# Fixed Linear Effects Models

$\rightarrow$ general form

Peter von Rohr

regression

2022-03-07

- regression on dummy variables
- models of non full rank

rank of X

# Extension of Dataset on Body Weight

*response y* (annotation)

*additional predictor* (annotation)

*cost+1 column in X* (annotation)

| Animal | BC | Body Weight | BCS | HEI | Breed |
|---|---|---|---|---|---|
| 1 | 176 | 471 | 5.0 | 161 | Angus |
| 2 | 177 | 463 | 4.2 | 121 | Angus |
| 3 | 178 | 481 | 4.9 | 157 | Simmental |
| 4 | 179 | 470 | 3.0 | 165 | Angus |
| 5 | 179 | 496 | 6.8 | 136 | Simmental |
| 6 | 180 | 491 | 4.9 | 123 | Simmental |
| 7 | 181 | 518 | 4.4 | 163 | Limousin |
| 8 | 182 | 511 | 4.4 | 149 | Limousin |
| 9 | 183 | 510 | 3.5 | 143 | Limousin |
| 10 | 184 | 541 | 4.7 | 130 | Limousin |

# Include Breed into Model

predictors are columns in $X$ ; components of $X$ are numeric

$$X = \begin{bmatrix} 1 & 176 & \cdots & Angus \\ 1 & 178 & & \\ 1 & & & \\ \vdots & \vdots & & \\ 1 & & & \end{bmatrix}$$

- ▶ Breed has an influence on body weight
- ▶ Predictor variables must be numeric
- ▶ Breed must be converted to numeric code
- ▶ Assignment of codes to breeds is rather arbitrary

# Breed Codes

$$X = \begin{bmatrix} 1 & - & \cdots & 1 \\ 1 & & & 1 \\ & & & 3 \end{bmatrix} \xleftarrow{\text{vce}}$$

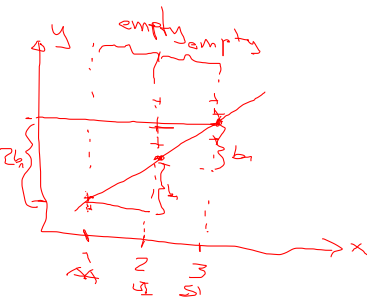| Code | Breed |
|------|-------|
| 1 | Angus — less BW |
| 2 | Limousin |
| 3 | Simmental |

Assignment    arbitrary

# Modelling Effect of Breed

$$X = \begin{bmatrix} 1 & 1 \\ 1 & 1 \\ \vdots & 3 \\ 1 & 2 \end{bmatrix}$$

→ numeric code

▶ Simplification: "breed" is the only predictor
▶ Expected body weight ($y_i$) for animal $i$

$$E(y_i) = b_0 + b_1 x_i$$

→ reg. coef

↓ intercept

→ breed code

# Problems

- Nothing wrong with previous model
- But the following relations might give a hint to some problems

Expected value for Body Weight

$$E(\text{BW Angus}) = b_0 + 1 \cdot b_1$$
$$E(\text{BW Limousin}) = b_0 + 2b_1$$
$$E(\text{BW Simmental}) = b_0 + 3b_1$$

$E[y_i] = b_0 + b_1 x_i$

breed code

This means, for example, that

$b_0 + 2b_1 \qquad - [b_0 + b_1] = b_0 - b_0 + 2b_1 - b_1 = b_1$

$$E(\text{BW Limousin}) - E(\text{BW Angus}) =$$

$b_0 + 3b_1 \qquad - [b_0 + 2b_1] = b_1$

$$E(\text{BW Simmental}) - E(\text{BW Limousin})$$
$$E(\text{BW Simmental}) - E(\text{BW Angus}) =$$
$$2\left[E(\text{BW Limousin}) - E(\text{BW Angus})\right]$$

# Consequences

*(handwritten annotations at top)*

Regression → { y — BW
                      continuous
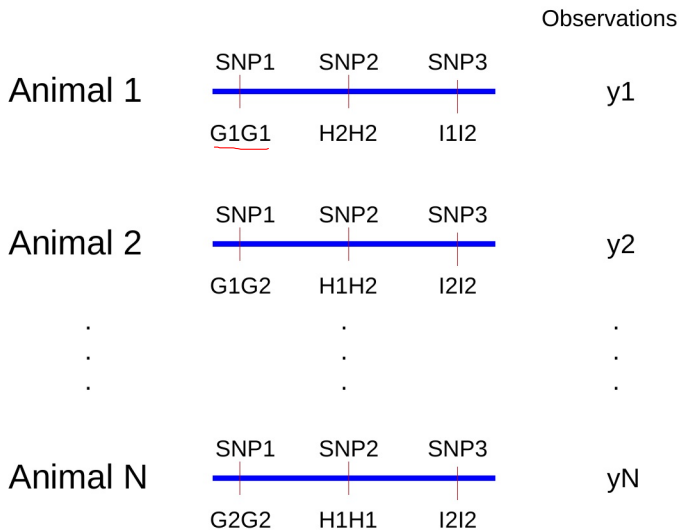              X — ABC, BC, HE1 }

- ▶ Allocation of numerical codes imposes relations between expected values
- ▶ Relations might be unreasonable
- ▶ Regression analysis only yields estimates for $b_0$ and $b_1$, effects of other breeds are determined
- ▶ Conclusion: regression on numerical codes of discrete variables are in most cases unreasonable
- ▶ Exception: Estimation of marker effects

*(handwritten annotations)* Breed ≠ 2.75

# Linear Regression Analysis for Genomic Data

For a single SNP: Assuming G1 to be the favorable allele, the numeric code for a given genotype will be the number of G1 alleles



only assuming add. Model

$G_2G_2$  $G_1G_2$  $G_1G_1$

Code  0  1  2