

1 Practical data skills for the farm animal vet

2 Robert Hyde, Martin Green, Chris Hudson, John Remnant

3 Abstract

4 Whilst farm animal veterinarians receive extensive training in the diagnosis and treatment of
5 individual animals, population medicine or “herd health” have also become commonplace in modern
6 production animal veterinary medicine, with clinicians using routinely collected farm data to identify
7 epidemiological patterns to implement preventive management changes. Production animal
8 veterinarians are in a unique position having to large volumes of farm data available and having the
9 means to enact necessary changes on farms based on data analysis. Whilst data handling and
10 statistical techniques can initially be challenging to learn, a relatively small investment in time spent
11 learning can result in dramatically shorter times analysing farm outcomes, and a useful evidence
12 base on which to base on-farm management decisions. This article aims to provide a practical guide
13 to data skills for the production animal veterinarian, briefly covering basic descriptive statistics and
14 graphical representations of data, with further examples of more advanced techniques such as
15 statistical modelling.

16

17 Introduction

18 Farm animal veterinarians are frequently faced with disease or productivity issues that go far beyond
19 the clinical treatment of a single animal. There are often many potential factors influencing these
20 issues, and often the relative importance of these factors are unknown. Whilst farm animal
21 veterinarians receive extensive training in the diagnosis and treatment of individual animals,
22 population medicine or “herd health” have also become commonplace in modern production animal
23 veterinary medicine, with clinicians using routinely collected farm data to identify epidemiological
24 patterns to implement preventive management changes.

25 The use of statistical analysis and data science have become more common, particularly as
26 computing power has increased the availability of statistical methods. Whilst farm animal
27 veterinarians are unlikely to have received extensive training in statistical methods, they will often
28 be utilising many statistical principles already within their clinical work. Determining whether a
29 change in treatment protocol has altered clinical cure rates, a change in calf milk powder has altered
30 growth rates or a new fertility protocol has improved fertility efficiency, all require an understanding
31 of basic statistics. Many modern farmers now require this information to be analysed and
32 communicated to a high level from their routine veterinarian. This has been identified as important
33 role for farm vets in the future (Woodward et al., 2019).

34 Despite a large volume of medical research exploring advanced analytical techniques, examples of
35 translation into an effect in a clinical setting are relatively rare. The use of big data in farm animal
36 medicine is increasingly being discussed (Hudson et al., 2018), and production animal veterinarians
37 are in a unique position having both access to large volumes of farm data and the means to enact
38 necessary changes on farms based on the data analysis. This article aims to introduce the reader to
39 some of the underlying concepts used in data analysis as well as providing some practical tips to
40 apply these concepts in a clinical setting. The principles outlined may also be of value when
41 conducting research in practice or when interpreting the results of published research. Whilst
42 extensive training statistical techniques would be beyond the scope of a single article, the purpose of
43 this article is to highlight a selection of useful methods, and provide a basic overview of statistical
44 techniques, and an introduction to statistical coding that will have practical value to the farm animal
45 veterinary clinician.

46

47 Sources and types of data

48 In addition to on-farm records of reproductive management, treatment records and disease
49 incidences, many other sources of electronic data are available. These include milk recordings,
50 milking parlour software and environmental or animal sensors. It is rare to find a solution that
51 combines all these data into an easily accessible format for all analyses. Where more analysis is
52 required beyond those available in proprietary herd health software it will usually be necessary to
53 export sensor or on-farm records into a common format (such as a spreadsheet or comma separated
54 values file) for analysis. The detail of how to do this is beyond the scope of this article, however
55 sensor and machinery manufacturers may be willing to provide advice in individual cases. One
56 downside to data analysis is the manual entry of data. Paper based recording obviously necessitates
57 manual data entry into electronic form which can be time consuming. If the veterinarian is expected
58 to manually type up paper-based records this should be encompassed within the consultancy
59 charge.

60 Once data has been collated/exported for analysis, it is important to understand what type of data is
61 being analysed. This will inform the statistical or data science methods that can be used (Figure 1).

62

63

64 As well as understanding the source and type of data, the data should also be critically assessed for
65 potential flaws. The same is true of any metrics or key performance indicators (KPIs) calculated from
66 the data. A detailed description of interpreting metrics is provided by Hermans and others (2018),
67 and lag, momentum, bias and variation may all affect interpretation (Overton, 2009).

68 Data analysis overview

69 **Summary statistics**

71 Simple descriptive statistics are often a sensible starting point in data analysis and can provide a
72 useful overview of your data. Descriptive statistics often provide some sort of measurement of
73 central tendency (Figure 2) and some measurement of the spread of data (e.g. variance/standard
74 deviation or range, see Box 3 and Figure 7). Data can also be tabulated to show the number or
75 proportions falling within different categories.

76
77
78 There are several pitfalls to avoid when interpreting basic statistics. Consider a 100-cow dairy herd,
79 where you are asked to interpret results from a milk recording. Figure 3 shows the somatic cell count
80 (SCC) results for each cow. A relatively high mean of 208,000 cells/ml might seem relatively
81 concerning at a herd level. On further analysis however, most cows have a relatively low SCC, and
82 the mean is being largely influenced by a small number of cows with a particularly high SCC. This is
83 an example of a right “tail” or “skew”, as the mean (sum of all scores divided by number of samples)
84 is higher than the median (“middle” score when ranked) SCC is far lower at 115,000 cells/ml.

85
86 This example highlights the downsides of utilising single descriptive metrics for data analysis in
87 isolation. A simple histogram (Figure 3) shows immediately that there is a skew to the data, allowing
88 the clinician to focus any interventions in a more targeted area.

89 **Visualising data**

90 Graphical representation can be extremely powerful in data analytics, and even the most advanced
91 statistical techniques often start with a simple visualisation of the raw data distributions. Carefully
92 organised graphs can be extremely powerful knowledge exchange tools, taking relatively dull data
93 and transforming it into something intuitive and engaging. Whilst there are a great number of
94 statistical tests available to assess the structure of data, graphical representation of data can be one
95 of the easiest methods of getting to know the dataset. Visualising data can often be more
96 informative than calculating summary statistics and can also be particularly useful when assessing
97 data quality.

98 There are an enormous number of potential graphical software options to choose from, and a wide
99 range of techniques available depending on whether you are comparing one, two, three or even
100 more variables. Graphs can be created in Excel, and pivot charts can be useful to quickly investigate
101 relationships between many variables. The use of Excel based graphs however are often relatively
102 time consuming to repeat when analysing many farms compared with code-based approaches that
103 can be replicated for multiple farms in a fraction of a second, presuming data is recorded in the
104 same format. This example dataset involves weights taken from preweaned calves at a range of
105 ages. Graphs shown in Figure 4 were created using the ggplot2 package (Wickham, 2016) in R (R core
106 team, 2020) (see *Software*). A tutorial in downloading and getting started in R is available in
107 *Appendix i*, and code for the creation of these graphs is available in *Appendix ii*.

108

109 Box 1: Graphical representations of data (Figure 4)

110 Scatter: A scatter plot is a simple way of examining two continuous variables, in this case age (d) vs
111 weight (kg). As age increases weight unsurprisingly follows, and there seems to be a strong pattern.
112 By visualising an entire population in one graph however, more subtle patterns may be missed.

113 Regression line: Adding a “line of best fit” effectively demonstrates the association between age and
114 weight, i.e. for each increase in 1d, how much does weight increase. By adding an equation to the
115 graph, we can see that the formula for weight is best represented by:

$$116 \text{ weight} = 42 + (0.8 \times \text{age})$$

117 This means that the average daily liveweight gain is 0.8kg/d, with a y-intercept of 42 (the weight at
118 0d of age – an estimation of birthweight). This is useful to determine a rough daily liveweight gain
119 (DLWG) for a farm, however estimations from this method will be inaccurate if a small number of
120 weights are used. A more accurate method for calculating DLWG would be to take birthweights and
121 weaning weights for each individual animal.

122 Colour: By adding a grouping variable “sex” in this case, it now becomes apparent that whilst the
123 overall DLWG may be around 0.8kg/d there are clearly differences between sex in terms of growth
124 rate. Whilst heifers (blue) are generally born lighter than bulls (red), their DLWG appears to be
125 higher than the bulls, which may be suggestive of management differences between the groups to
126 be explored in more detail on the next farm visit.

127 Histogram: By “binning” observations into chunks it is possible to visualise distributions of single
128 continuous variables. This graph shows that whilst most calves are weaned around 100kg, there are
129 a small population of calves that are being weaned at 110-120kg. It is possible these animals could
130 have been weaned earlier and would be worth investigating why these animals were this heavy by
131 the time they were weaned.

132 Boxplot: If it is possible to calculate DLWG for each calf from a birthweight and a weaning weight,
133 DLWG can be expressed in boxplots to indicate the median (middle line), and interquartile range
134 (50% of calf DLWGs are within the box, see Box 3) and range (top to bottom of lines) to allow
135 visualisation of spread in data. Whilst it is important to achieve a high DLWG, it is also important to
136 do so consistently, and it appears in this case that heifers achieve both a higher and more consistent
137 DLWG than bull calves.

138 Violin plots: A version of the boxplot, where the width of the “violin” demonstrates the number of
139 observations at a given value. Whilst there are a small number of bull calves that have a higher
140 DLWG than heifer calves, they are relatively inconsistent, with heifer calves being much more
141 consistently close to 1kg/d DLWG.

142 Regression

143 There are various types of statistical “model”, which are methods of exploring relationships between
144 variables that we might use in a clinical veterinary environment. Regression models are often useful
145 in both clinical and research settings. Regression allows the strength of the association between two
146 variables to be measured. Regression also has the benefit that multiple variables can be investigated
147 simultaneously (multiple regression). It is also possible to carry out regression with a binary outcome
148 (logistic regression), where the binary variable is transformed to a probability-based scale. Linear
149 regression is a commonly used statistical method within research but is relatively easy to perform in
150 practice to determine effects at a farm level.

151
152 Consider a dairy farm that manages their beef and dairy calves separately and wish to analyse the
153 performance of both groups with their veterinarian. The farm has been weighing calves >21d, and
154 after plotting the ages and weights of calves, it appears there may be some difference between the
155 two groups. A simple regression model can put a numeric figure on this difference, which may aid in
156 clinical decisions.

157
158 It is vital for regression models that all data points are independent (the *assumption of*
159 *independence*). For example, if the only Belgian blue bull calf were weighed 5 separate times, the
160 effect of being a beef calf might be exaggerated through multiple inclusions of a single fast-growing
161 animal. For these models it is important that each data point comes from an independent unit (in
162 this case a single calf weighed once) with as few biases between groups (such as breed, sex or
163 housing) as possible. The failure to achieve independence of observations may lead to incorrect
164 model results, and whilst in research this can be achieved by prospective random allocation to
165 groups, this will not always be practical in a clinical setting. “Mixed-effects” models can be used to
166 overcome issues with independence and can account for example multiple measurements per calf,
167 or multiple calves per farm within the model to account for the effect of “calf” or “farm” on the
168 outcome of interest. Even with appropriate study design and model analysis, correlation does not
169 necessarily imply causation, and it is important to fully understand the dataset through graphical
170 analysis as well as statistical modelling.

171 **Box 2: Linear regression**

172 To analyse the weights of beef and dairy calves, a clinician might first examine the weight of calves
173 by age. This can easily be done in Excel by highlighting the “Age” and “Weight” columns and
174 inserting a scatter plot. A trendline can be added to the graph, and a linear regression equation can
175 be added to the line. In this example, calves gain on average 0.85kg for every 1d increase in age (a
176 DLWG of 0.85kg), and the y-intercept (weight at age 0d) is 31.9kg on average. An R² of 0.87 means
177 that around 87% of the variation in weight is explained by the predictor variable in the linear
178 regression model.

179

180 Whilst plotting data in Excel (Figure 5) allows the analysis of weight by age, it would not be easy to
181 add more variables for analysis. For example, we might be interested in whether beef or dairy calves
182 grow faster (Figure 6). By using coding software, it is extremely easy to add or remove as many
183 variables as required as shown below.

184

185 The formula for a regression model is shown below. This model will attempt to model weight, by
186 using the variables age (d) and breed type (either beef or dairy), which would have been challenging
187 to analyse in Excel.

188 `weight ~ intercept + age + breed_type`

Variable	Coefficient	95% confidence interval	P value
Intercept	40.62	37.22-44.03	<0.001
Age	0.80	0.76-0.85	<0.001
Breed: Beef	7.15	4.52-9.78	<0.001

189 Table 1: Results from regression model exploring effect of age (d) and breed type (beef or dairy) on
190 weight (kg).

191 Results from this model are shown in table 1. An “Intercept” of 40.62 suggests that at day zero
192 (birth), calves weigh on average 40.62kg (and we are 95% sure that this value lies between 37.22 and
193 44.03kg). A coefficient of 0.80 for variable “age” means that for every 1d increase in age, calves can
194 be expected to grow 0.80kg (effectively a DLWG of 0.80kg/d). A 95% confidence interval effectively
195 means that whilst our best guess is 0.80kg/d, we are 95% certain that this figure is between 0.76-
196 0.85. Whilst model outcomes suggest an association between age and weight, this does not always
197 indicate causality. In this model every 1d increase in age is associated with a 0.8kg increase in
198 weight, however the increase in weight is not directly caused by the increase in age. A coefficient of
199 7.15 for the categorical variable “Breed: Beef” suggests that beef calves are 7.15kg heavier at any
200 given time point compared with dairy calves (the default or “Reference” category “Breed: Dairy”).
201 Again, our confidence interval shows that we are 95% sure beef calves are between 4.52-9.78kg
202 heavier than dairy calves. In this example, it appears we could be quite confident that beef calves
203 are significantly heavier than dairy calves. Whilst this model is likely to provide a robust estimate of
204 calf weight within the age ranges shown, it is important to be cautious when extrapolating model
205 predictions to areas with no data. In this example we have data from ~21-120d of age and the
206 effects of age and breed type on weight reported in the regression model be accurate if applied
207 earlier or later in life. The R² for this model is 94%, which indicates more of the variation in weight is
208 explained by the predictor variables when age and breed are accounted for than the previous Excel
209 model using only age. Multiple variables can be added to regression models such as the housing
210 type, feeding and other management factors, and regression models provide a powerful method of
211 determining the impact of multiple factors on an outcome of interest such as DLWG. The concept of

212 “overfitting” in model building is largely beyond the scope of this article and describes analyses that
213 correspond too closely to a particular set of data but failing to fit with additional data. For example,
214 by including “Calf ID” in our model, we might well be able to predict DLWG perfectly, but the model
215 would not be very useful for the prediction of DLWG from future calves, or different farms. Graphical
216 analysis of data in conjunction with model outcomes should help prevent errors in the over
217 interpretation of model results in relatively small datasets.

218

219 Box 3: Statistical definitions

220 Correlation: The degree of association between two variables.

221 Interquartile range: A measure of the variability or “spread” of data representing the “middle” 50%
222 of the data; between the 25th (1st quartile) and 75th (3rd quartile) percentiles.

223 Variance and standard deviation: Measures of the variation around the mean. A lower variance or
224 standard deviation suggests datapoints are closer to the mean value (see Figure 7). Standard
225 deviation is calculated as the square root of the variance, and variance as standard deviation
226 squared.

227 Linear regression: A statistical modelling technique that can be used to analyse the strength and
228 nature of an association between single or multiple variables and a continuous outcome such as
229 weight.

230 Logistic regression: A statistical modelling technique that can be used to analyse the strength and
231 nature of an association between single or multiple variables and a binary outcome such as
232 conception to a given insemination.

233 Multivariable regression: Either linear or logistic regression where there are more than one predictor
234 variables included such as the effect of both powder type and age on weight.

235 Machine learning: A set of algorithms that can be used to classify or predict an outcome by
236 “learning” patterns from data, often outperforming more traditional approaches such as regression
237 in terms of predictive performance.

238 Black-box: Where the effect of individual variables on the outcome are challenging to interpret
239 (often used to describe types of machine learning algorithms)

240

241 **Box 4 – Useful spreadsheet functions**

242 Spreadsheets can be a useful and accessible way to store, analyse and visualise data. The internet is
 243 full of discussion boards and blogs giving advice on functions and techniques for spreadsheets and
 244 other analysis software. Using a search engine to search for what you want to achieve is often the
 245 quickest way of finding an effective approach. The help features in the software and manufacturers
 246 websites and instructions also usually provide details examples of using functions. The following
 247 examples from Microsoft Excel (Table 2) are functions or features that the authors have found
 248 valuable when reviewing data, equivalent functions will be available in other software packages.

Function or feature	Description	Example usage
"=" or "Fx"	Used to enter a function into a cell, you can type "=" at the start or click on the <i>Fx</i> symbol in the formula bar. Using the formula bar allows you to search for functions and see a description of the correct usage more easily.	To enter a formula or perform a calculation in any cell. The <i>Fx</i> bar can be used to search for a formula and check how to enter it.
Fill handle	The small box in the bottom right corner of the cell can be dragged down or across to fill the following cells. Double clicking the box fills the column downwards to the end of the data. Excel recognises a range of options here from days of the week or dates to sequences of numbers or letters. The fill handle will also allow you to drag a formula across lots of rows.	To drag a formula down across multiple rows of data to perform the same calculation on lots of lines of data.
Static and dynamic references	Dragging a column down in Excel usually updates the "reference" (the other cells the formula uses) with each cell i.e. they are dynamic (cell A1+ cell B1 in cell C1 will change to A2+B2 in C2 etc). You can use a \$ sign (or press F4) to make them static (cell A\$1+ cell B\$1 in C1 will still be A\$1+B\$1 in C2 etc when dragged down).	To fix a cell in a formula so that when you drag and fill the formula this part doesn't change. For example, a mating start date or average milk yield value that will be used lots of times on different rows of data.
Custom sort	The custom sort option allows data to be sorted by multiple columns.	To sort calf weight data by eartag number and then date so that all weights for the same calf are listed, in date order, together
IF	Returns different values or text depending on whether a	To recode continuous data (e.g. Milk somatic cell count)

	<p>logical statement is true or false e.g. =IF(logical argument, value if true, value if false).</p> <p>If functions can be “nested” within each other to allow further logical arguments to be used.</p>	<p>as a category (e.g. infected or not infected).</p> <p>To apply a formula only if certain criteria are met e.g. if this eartag number matches the one above calculate the difference in weight, if they are not the same leave the cell blank.</p>
COUNTIF	<p>This counts the number of cells in a given range that meet a specific condition e.g. =COUNTIF(range, criteria)</p>	<p>To count the number of cows with a somatic cell count of over 200,000.</p>
AVERAGE	<p>The AVERAGE function calculates the arithmetic mean of a selection of data, other measures of central tendency are available and include MEDIAN, MODE. Other summary statistics are also available for example MIN and MAX.</p>	<p>To calculate the average age at first calving from a list of calculated age at first calving dates.</p>
VLOOKUP	<p>VLOOKUP can be used to lookup data from another worksheet using an identifier common to both sheets e.g. =VLOOKUP(lookup_value, table_array, col_index_num, range_lookup). The lookup_value is the cell in the first table that is being referred to, the table array is the second table, the col_index_num is the column number of the second table that contains the data of interest. range_lookup is either TRUE (for an approximate match) or FALSE (for an exact match). INDEX and MATCH formulas are more advanced but can be used to provide greater flexibility than VLOOKUP.</p>	<p>To add the date of birth of calves to a table containing weight data. The calf eartag number can be used as the identifier and a column added to a table of weights that uses VLOOKUP to find the date of birth in a different table containing calf eartag numbers and birth dates. This would allow the age at each weight recording to be calculated.</p>
Pivot tables	<p>Pivot tables are a way to summarise findings from longer tables and lists of data. There are lots of ways of using</p>	<p>To summarise individual cow data to produce a summary plot of average milk yield by calendar year or to summarise milk recording data by</p>

	and combining data in to a pivot table (or pivot chart).	displaying counts of infected cows by calendar month.
--	--	---

249 Table 2: There are a number of useful tools and functions available in Excel that can be used to
250 quickly analyse and report farm data.

251

252 P-values, hypothesis testing and sample sizes

253

254 P-values are usually calculated for hypothesis testing. This is when a research or clinical question is
 255 rephrased into a hypothesis – a statement that can be tested. The alternative hypothesis is usually
 256 the statement that would be true if there is a difference, for example “calves fed milk replacer A will
 257 have a higher growth rate than calves fed milk replacer B” or “Cows artificially inseminated by
 258 operator X are more likely to become pregnant than those served by operator Y”. The null
 259 hypothesis is the opposite statement, in this case that there is no difference between milk replacer A
 260 and B or operator X and Y. There are lots of hypothesis tests available, an extensive description of all
 261 possible tests is beyond the scope of this article. When using statistical tests on data, it will often be
 262 necessary to decide whether to use parametric or non-parametric tests. Data that are parametric
 263 follow a typical distribution that can be defined with “parameters” that follow certain patterns. The
 264 most common parametric distribution described is the “normal” or “gaussian” distribution. A normal
 265 distribution can be described by its mean and standard deviation, we know that 95% of the data will
 266 be spread symmetrically 1.96 standard deviations either side of the mean value (Figure 7).

267 Many parametric statistical tests will have an “assumption of normality”. Whilst the software may
 268 give a value, if the assumptions are not met the results may be inaccurate. If in doubt, non-
 269 parametric tests are usually considered a more conservative approach, being less likely to
 270 overestimate statistical significance.

271 Common hypothesis tests and the type of data they are used on are shown in table 3.

Type of data	Hypothesis test	Example
Comparing two continuous, parametric, normally distributed groups	T-test	To test for differences in daily live weight gain between calves fed two different milk replacers
Comparing continuous non parametric outcomes between two groups	Mann-Whitney test	To test for differences in age at first calving between heifers reared on two different diets
Comparing multiple groups with a continuous outcome	ANOVA	To test for differences in daily live weight gain between three different breeds of calf
Comparing categorical data e.g. proportions	Chi-squared tests	To test for differences in the proportion of cows becoming pregnant following insemination by two different operators

272 Table 3: Common hypothesis tests for various data types.

273 In a research setting, P-values being <0.05 is a conventionally accepted threshold to indicate the
 274 statistical significance of a result. The definition of a p-value is the likelihood that the results would
 275 have been obtained if the null hypothesis was true (for example that there was no difference
 276 between two interventions being tested). In simple terms, a p-value of 0.05 effectively means that if
 277 a study was repeated 100 times, we might expect this result to occur purely by chance five times,
 278 even if there was no difference. The somewhat arbitrary cut-off of $p < 0.05$ suggesting *statistical*
 279 significance does not guarantee *clinical* importance, and similarly $p \geq 0.05$ does not mean there is *no*
 280 clinical importance. P-values do not describe effect size, and clinicians must consider the likelihood
 281 of a true difference between interventions, likely effect sizes and the cost of potential interventions.

282 For example, if a fertility treatment were shown to increase pregnancy rate with a p-value of 0.04,
283 would it be worth the clinician using the treatment purely because the result was “statistically
284 significant”? This result may have occurred by chance, although it is much less likely than if the p-
285 value were 0.4. There may be a true difference between treatment and no treatment, but it is also
286 worth considering the effect size and cost of the intervention. If the intervention cost £10 and
287 resulted in only a 1% increase in conception rate this may no longer be worth considering even
288 though there was a “statistically significant” difference. Similarly, an intervention costing £10, with a
289 20% increase in conception rate, but with a p-value of 0.07 might still be deemed clinically
290 worthwhile, (if the sample size of animals tested was larger, a true difference might have been
291 found).

292 This concept can be applied to clinical examination of animals; if examining a dairy cow that is
293 clinically normal with a temperature 0.1°C higher than the clinicians’ “normal” range, this would
294 might constitute a statistically significant finding (being outside of the normal 95% coverage range
295 for normal individuals) but the clinician may not consider this clinically significant in terms of
296 deciding on treatment. Conversely, a clinician detecting temperature at the upper end of “normal”
297 in an animal demonstrating clinical signs suggestive of an infectious disease might place less ‘clinical
298 weighting’ on temperature in terms of clinical decision making

299 There has been a move away from p-values within scientific literature, to end the string of hyped
300 claims and dismissal of potentially crucial effects due to this arbitrary threshold that is commonly
301 referred to (Amrhein et al., 2019; Ioannidis, 2005). Sample size calculations (also known as power
302 calculations) are essential to determine the number of animals required to have a reasonable chance
303 of detecting a statistically significant difference. This largely depends on the size of the expected
304 difference between two groups. For example, more calves would be required to detect a difference
305 of 0.05kg/d in growth rate between two groups of calves than to detect a 0.1kg/d difference in
306 growth rates. Typically, sample size calculations are used to calculate the numbers of animals
307 needed to have an 80% chance of detecting a true difference between groups (the “power”) at a
308 significance level of 0.05. For example, to have an 80% chance of detecting a statistically significant
309 difference of 0.05kg/d between two groups of calves, we would require at least 63 calves to be
310 present in each of the two groups. If we have fewer calves than this and do not detect a statistically
311 significant difference then it does not mean there is no difference between the group, as this study
312 would be underpowered to detect a difference.

313 When analysing whether results from analysis are likely to be clinically important, it is worth
314 considering the effect size, p-value and cost/benefit of a given intervention. It is important
315 veterinary clinicians are involved in data analysis and interpretation as these complex interactions of
316 client expectation and aspiration, uncertainty, cost-benefit and practicality go far beyond a simple
317 (or even a complex) mathematical equation, even before interventions are planned and
318 implemented. It is important to understand the various pitfalls of p-values, and not place too much
319 reliance on a somewhat arbitrary threshold, particularly if the study is underpowered.

320

321 **Box 5: Epitools**

322 Whilst statistical significance and sample size calculations can be performed in Excel to calculate
323 statistical significance, it can often be complex to arrange. Epitools is one of a number of online tools
324 where a number of statistical tests can be performed and is freely available. For example, whether a
325 fertility intervention would be statistically significant could be calculated with using the following
326 2x2 table in figure 7 (<https://epitools.ausvet.com.au/twobytwotable>).

327 Using Epitools, a 2x2 table can be filled out to calculate statistical significance between two
328 treatment groups (Figure 8). In this example, 18/50 cows were pregnant in the treatment group,
329 with only 10/50 cows being pregnant in the control group. This resulted in a p-value of 0.07, and the
330 clinician will have to decide whether this result is clinically significant, even if not “statistically
331 significant” by conventional thresholds of $p < 0.05$.

332 A sample size calculation can also be performed using Epitools using the sample size calculator
333 (<https://epitools.ausvet.com.au/twomeansone>)

334 Using Epitools, a sample size can be calculated (Figure 9) to see how many animals would be
335 required to have an 80% (power 0.8) chance of detect a difference in growth rates of 0.8kg/d and
336 0.85kg/d at a significance level of $p < 0.05$ (confidence level 0.95). This resulted in a requirement of 63
337 calves per group (126 calves in total), and any fewer calves means that a difference of 0.8kg/d and
338 0.85kg/d might not be detected.

339 Software

340 Even with limited experience, the calculation of basic statistics and the creation of powerful
341 visualisations are relatively straightforward and can be valuable tools for on farm investigations and
342 routine herd health analysis. A wide range of data analysis tools are available, and some of the most
343 powerful are free to download. Most veterinarians will be familiar with Microsoft Excel, which can
344 certainly provide adequate graphs and basic statistical analysis for most situations. Microsoft Excel
345 or other spreadsheet packages can also be useful for collecting and storing data for analysis where
346 this is not possible in on farm software. Learning a coding language provides a far more powerful
347 data analysis toolkit to the farm animal vet, which ultimately will result in higher quality, faster
348 reports for farms. Proprietary on-farm or vet software will often provide tools to create graphs and
349 reports, but these are often limited in their options and do not provide the flexibility that Excel or
350 coding languages provide.

351 Whilst coding requires more effort to learn initially, there are enormous benefits in the long term, in
352 terms of flexibility and repeatability. Analysis that might take 10 minutes per farm to create in a
353 spreadsheet can be replicated almost instantly for any number of farms once code has been written.
354 There is also vastly improved flexibility with coding, with almost limitless opportunities compared
355 with the more limited analytical potential of spreadsheet software. Whilst a simple graph can be
356 easily created with either option, the potential for more advanced analytics such as regression or
357 even advanced machine learning techniques are greatly limited in spreadsheet software, compared
358 with enormous potential within coding environments. A range of coding languages available such as
359 R and Python, and whilst most languages allow for powerful data analysis, some can be easier to
360 learn than others. Whilst coding is initially more difficult, an investment in learning how to write
361 efficient and powerful code will pay dividends in the long term for a clinician keen to develop more
362 data skills (Figure 10).

363 There is also the advantage that both R and Python are freely available, with no licensing costs at all,
364 in contrast to many proprietary statistics packages and spreadsheet programs, where the software
365 must be purchased. There are also freely available online tools such as the Nottingham Herd Health
366 Toolkit (www.nottingham.ac.uk/herdhealthtoolkit), which allows upload of farm-based data with a
367 predefined output for those less keen to write their own code or taking their first steps in this area.

368

369 **Box 6: Tidy data**

370 Data often comes to veterinarians in a variety of formats. The principles of “Tidy data” (Grolemund
371 and Wickham, 2016) is extremely useful in data analysis, as it ensures that analytical techniques are
372 easily repeatable between datasets, and aids in creating an efficient data analysis workflow.

373 “Tidy data” essentially describes a dataset that involves one row per observation, and one column
374 per variable. Data being recorded on farm will often come in a wide range of formats, and it is not
375 often tidy. A sensible first step in data analysis is to ensure your data fits the “Tidy data” framework
376 before attempting more advanced analytics.

377

378 **Box 7: Calf growth rate graphs in Excel**

379 Monitoring growth rates in calves can be extremely valuable, in benchmarking farms and identifying
380 performance issues to ensure maximum efficiency in terms of health and productivity.

381 Creating box plots is relatively straightforward in Excel (Figure 12). The disadvantage is that to repeat
382 this analysis for a farm you generally have to go through all these steps every single time. By using a
383 coding language, you simply re-run the code you have already written to create a bespoke graph in a
384 fraction of a second. Note also that this dataset is not in the “Tidy data” format (See box 6, figure 11)
385 and would require a lot of manual work in Excel to reorganise to create other graphs, for example
386 scatter plots.

387 Box 8: Calf growth rate graphs in R

388

389 Whilst there is not much difference between the plot in Figure 13 and the excel plot in Figure 12, the
390 real advantage comes when this process needs to be repeated for future analysis. Next time you
391 wish to plot this graph for a farm, all that is needed is to change the file path to the new .csv file and
392 rerun the code. Whilst the excel graph would need to be constructed from scratch for a new set of
393 data, this code takes under a second to import data, calculate DLWG and produce the graph shown.
394 A small amount of effort learning how to write 13 lines of code has resulted in a graph that can be
395 produced for any farm that records calf weights in less than a second.

396 This graph can be pasted into a word document to add comments etc before sending as part of a
397 report. There is also the option in R to create an automated report that includes whichever graphs
398 and data analytics you require, enabling a practice to write detailed reports for multiple farms, or
399 even also benchmark multiple farms against each other in a matter of seconds once the code has
400 initially been written. Further coding examples are provided in *Appendix i*.

401 In this example, it seems that bull calves generally grow faster than heifers. If there are differences
402 in management between bull and heifer calves this is worth exploring more on farm.

403

404 **Summary**

405 Data skills are becoming increasingly important in farm animal population medicine. A wide range of
406 “big data” sources are available on farms, and the ability to utilise this data to improve animal health
407 outcomes can be a useful tool to the modern production animal veterinarian.

408 Whilst data skills can be initially be challenging to learn, a relatively small investment in time spent
409 learning can result in dramatically shorter times analysing farm outcomes and writing reports in the
410 future when analysing data for multiple farms. The statistical and analytic techniques discussed in
411 this article are freely available and should provide veterinarians with new tools to prevent disease at
412 a population level.

413

414 Further resources

415 The following resources are available for free:

416 For a tutorial in downloading and getting started in R see *Appendix i*

417 Further coding examples for common herd health problems are provided in *Appendix ii*

418 Epitools epitools.ausvet.com.au

419 Nottingham herd health toolkit nottingham.ac.uk/herdhealthtoolkit

420 R for data science book r4ds.had.co.nz YaRrr! The Pirate's Guide to R bookdown.org/ndphillips/YaRrr

421 Teacups Giraffes and Statistics tinystats.github.io/teacups-giraffes-and-statistics

422 Stack overflow stackoverflow.com

423 Quick-R statmethods.net

424

425 **References**

426 Amrhein, V., Greenland, S., McShane, B., 2019. Scientists rise up against statistical significance.
427 Nature. <https://doi.org/10.1038/d41586-019-00857-9>

428 Grolemund, G., Wickham, H., 2016. R for Data Science. O'Reilly.

429 Hermans, K., Opsomer, G., Waegeman, W., Moerman, S., De Koster, J., Eetvelde, V., Ranst, V.,
430 Hostens, M., 2018. Interpretation and visualisation of data from dairy herds 40, 195–203.
431 <https://doi.org/10.1136/inp.k2166>

432 Hudson, C., Kaler, J., Down, P., 2018. Using big data in cattle practice. In Pract. 40, 396–410.
433 <https://doi.org/10.1136/inp.k4328>

434 Ioannidis, J., 2005. Why Most Published Research Findings Are False. PLoS Med. 2, e124.
435 <https://doi.org/10.1371/journal.pmed.0020124.t004>

436 Overton, M.W., 2009. Using Reproductive Records: Basics of Monitoring.

437 Team, R.C., 2020. R: A Language and Environment for Statistical Computing.

438 Wickham, H., 2016. ggplot2: Elegant Graphics for Data Analysis.

439 Woodward, H., Cobb, K., Remnant, J., 2019. The future of cattle veterinary practice: Insights from a
440 qualitative study. Vet. Rec. 185, 205. <https://doi.org/10.1136/vr.105321>

441

442

443 Figure legends:
444
445 Figure 1: Types of data commonly encountered in farm animal data
446
447 Figure 2: Commonly used summary statistics
448
449 Figure 3: Somatic cell counts for 100 dairy cows showing a “right tail” with a small number of very
450 high cell count animals resulting in an increased mean relative to the median.

451 Figure 4: A variety of graphical options to describe the calf performance
452
453 Figure 5: Plotting data in Excel allows for the simple analysis of weight by age.
454
455 Figure 6. By colouring by breed type, it is clear that in addition to the effect of age, there is also an
456 effect of breed type on weight.

457 Figure 7: Daily liveweight gain of 500 calves from two farms. This dataset has a “normal” or
458 “gaussian” distribution, with a mean growth rate of 0.8kg/d (dashed line), however the calves from
459 farm A (red) have a far more consistent growth rate (standard deviation 0.1kg/d) compared with the
460 wider spread in growth rates of farm B (blue, standard deviation 0.2kg/d). 95% of the data is spread
461 between 1.96 standard deviations either side of the mean for both farms (0.196kg/d for farm A and
462 0.392kg/d for farm B).

463 Figure 8: Epitools can be used to calculate statistical significance between two treatment groups
464
465 Figure 9: Epitools can be used to perform sample size calculations. Variance is the standard deviation
466 squared (see Box 3).
467
468 Figure 10: Whilst statistical coding using R/Python initially takes longer, there are significant
469 productivity benefits to more advanced users in both analytical performance and speed.
470
471 Figure 11: An example of “Tidy data”
472
473 Figure 12: Excel allows relatively straightforward plotting of calf performance by sex.
474
475 Figure 13: An example of calf performance analysis using a programming language such as R.

471 Questions

472 What statistical test could be used to compare two continuous parametric groups with normal
473 distribution?

- 474 • T-test
- 475 • Mann-Whitney test
- 476 • ANOVA
- 477 • Chi-squared test

478 With normally distributed data how many standard deviations either side of the mean would contain
479 95% of the data?

- 480 • 0.96
- 481 • 1.96
- 482 • 2
- 483 • 3.96

484 What type of model might be used to analyse the effect of age on calf weight?

- 485 • Logistic regression
- 486 • Linear regression
- 487 • ANOVA
- 488 • Mann-Whitney

489 In a regression model, what term describes the percentage of variation explained by the model?

- 490 • R^2
- 491 • MAE
- 492 • RMSE
- 493 • Coefficient

494 What advantages do statistical coding languages provide over spreadsheets when performing data
495 analysis?

- 496 • Speed
- 497 • Repeatability
- 498 • Performance
- 499 • All of the above

500