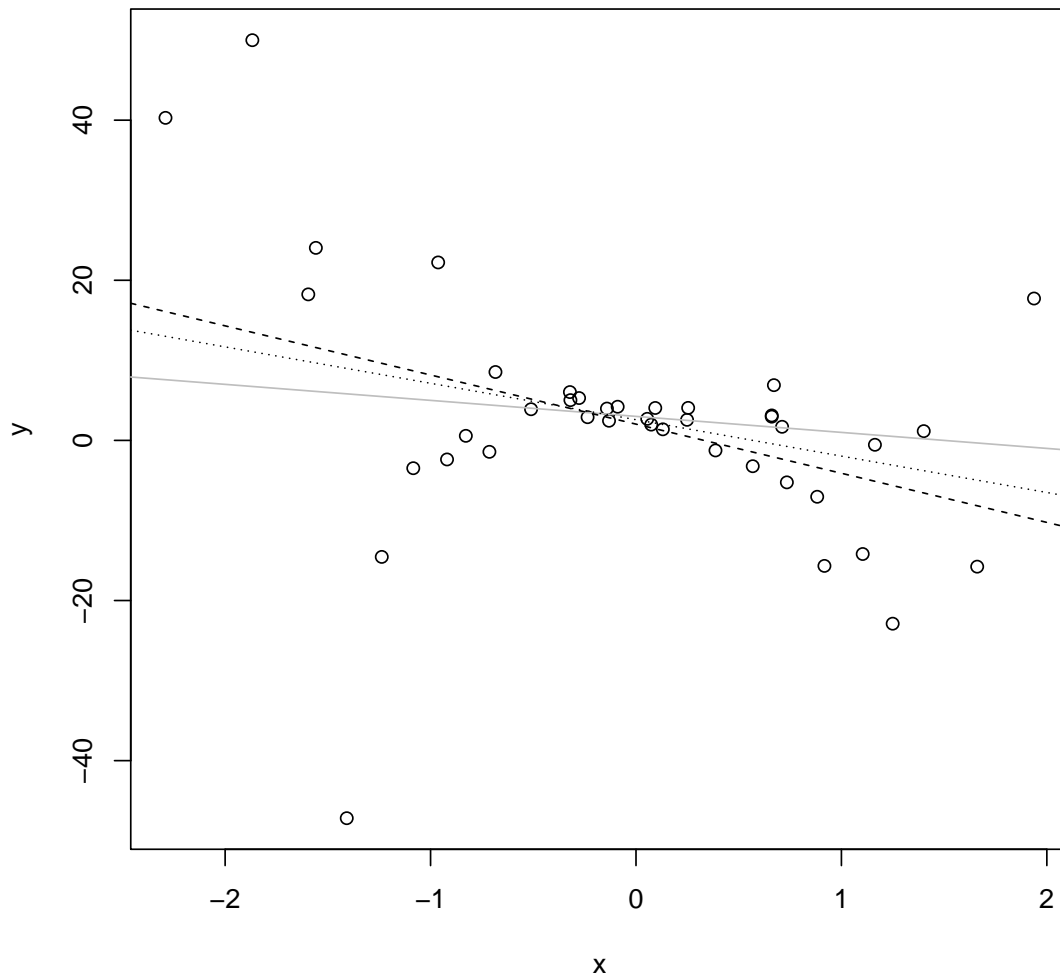
$\beta_{OLS} \approx (2, -6)$, equivalent to $y = 2 - 6x$.

Case Study 1: Heteroscedasticity - Step 3



This plot shows how the OLS residuals ϵ_{OLS} increase rapidly away from 0, as expected (since this was how the data was generated).

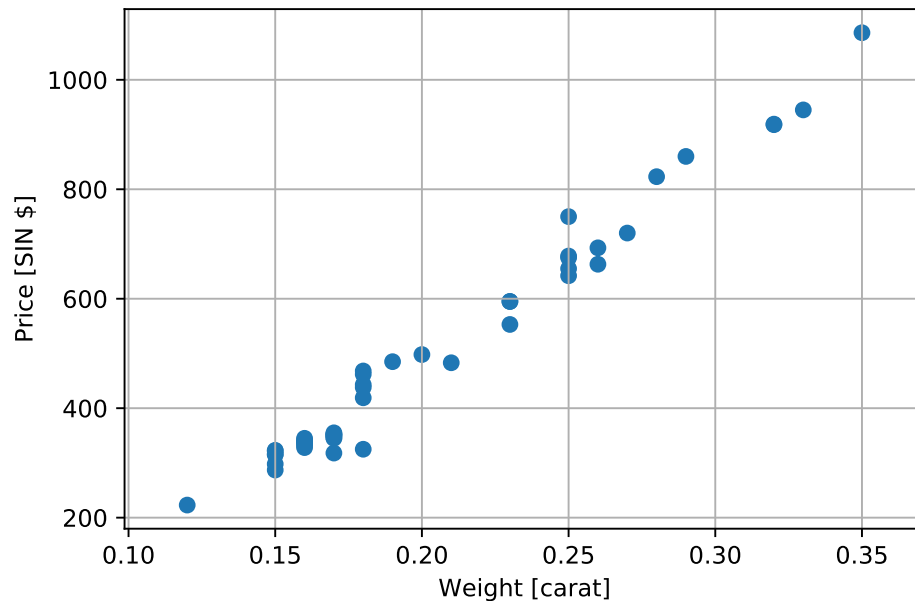
Case Study 1: Heteroscedasticity - Step 4



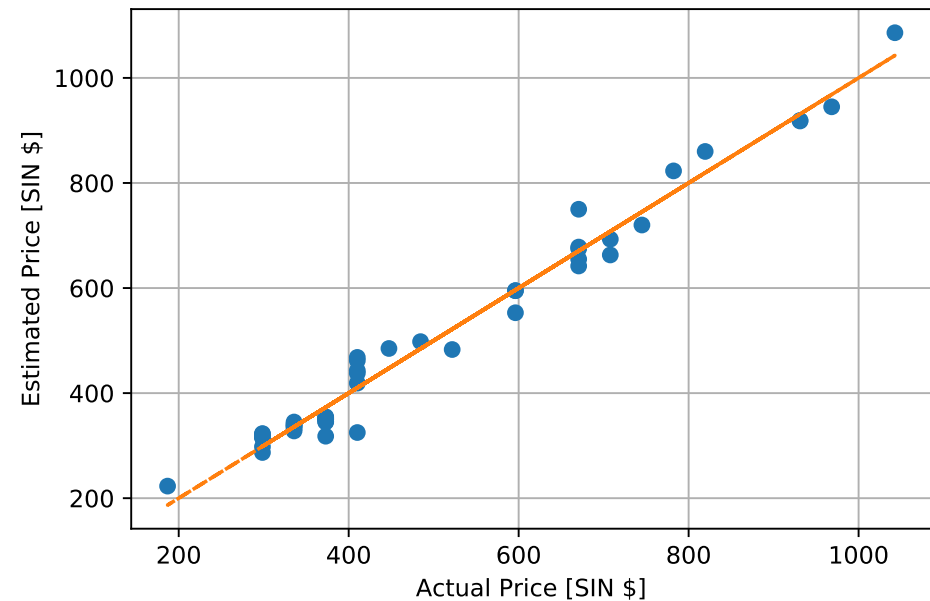
- By inspecting the previous residual plot I estimated a weighting function so that the residuals would be “more constant”. When this was used to scale the residuals, the resulting Weighted Least Squares estimates were $\beta \approx (2.6, -4.5)$ (shown as a dotted line) and hence closer to the “true” $\beta = (3, -2)$.
- So we were only partially successful at stripping away the noise and recovering the original line.
- **Can you see a problem with finding the weights?**
- *Iteratively Reweighted Least Squares* has been proposed to optimise regression models.

Case Study 2: Diamonds - Check relationship

Relation between diamonds' price and weight



Relation between estimated and actual diamonds' prices

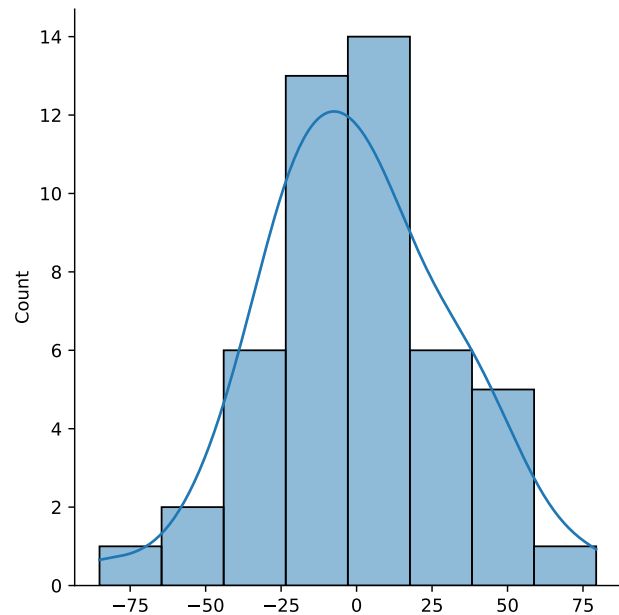


Clearly there is a linear relationship between a diamond's weight (in carats) and its price (in Singapore dollars, as here). So that is one assumption satisfied!

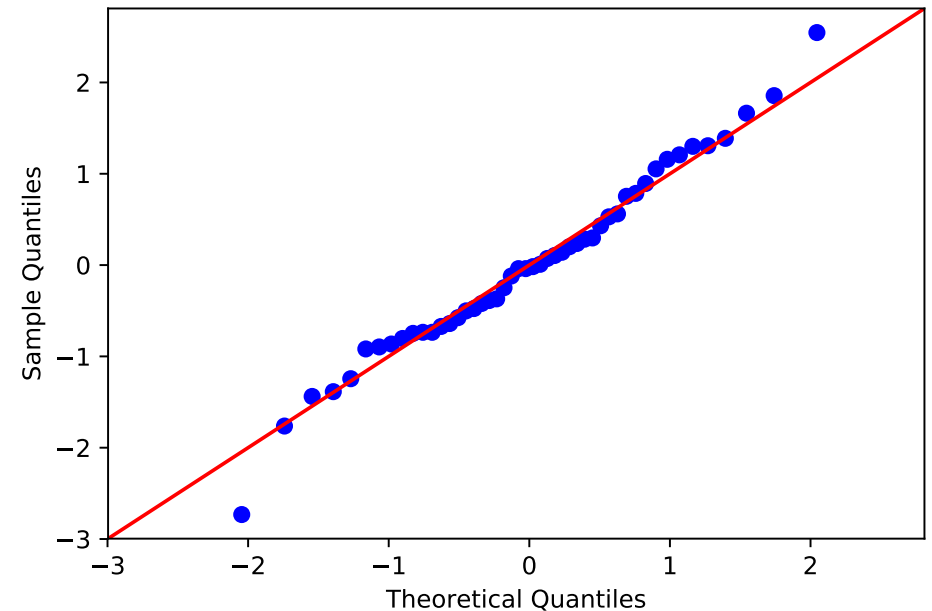
Sometimes the dependent variable has a linear dependence on a **feature that is a function of an attribute (column)**. Example functions include log, exp, sqrt, polynomial, etc. Even if the **function** is nonlinear in the attribute (e.g., the polynomials used in Kieran's model building example), that does not matter, as long as the model is **linear in the regression parameters β** .

Case Study 2: Diamonds - Check residual distribution

```
import seaborn as sns
resFig = "res/residHist.pdf"
sns_plot = sns.displot(x = residuals, kde=True)
sns_plot.savefig(resFig)
```



```
# Q-Q plot to verify the residuals distribution
resFig = "res/residualsqq.pdf"
fig = sm.qqplot(residuals, fit=True, line = '45')
fig.savefig(resFig)
```



Both diagnostic plots indicate the residuals are reasonably close to Normal distribution centred on 0. The qqplot is perhaps more informative. Looking good so far!

Is the standardised residual distribution heavy-tailed or light-tailed relative to the Normal distribution? Any other features?

Case Study 2: Diamonds - model summary

Dep. Variable:	price	R-squared:	0.978
Model:	OLS	Adj. R-squared:	0.978
Method:	Least Squares	F-statistic:	2070.
Date:		Prob (F-statistic):	6.75e-40
Time:		Log-Likelihood:	-233.20
No. Observations:	48	AIC:	470.4
Df Residuals:	46	BIC:	474.1
Df Model:	1		
Covariance Type:	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
const	-259.6259	17.319	-14.991	0.000	-294.487	-224.765
carats	3721.0249	81.786	45.497	0.000	3556.398	3885.651

Omnibus:	0.739	Durbin-Watson:	1.994
Prob(Omnibus):	0.691	Jarque-Bera (JB):	0.181
Skew:	0.056	Prob(JB):	0.913
Kurtosis:	3.280	Cond. No.	18.5

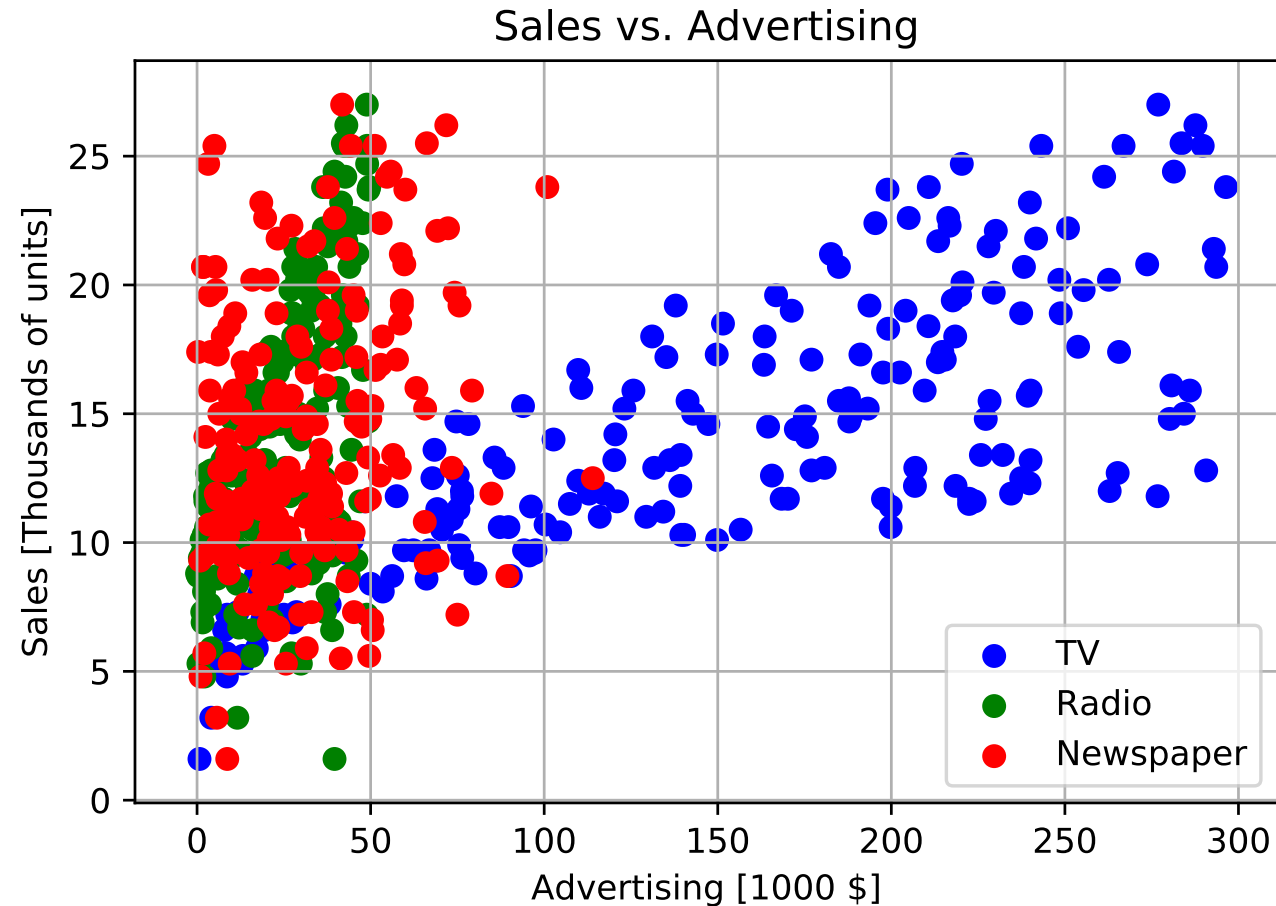
Dep. Variable:	price	R-squared:	0.978
Model:	OLS	Adj. R-squared:	0.978
Method:	Least Squares	F-statistic:	2070.
Date:		Prob (F-statistic):	6.75e-40

Case Study 3: Advertising: Data and Hypotheses

	TV	Radio	Newspaper	Sales
1	230.1	37.8	69.2	22.1
2	44.5	39.3	45.1	10.4
3	17.2	45.9	69.3	9.3
4	151.5	41.3	58.5	18.5
5	180.8	10.8	58.4	12.9

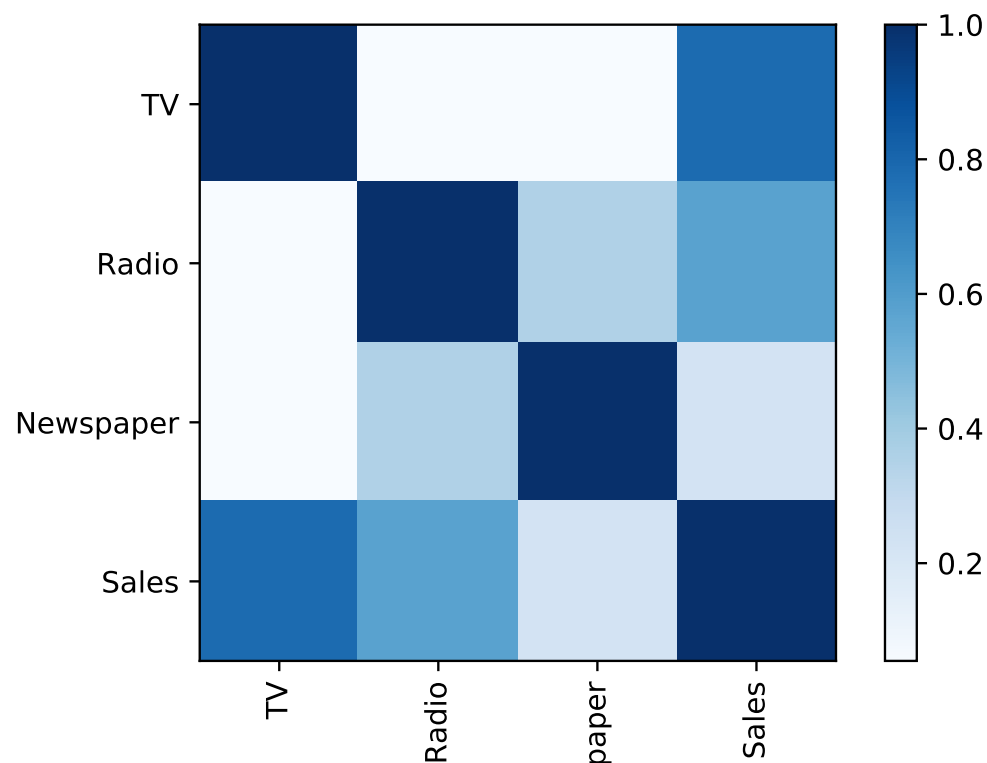
In this data set, the sales figure captures how many thousands of widgets of a particular type was sold in a year. Newspaper, Radio and TV represent the annual spend per widget type on the associated advertising channel. The hypothesis is that spend on advertising is a good predictor of sales performance. Since marketing budgets are limited, where should the adverts be placed for maximum sales?

Case Study 3: Advertising: Looking at the data



Which of the advertising channels appear to have a linear relationship with Sales?

Case Study 3: Advertising: Collinearity?



- Correlation matrix can indicate which attributes should participate in the model as predictors.
- A good predictor should have a **high correlation with the target** (Sales in this case) and should have **low correlation with other candidate predictors**.
- **What are expected to be good predictors for this data?**
 - Sales (the target) is placed in the last row (or column).
 - $TV > Radio > Newspaper$, with moderate correlation between Radio and Newspaper.

Sidebar: specifying models

The statsmodels way

- The dataframe contains the observed variables
- The model is specified separately
- Easier to change the model when experimenting

The sklearn way

- The dataframe contains the (computed) features
- The model is defined implicitly
- Standard interface across all sklearn

- statsmodels models are expressed like `"Sales ~ TV * Radio + poly(Newspaper,2)"`. This notation came from the applied statistics community.
- In words: “Sales depends on TV spending, Radio spending, the interaction between TV and Radio spending, Newspaper spending and Newspaper spending squared (5 features from 3 measured attributes).”
- statsmodels offers its own plotting (like seaborn but not as good). Its model summary is very convenient.
- sklearn exposes more of the details (e.g., choice of algorithm and configuration parameters).
- Both statsmodels and sklearn use the same libraries (scipy, numpy, etc.) underneath.

Case Study 3: Advertising: Model Building (“stats” way)

- Start from a “full model” and prune, versus from an “empty model” and add
- We choose the latter, as it is often easier to avoid overfitting

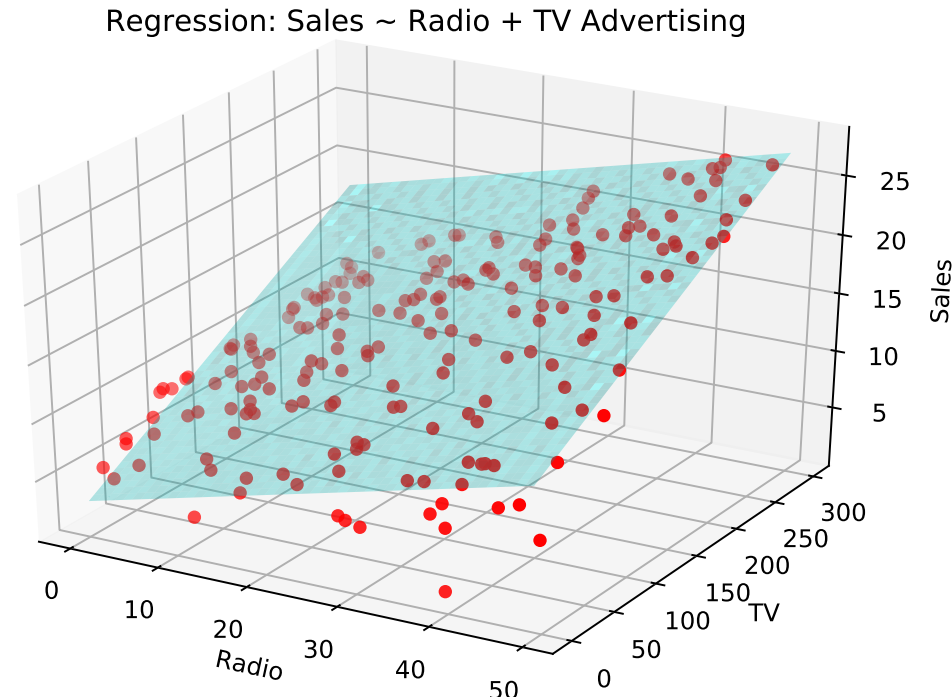
Example 8 (Forward Selection for Advertising Data)

Define: regression model score: (adjusted) R^2 (or MSE or similar metric) for a given model.

- 1 Fit “Sales \sim Newspaper”, “Sales \sim Radio”, “Sales \sim TV” and calculate their R^2 values.
- 2 Choose the best (largest R^2) single-term model (“Sales \sim TV” in this case) with $R^2 = 0.61$.
- 3 Fit “Sales \sim TV + Newspaper” and “Sales \sim TV + Radio” and choose the best by R^2 score, which is “Sales \sim TV + Radio” with $R^2 = 0.9$, which is significantly better.
- 4 Fit “Sales \sim TV + Radio + Newspaper”. Its R^2 score is also 0.9, so we favour the existing simpler two-term model (Occam’s Razor: other things being equal, choose the simplest model.).

So our preferred model is “Sales \sim TV + Radio” with adjusted $R^2 = 0.9$.

Case Study 3: Advertising: Viewing the Model



Since this two-term model ignores the contribution of the newspaper channel, the Newspaper spend as a contribution to Sales is just another component of the unmodelled (and apparently random) contribution to Sales.

However, the result is a model where every term is highly significant and the model “explains” 90% of the variance of the data, which is high for an observational study. **Why? Can we do better?**

Case Study 3: Advertising: Interactions; Interpretation

- Trying powers of the Radio and TV greater than 1 did not offer much more.
- However, by adding the TV, Radio interaction so that the model became “Sales \sim TV + Radio + TV:Radio” or equivalently “Sales \sim TV * Radio”, R^2 increased to 0.97 from 0.9, which is a significant improvement.
- All β terms have t –statistic significance of approximately 0.001 which is extremely significant.
- $\beta_0 = 6.75$, $\beta_{\text{TV}} = 0.019$, $\beta_{\text{Radio}} = 0.029$ and $\beta_{\text{TV:Radio}} = 0.001$, indicating that there is a favourable relationship between TV and Radio advertising ($\beta_{\text{TV:Radio}} > 0$), and that additional spending on Radio results in more Sales than the same spending on TV ($\beta_{\text{Radio}} > \beta_{\text{TV}}$).
- Spending on Newspaper advertising should be discontinued as its contribution to Sales is insignificant (indistinguishable from random noise).

Case Study 4: Credit balances - overview

Introducing

- the sklearn approach to regression (we used statsmodels with the Diamonds and Advertising data)
- non-numeric explanatory variables like gender and ethnicity
- more advanced regression modelling, e.g., handling correlated variables

Case Study 4: Credit balances - introduction

	Income	Limit	Rating	Cards	Age	Education	Gender	Student	Married	Ethnicity	Balance
1	14.891	3606	283	2	34	11	Male	No	Yes	Caucasian	333
2	106.025	6645	483	3	82	15	Female	Yes	Yes	Asian	903
3	104.593	7075	514	4	71	11	Male	No	No	Asian	580
4	148.924	9504	681	3	36	11	Female	No	No	Asian	964
5	55.882	4897	357	2	68	16	Male	No	Yes	Caucasian	331

- Note the presence of some categorical attributes (*Gender*, *Student*, *Married*, *Ethnicity*).
- These can participate in linear regression models to predict a numeric response, but must be coded first.
 - For example, *Gender* can become an indicator (0,1)-valued variable of the form *IsFemale*.
 - *Ethnicity* has 3 levels and is replaced by 3-1=2 indicator variables.

➤ A single categorical feature with n levels becomes $n-1$ (0,1)-coded “dummy” features.

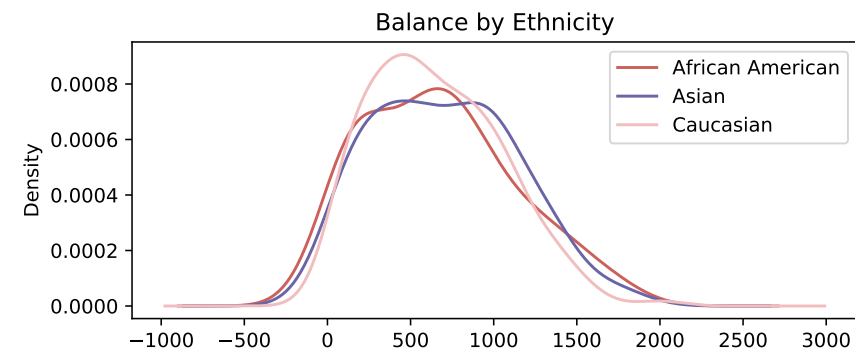
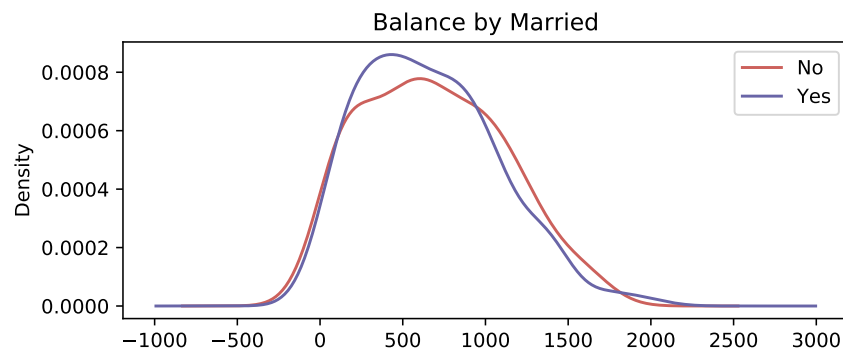
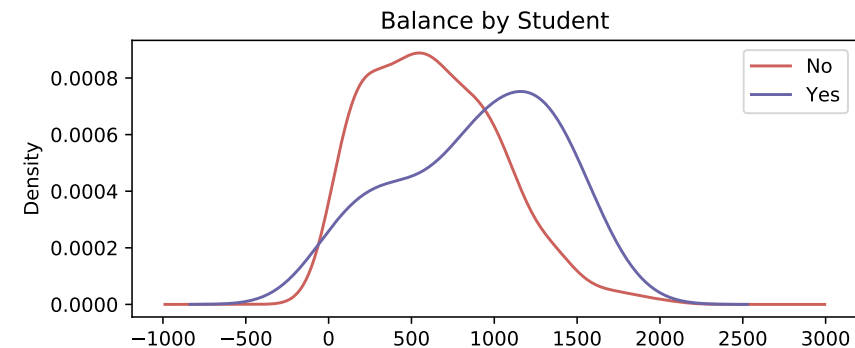
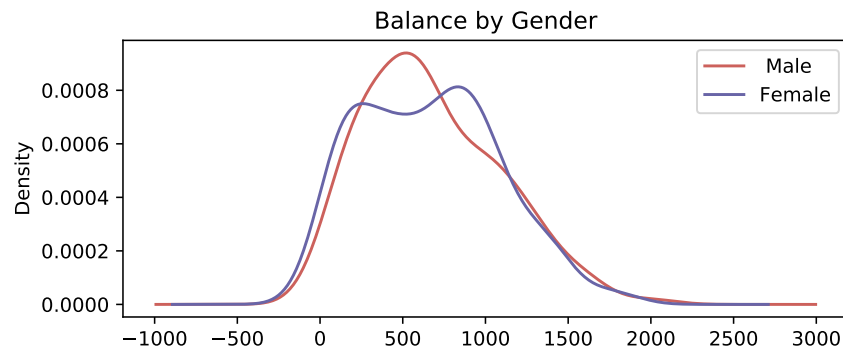
Case Study 4: Credit balances - Removing Data

- the purpose of the analysis is to predict credit balances.
- Basic exploratory data techniques (histograms) soon indicated that there were 2 cohorts
 - ① those who do not use their cards and/or clear their balance each month
 - ② those who use their cards and have nonzero balances
- Removing data relating to the first cohort meant that the remaining data looked more cohesive and also made linear regression easier
- Take-away: look for inconsistent subsets in the data, remove them if possible

Case Study 4: Credit balances - Removing correlated attribute

- Correlations between predictors are relatively high, but that between “Limit” and “Rating” is 1
- Generally, customers with a high rating are allowed to have high credit limits
- Conversely, customers will not be allowed high credit limits unless they have a high credit rating
- “Limit” was removed from the data used for analysis
- Take-away: remove all but 1 of correlated attributes from a set of such attributes, because they increase the standard error (hence variance) and make the solver’s job much more difficult (larger condition number)

Case Study 4: Credit balances - Contribution of Categorical Variables



Which of these categorical attributes has a significant effect on Balance?

Case Study 4: Credit balances - Model building

- Using forward selection as before, the best model was found to be “Balance \sim poly(Income,2) + Rating + Age + Student + Income:Rating”
- Could also use Backward Elimination to prune from a complex model
- For this data, high correlations between features can cause difficulties - we need techniques to handle this

Difficulties caused by correlated features

The Problem : Several features are highly correlated, so the solver has difficulty assigning an importance independently to each.

How it shows up : The condition score is large and several model coefficients take large values with opposite signs. Sometimes the solver gives up.

Solution options :

- ① Remove selected features from the model (simple, does not always work and requires care)
- ② Use *regularisation*, to “penalise” large model coefficients (solve a related problem with a different loss function)
- ③ Use *dimensionality reduction* (linear PCA) to derive an uncorrelated subset of the features with least loss in explanatory power (principal components can be opaque)

Regularisation introduction

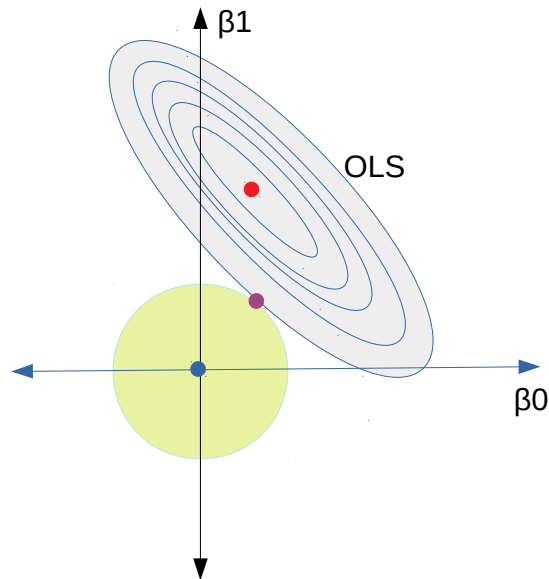
Add *regularisation* constraints to make the model work: $\min_{\beta} \|\epsilon\|_2^2 + \lambda p(\beta)$

- Options are
 - ① *Ridge Regression* where the penalty term takes the form $p(\beta) = \|\beta\|_2$
 - ② *Lasso* where the penalty term takes the form $p(\beta) = \|\beta\|_1$
- Regularisation has a metaparameter λ - the challenge is to choose a suitable value
 - if too large: tries less to match the data, increases the bias.
 - if too small: tries too hard to match the data so $\beta \rightarrow \infty$ and increases the variance

Ridge vs Lasso Regression

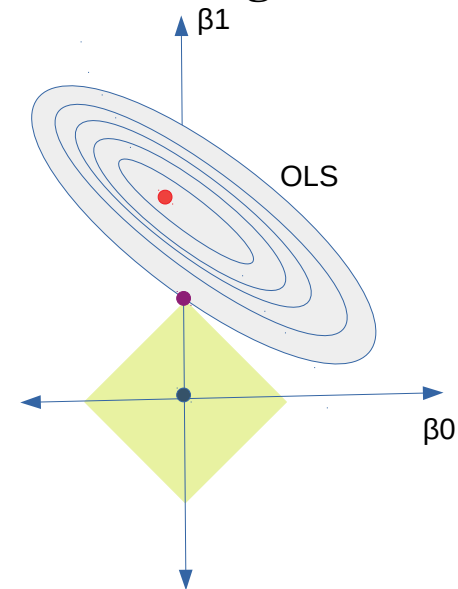
Because lasso regression favours the “corners” in parameter space, it tends to set some parameter values to 0 (essentially dropping the associated features). This has the added benefit of making the model smaller and easier to interpret.

Ridge Regression



Intersection point has $\beta_0 \neq 0$ and $\beta_1 \neq 0$ so both features are needed.

Lasso Regression



Intersection point has $\beta_0 = 0$ so its feature is no longer needed.

Case Study 4: Credit balances - Regularisation - Searching for λ

- 1 Choose a set of candidate λ values
- 2 For each candidate λ , use K-fold cross-validation (see Kieran's notes) on data subsets to estimate the prediction error for the regularised fit with that λ
- 3 Choose the λ for which the expected error is least
- 4 Now fit all the training data again with this choice of λ

Note that **lasso** (but not ridge) regression can set particular β_j to 0 (effectively removing them from the model), so it operates more like backwards elimination in terms of creating a more frugal model having fewer terms.

Ridge regression downweights certain terms but does not set them to zero. However, it can be more performant.

Diagnostic plots, statistics and worked examples

See notebooks accompanying this weeks lecture notes.

Review and summary

- Linear regression is one of the classic machine learning techniques.
- Compared to other techniques, statistics have more to offer, but ML objective (**minimise prediction error on test set**) is still as important!
- It has two phases, of which the first (learning from the training set) is generally the most challenging.
- It has many variants, so is quite flexible, but flexibility can be abused!
- Careful validation and model building is essential for success - it is an extension of the exploratory work done earlier in the process.
- In machine learning, prediction error is the main focus, but you need to be aware of other considerations such as
 - 1 model parsimony (keep model as small/simple as possible!): faster at both training and evaluation time
 - 2 the bias-variance dilemma: avoid overfitting and underfitting - remember, your model needs to generalise well from the training to the test set
 - 3 model interpretability: some models are easier to understand because the terms in the model represent concepts from the domain the data is from

Some Additional Resources

- Book: Introduction to Statistical Learning with R (2013) by James, Gareth and Witten, Daniela and Hastie, Trevor and Tibshirani, Robert.
I strongly recommend that you read Chapter 3 of the book, as it is very well written and available online for free.
- Kaggle notebooks relating to the datasets addressed this week. There are many, but searching Kaggle should provide nice examples of data mining in action.