# Model Selection

Peter von Rohr

**Last week:** 20.04.2020
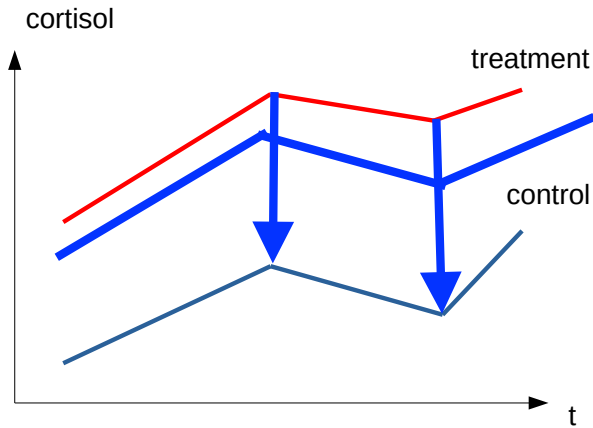**2 sides in breeding program:**
**1. economic evaluation**
**2. prediction of breeding values using statistical modelling**

# Why Statistical Modelling?

Some people believe, they do not need statistics. For them it is enough to look at a diagram
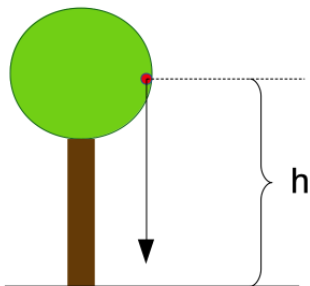
# Statistical Modelling Because . . .

Two types of dependencies between physical quantities

1. deterministic ⟶ **no sources of uncertainty**
2. stochastic

# Deterministic Versus Stochastic



**deterministic**

**stochastic**

h

Law of gravity

**time t that it takes the apple to fall can be computed without sources of uncertainty**

genotype

+

environment

**many sources of uncertainty**
**> genes important**
**> post-translational processes**
**> env**

phenotype

# Statistical Model

- stochastic systems contains many sources of uncertainty
- statistical models can handle uncertainty
- components of a statistical model
  - response variable $y$ **observed phenotypes**
  - predictor variables $x_1, x_2, \ldots, x_k$ **fixed effects and random breeding values**
  - error term $e$
  - function $m(x)$

  **statistical model**

# How Does A Statistical Model Work?

- predictor variables $x_1, x_2, \ldots, x_k$ are transformed by function $m(x)$ to explain the response variable $y$
- uncertainty is captured by error term.
- as a formula, for observation $i$

$$y_i = m(x_i) + e_i$$

**x_i predictors are taken as input to the function m()**

# Which function $m(x)$?

- class of functions that can be used as $m(x)$ is infinitely large
- restrict to linear functions of predictor variables

# Which predictor variables?

▶ Question, about which predictor variables to use is answered by model selection

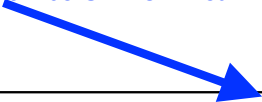**predictor variables: x_1, x_2, ..., x_k**

# Why Model Selection

**with respect to explaining the differences in the response variables**

- Many predictor variables are available
- Are all of them `relevant`?
- What is the meaning of `relevant` in this context?

## Example Dataset

**random numbers = non-meaningful**

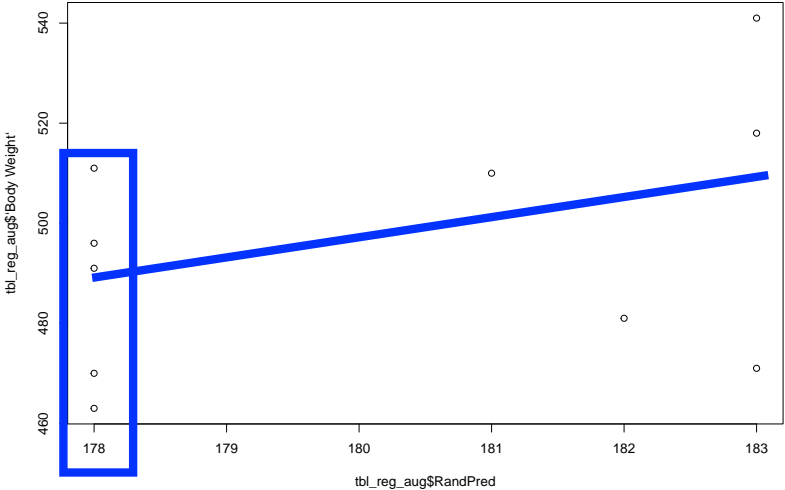| Animal | Breast Circumference | Body Weight | RandPred |
|--------|---------------------|-------------|----------|
| 1 | 176 | 471 | 183 |
| 2 | 177 | 463 | 178 |
| 3 | 178 | 481 | 182 |
| 4 | 179 | 470 | 178 |
| 5 | 179 | 496 | 178 |
| 6 | 180 | 491 | 178 |
| 7 | 181 | 518 | 183 |
| 8 | 182 | 511 | 178 |
| 9 | 183 | 510 | 181 |
| 10 | 184 | 541 | 183 |

**predictor = meaningful**

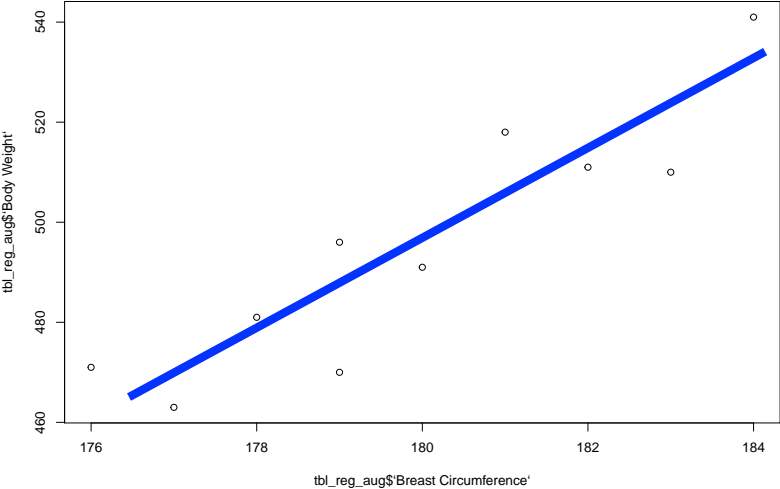**response variable $y_i$**

# No Relevance of Predictors



**random pattern**

# Relevance of Predictors



**pattern: points all grouped around regression line**

# Fitting a Regression Model

```
## 
## Call:
## lm(formula = `Body Weight` ~ RandPred, data = tbl_reg_aug)
## 
## Residuals:
##     Min      1Q  Median      3Q     Max
## -35.574 -20.200   7.236  11.519  34.426
## 
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) -236.775    608.880  -0.389    0.708
## RandPred       4.062      3.379   1.202    0.264
## 
## Residual standard error: 24.27 on 8 degrees of freedom
## Multiple R-squared:  0.153,  Adjusted R-squared:  0.04716
## F-statistic: 1.445 on 1 and 8 DF,  p-value: 0.2636
```

**Low Adj. R-sq ==>**
**Model does not explain**
**a high rate of variation**
**of response variables**

**std. error has about the same magnitude as the estimate ==> problem**

# Fitting a Regression Model II

```
##
## Call:
## lm(formula = `Body Weight` ~ `Breast Circumference`, data = tbl_reg_aug)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -17.3941  -6.5525  -0.0673   9.3707  13.2594
##
## Coefficients:
##                         Estimate Std. Error t value Pr(>|t|)
## (Intercept)            -1065.115    255.483  -4.169 0.003126 **
## `Breast Circumference`     8.673      1.420   6.108 0.000287 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 11.08 on 8 degrees of freedom
## Multiple R-squared:  0.8234, Adjusted R-squared:  0.8014
## F-statistic: 37.31 on 1 and 8 DF,  p-value: 0.000287
```

## Multiple Regression

```
## 
## Call:
## lm(formula = `Body Weight` ~ `Breast Circumference` + RandPred,
##     data = tbl_reg_aug)
## 
## Residuals:
##      Min      1Q   Median      3Q      Max
## -13.1363  -3.0404   0.7548   4.3149  14.3068
## 
## Coefficients:
##                        Estimate Std. Error t value Pr(>|t|)
## (Intercept)            -1492.865    295.360  -5.054 0.001473 **
## `Breast Circumference`     8.304      1.202   6.909 0.000229 ***
## RandPred                   2.742      1.306   2.100 0.073839 .
## 
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## Residual standard error: 9.278 on 7 degrees of freedom
## Multiple R-squared:  0.8917,	Adjusted R-squared:  0.8607
## F-statistic: 28.81 on 2 and 7 DF,  p-value: 0.0004183
```

# Which model is better?

Why not taking all predictors?

▶ Additional parameters must be estimated from data
▶ Predictive power decreased with too many predictors (cannot be shown for this data set, because too few data points)
▶ Bias-variance trade-off

# Bias-variance trade-off

▶ Assume, we are looking for optimum prediction

**Example with BW:**
**q_1 = {BC, RandPred}**
**q_2 = {BC}**

$$s_i = \sum_{r=1}^{q} \hat{\beta}_{j_r} x_{ij_r}$$

with $q$ relevant predictor variables

▶ Average mean squared error of prediction $s_i$

$$MSE = n^{-1} \sum_{i=1}^{n} E\left[(m(x_i) - s_i)^2\right]$$

where $m(.)$ denotes the linear function of the unknown true model.

# Bias-variance trade-off II

**the more predictors x are included in s_i, the smaller the bias will be**

▶ MSE can be split into two parts

**variance**

$$MSE = n^{-1} \sum_{i=1}^{n} \left(E\left[s_i\right] - m(x_i)\right)^2 + n^{-1} \sum_{i=1}^{n} var(s_i)$$

where $n^{-1} \sum_{i=1}^{n} \left(E\left[s_i\right] - m(x_i)\right)^2$ is called the squared **bias**

▶ Increasing $q$ leads to reduced bias but increased variance ($var(s_i)$)
▶ Hence, find $s_i$ such that MSE is minimal
▶ Problem: cannot compute MSE because $m(.)$ is not known

$\rightarrow$ estimate MSE

# Mallows $C_p$ statistic

**M: y_i = b_0 + b_1 * breast_circumference**

- For a given model $\mathcal{M}$, $SSE(\mathcal{M})$ stands for the residual sum of squares.
- MSE can be estimated as

**number of predictors in model M**

$$\widehat{MSE} = n^{-1}SSE(\mathcal{M}) - \hat{\sigma}^2 + 2\hat{\sigma}^2|\mathcal{M}|/n$$

where $\hat{\sigma}^2$ is the estimate of the error variance of the full model, $SSE(\mathcal{M})$ is the residual sum of squares of the model $\mathcal{M}$, $n$ is the number of observations and $|\mathcal{M}|$ stands for the number of predictors in $\mathcal{M}$

$$C_p(\mathcal{M}) = \frac{SSE(\mathcal{M})}{\hat{\sigma}^2} - n + 2|\mathcal{M}|$$

**Values of Mallow C_p should be as small as possible**

# Searching The Best Model

- Exhaustive search over all sub-models might be too expensive
- For $p$ predictors there are $2^p - 1$ sub-models
- With $p = 16$, we get $6.5535 \times 10^4$ sub-models

$\rightarrow$ step-wise approaches

**2 ways to do step-wise approach:**
  **1. forward selection**
  **2. backward elimination**

# Forward Selection

**M_0: y_i = b_0 + e_i ==> just use an intercept b_0 ==> compute C_p**
  **Step 2: Question: would it be better to include any of the available pred?**
  **==> Constructing model M_1: In our example with BW ==>**
  **Should M_1 contain BC or RandPred as its preditor?**

1. Start with smallest sub-model $\mathcal{M}_0$ as current model
2. Include predictor that reduces SSE the most to current model
3. Repeat step 2 until all predictors are chosen

$\rightarrow$ results in sequence $\mathcal{M}_0 \subseteq \mathcal{M}_1 \subseteq \mathcal{M}_2 \subseteq \ldots$ of sub-models

4. Out of sequence of sub-models choose the one with minimal $C_p$

**For each sub-model M_0, M_1, M_2, … we have computed C_p**
**From all submodel select the one with lowest C_p value**
**This will be the best model.**

# Backward Selection

**M_0: full model, for example of BW:**
**M_0: y_i = b_0 + b_1 * breast_circum. + b_2 * randpred + e_i**

**Step 2: exlude predictors from M_0**

1. Start with full model $\mathcal{M}_0$ as the current model
2. Exclude predictor variable that increases SSE the least from current model
3. Repeat step 2 until all predictors are excluded (except for intercept)

$\rightarrow$ results in sequence $\mathcal{M}_0 \supseteq \mathcal{M}_1 \supseteq \mathcal{M}_2 \supseteq \ldots$ of sub-models

4. Out of sequence choose the one with minimal $C_p$

# Considerations

- Whenever possible, choose **backward** selection, because it leads to better results
- If $p \geq n$, only forward is possible, but then consider LASSO

**backward eleminiation is not possible because full model cannot be fitted**

# Alternative Selection Criteria

**AIC: Akaike Information Criterion**
**BIC: Bayesian Information Criterion**

- ▶ AIC or BIC, requires distributional assumptions.
- ▶ AIC is implemented in `MASS::stepAIC()`
- ▶ Adjusted $R^2$ is a measure of goodness of fit, but sometimes is not conclusive when comparing two models
- ▶ Try in exercise

**R-package: olsrr, MASS uses just AIC**

**For Genetic Evaluation:**
**\* In our database: many different predictors are available for a given trait**
**\* Do model selection to find good balance between bias and variance**
**\* Model selection is used to identify fixed effects in our models to estimated variance components and to predict breeding values.**