

MASM22/FMSN30: Linear and Logistic Regression, 7.5 hp

FMSN40: ... with Data Gathering, 9 hp

Lecture 1, spring 2024

Linear regression: assumptions and estimates

<https://canvas.education.lu.se/courses/xxx>

Mathematical Statistics / Centre for Mathematical Sciences
Lund University

18/3-24

Introduction

Simple linear regression

Multiple linear regression

Assumptions

Model

Multivariate distributions

Covariance and correlation

Covariance matrix

Bivariate distributions

Regression cont'd.

Matrix formulation

Least squares estimation

Properties of estimates

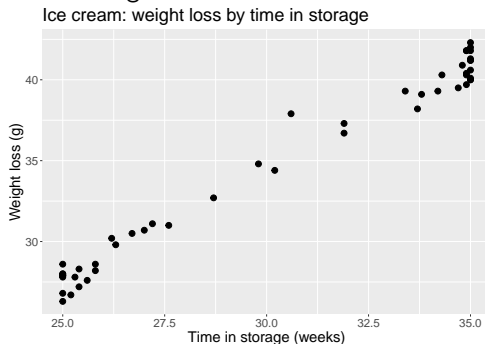
Example: Ice cream

Simple Linear Regression

We measure two variables, x and Y . How does the value of Y depend on the value of x ? Is there a linear relationship? How can we estimate this relationship using observed data?

Example: Ice cream

An ice cream manufacturer suspects that storing ice cream at low temperatures leads to weight loss.

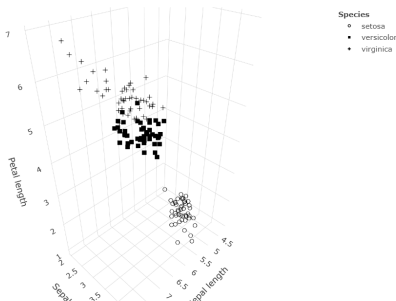


Multiple linear regression

We measure several variables, x_1, \dots, x_p , and Y . How does the value of Y depend on the values of x_1, \dots, x_p ?

Example: Iris

The petal length depends on sepal length, sepal width and species.



Basic assumptions

- ▶ Y : **continuous** dependent variable, “response” or “outcome”, assumed **random**.
- ▶ x_1, \dots, x_p : explanatory variables, “covariates”; assumed **non-random**.
- ▶ We *hypothesize* that Y has a linear relationship with x_1, \dots, x_p , on average, and follows the **linear model**:

$$E(Y) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

- ▶ $\beta_0, \beta_1, \dots, \beta_p$: unknown parameters, assumed **non-random**.
- ▶ β_0 = intercept; $E(Y)$ when $x_1 = \dots = x_p = 0$,
- ▶ β_1, \dots, β_p = slopes in the corresponding x -directions; the additive change in $E(Y)$ when the corresponding x -variable is increased by 1 unit and the others are held fixed.

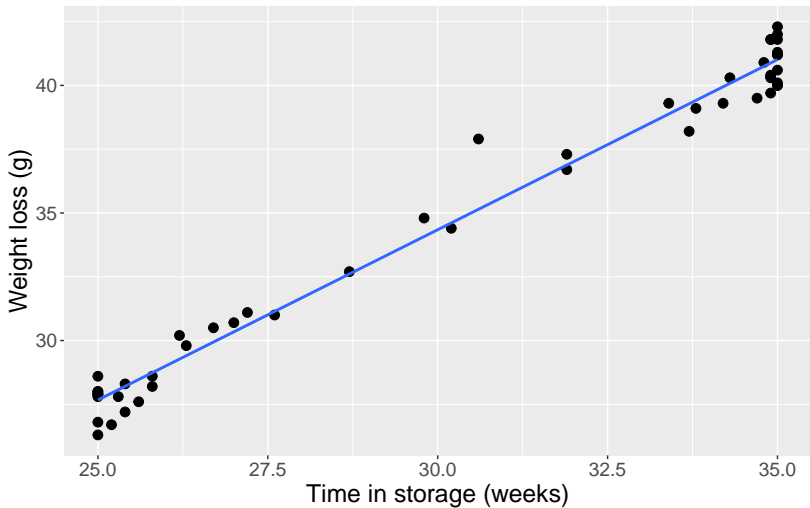
Model: multiple linear regression

- ▶ We denote with Y_i , where $i = 1, \dots, n$, the i th observation from a set of n measurements of Y .
- ▶ We denote with x_{ij} , where $i = 1, \dots, n$ and $j = 1, \dots, p$, the corresponding i th observation of the j th x -variable.
- ▶ The model for a generic observation Y_i is

$$Y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \epsilon_i, \quad i = 1, \dots, n$$

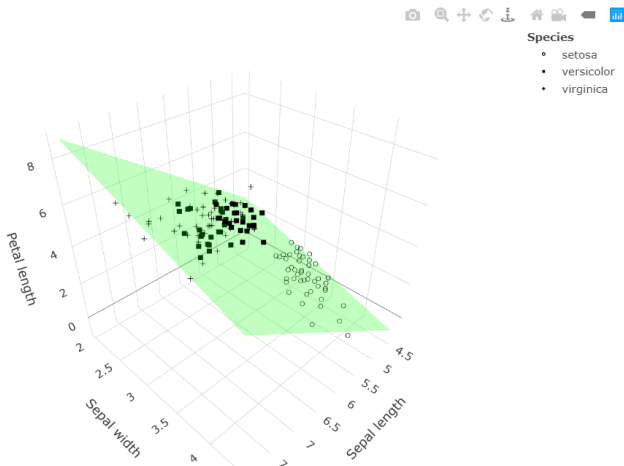
where ϵ_i is the “measurement error” which contains all the random variation not explained by the linear model.

Ice cream: weight loss by time in storage



data and fitted line

Iris: petal length as function of sepal length and width, $p = 2$



Assumptions for the measurement error

Besides linearity, we also assume for all $i = 1, \dots, n$:

$$\begin{array}{ll} E(\epsilon_i) = 0 & E(Y_i) = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} = \mu_i \\ V(\epsilon_i) = \sigma^2 & V(Y_i) = \sigma^2 \\ \epsilon_i \sim N(0, \sigma^2) & Y_i \sim N(\mu_i, \sigma^2) \\ \epsilon_i \text{ are pairwise independent} & Y_i \text{ are pairwise independent} \end{array}$$

A note on notation

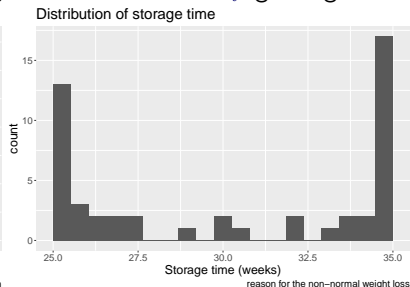
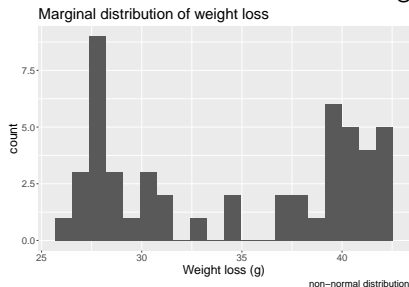
Strictly speaking, we assume properties of Y_i *conditional* on $X_j = x_{ij}$:

$$\begin{aligned} E(Y_i \mid X_1 = x_{i1}, \dots, X_p = x_{ip}) &= \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} = \mu_i, \\ V(Y_i \mid X_1 = x_{i1}, \dots, X_p = x_{ip}) &= \sigma^2, \\ Y_i \mid X_1 = x_{i1}, \dots, X_p = x_{ip} &\sim N(\mu_i, \sigma^2) \end{aligned}$$

We consider this notation as implicit and always write $E(Y_i)$, etc., in place of $E(Y_i \mid X_1 = x_{i1}, \dots)$... except on the next few slides...

Warning!

Our assumptions imply that the conditional distribution of $Y_i | X = x_i$ is Normal but we don't know the marginal distribution of Y_i ignoring X !



The marginal distribution of Y is clearly not Normal but this is just due to the strange distribution of the x -values.

To assess normality you should instead inspect **residuals** (introduced later), which means you cannot always see that your model will be wrong (or right!) until after you have fitted it!

Introduction

Simple linear regression

Multiple linear regression

Assumptions

Model

Multivariate distributions

Covariance and correlation

Covariance matrix

Bivariate distributions

Regression cont'd.

Matrix formulation

Least squares estimation

Properties of estimates

Example: Ice cream

Covariance and correlation

- ▶ The **covariance**, $C(U, W)$, between two random variables, U and W , measures the **linear** dependance between them:

$$C(U, W) \stackrel{\text{Def.}}{=} E((U - E(U))(W - E(W)))$$

Note: $C(W, U) = C(U, W)$.

- ▶ The **correlation**, $\rho(U, W)$, is a unitless, bounded, version:

$$\rho(U, W) \stackrel{\text{Def.}}{=} \frac{C(U, W)}{\sqrt{V(U)V(W)}}, \quad -1 \leq \rho \leq 1$$

- ▶ The **variance**: $V(U) \stackrel{\text{Def.}}{=} E((U - E(U))^2) = C(U, U)$

$$\rho(U, U) = \frac{C(U, U)}{\sqrt{V(U)V(U)}} = \frac{V(U)}{V(U)} \equiv 1$$

Covariance matrix

The **Covariance matrix** for U and W , $\text{Var}\left(\begin{bmatrix} U \\ W \end{bmatrix}\right)$, is defined as

$$\text{Var}\left(\begin{bmatrix} U \\ W \end{bmatrix}\right) = \begin{bmatrix} V(U) & C(U, W) \\ C(W, U) & V(W) \end{bmatrix}$$

The **Correlation matrix** for U and W , $\text{Corr}\left(\begin{bmatrix} U \\ W \end{bmatrix}\right)$, is defined as

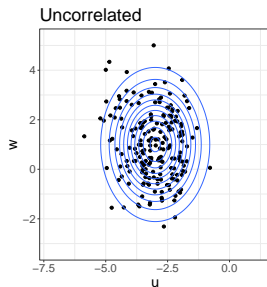
$$\text{Corr}\left(\begin{bmatrix} U \\ W \end{bmatrix}\right) = \begin{bmatrix} 1 & \rho(U, W) \\ \rho(W, U) & 1 \end{bmatrix}$$

Uncorrelated vs independant

- ▶ U and W are **uncorrelated** $\stackrel{\text{Def.}}{\Leftrightarrow} C(U, W) = 0$.

Thus $\text{Var}\left(\begin{bmatrix} U \\ W \end{bmatrix}\right)$ will be a diagonal matrix.

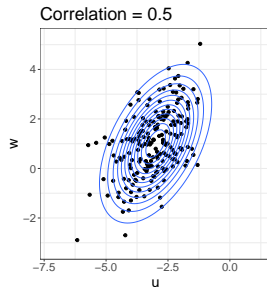
- ▶ U and W are **independant** $\Rightarrow U$ and W are **uncorrelated**.
- ▶ If U and W are normally distributed:
 U and W are **independant** $\Leftrightarrow U$ and W are **uncorrelated**.



$$E\left(\begin{bmatrix} U \\ W \end{bmatrix}\right) = \begin{bmatrix} -3 \\ 1 \end{bmatrix}$$

$$\text{Var}\left(\begin{bmatrix} U \\ W \end{bmatrix}\right) = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$$

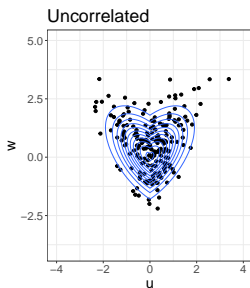
Uncorrelated and independent
(Bivariate normal)



$$E\left(\begin{bmatrix} U \\ W \end{bmatrix}\right) = \begin{bmatrix} -3 \\ 1 \end{bmatrix}$$

$$\text{Var}\left(\begin{bmatrix} U \\ W \end{bmatrix}\right) = \begin{bmatrix} 1 & 0.71 \\ 0.71 & 2 \end{bmatrix}$$

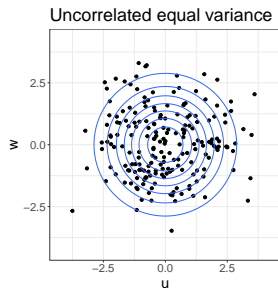
Correlated (and dependant)
(Bivariate normal)



$$E\left(\begin{bmatrix} U \\ W \end{bmatrix}\right) = \begin{bmatrix} 0 \\ 0.80 \end{bmatrix}$$

$$\text{Var}\left(\begin{bmatrix} U \\ W \end{bmatrix}\right) = \begin{bmatrix} 1 & 0 \\ 0 & 1.36 \end{bmatrix}$$

Uncorrelated but dependant
(Not bivariate normal)



$$E\left(\begin{bmatrix} U \\ W \end{bmatrix}\right) = \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \mathbf{0}$$

$$\text{Var}\left(\begin{bmatrix} U \\ W \end{bmatrix}\right) = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix} = 2\mathbf{I}$$

$N_2(\mathbf{0}, 2\mathbf{I})$
(Bivariate normal)

Introduction

Simple linear regression

Multiple linear regression

Assumptions

Model

Multivariate distributions

Covariance and correlation

Covariance matrix

Bivariate distributions

Regression cont'd.

Matrix formulation

Least squares estimation

Properties of estimates

Example: Ice cream

Multiple linear regression with matrices

$$\mathbf{Y} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}, \mathbf{X} = \begin{bmatrix} 1 & x_{11} & \cdots & x_{1p} \\ 1 & x_{21} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{np} \end{bmatrix}, \boldsymbol{\epsilon} = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}, \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix}$$

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} = \begin{bmatrix} 1 \cdot \beta_0 + x_{11} \cdot \beta_1 + \dots + x_{1p} \cdot \beta_p \\ 1 \cdot \beta_0 + x_{21} \cdot \beta_1 + \dots + x_{2p} \cdot \beta_p \\ \vdots \\ 1 \cdot \beta_0 + x_{n1} \cdot \beta_1 + \dots + x_{np} \cdot \beta_p \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

$$\sim N_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$$

(n -dimensional multivariate normal distribution)

We assume $\epsilon \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I})$ where $E(\epsilon) = \mathbf{0} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$ and

the covariance matrix $\text{Var}(\epsilon)$ is given by

$$\begin{aligned} \text{Var}(\epsilon) &\stackrel{\text{Def.}}{=} \begin{bmatrix} V(\epsilon_1) & C(\epsilon_1, \epsilon_2) & \dots & C(\epsilon_1, \epsilon_n) \\ C(\epsilon_2, \epsilon_1) & V(\epsilon_2) & \dots & C(\epsilon_2, \epsilon_n) \\ \vdots & \vdots & \ddots & \vdots \\ C(\epsilon_n, \epsilon_1) & C(\epsilon_n, \epsilon_2) & \dots & V(\epsilon_n) \end{bmatrix} \\ &= \begin{bmatrix} \sigma^2 & 0 & \dots & 0 \\ 0 & \sigma^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma^2 \end{bmatrix} = \sigma^2 \cdot \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{bmatrix} = \sigma^2 \mathbf{I} \end{aligned}$$

Least squares estimates: simple linear without matrices

Find β_0, β_1 that minimize the loss function

$$Q(\beta_0, \beta_1) = \sum_{i=1}^n (Y_i - E(Y_i))^2 = \sum_{i=1}^n (Y_i - (\beta_0 + \beta_1 x_i))^2.$$

Partial derivatives

Find the minimum by solving the linear equation system

$$\frac{\partial Q}{\partial \beta_0} = -2 \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 x_i) = 0 \Leftrightarrow n\beta_0 + \beta_1 \sum_{i=1}^n x_i = \sum_{i=1}^n Y_i$$

$$\frac{\partial Q}{\partial \beta_1} = -2 \sum_{i=1}^n x_i (Y_i - \beta_0 - \beta_1 x_i) = 0 \Leftrightarrow \beta_0 \sum_{i=1}^n x_i + \beta_1 \sum_{i=1}^n x_i^2 = \sum_{i=1}^n x_i Y_i$$

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{x},$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y})}{\sum_{i=1}^n (x_i - \bar{x})^2}.$$

Least squares estimates: with matrices

Find β that minimizes $Q(\beta) = (\mathbf{Y} - \mathbf{X}\beta)^\top (\mathbf{Y} - \mathbf{X}\beta)$.

Partial derivatives

Expand $Q(\beta)$ and use $\frac{\partial \beta^\top \mathbf{A}}{\partial \beta} = \mathbf{A}$, $\frac{\partial \mathbf{A} \beta}{\partial \beta} = \mathbf{A}^\top$ when β is a column vector and \mathbf{A} is a matrix.

$$\begin{aligned} Q(\beta) &= \mathbf{Y}^\top \mathbf{Y} - \mathbf{Y}^\top \mathbf{X} \beta - \beta^\top \mathbf{X}^\top \mathbf{Y} + \beta^\top \mathbf{X}^\top \mathbf{X} \beta \\ \frac{\partial Q(\beta)}{\partial \beta} &= \mathbf{0} - (\mathbf{Y}^\top \mathbf{X})^\top - \mathbf{X}^\top \mathbf{Y} + \mathbf{X}^\top \mathbf{X} \beta + (\mathbf{X}^\top \mathbf{X})^\top \beta \\ &= -2\mathbf{X}^\top \mathbf{Y} + 2\mathbf{X}^\top \mathbf{X} \beta = \mathbf{0}. \end{aligned}$$

The solution satisfies the normal equations: $\mathbf{X}^\top \mathbf{X} \beta = \mathbf{X}^\top \mathbf{Y}$.

$$\begin{bmatrix} n & \sum_{i=1}^n x_{i1} & \cdots & \sum_{i=1}^n x_{ip} \\ \sum_{i=1}^n x_{i1} & \sum_{i=1}^n x_{i1}^2 & \cdots & \sum_{i=1}^n x_{i1} x_{ip} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_{i=1}^n x_{ip} & \sum_{i=1}^n x_{i1} x_{ip} & \cdots & \sum_{i=1}^n x_{ip}^2 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^n Y_i \\ \sum_{i=1}^n x_{i1} Y_i \\ \vdots \\ \sum_{i=1}^n x_{ip} Y_i \end{bmatrix}$$

- **Parameter estimates:** $\hat{\beta}$ is the solution to the normal equations:

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

- **Predicted values** = fitted line (plane):

$$\hat{Y}_i = \hat{E}(Y_i) = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \dots + \hat{\beta}_p x_{ip} \quad \hat{\mathbf{Y}} = \mathbf{X} \hat{\beta}$$

- **Residuals:** the difference between observations and predictions:

$$e_i = Y_i - \hat{Y}_i \quad \mathbf{e} = \mathbf{Y} - \hat{\mathbf{Y}} = \mathbf{Y} - \mathbf{X} \hat{\beta}$$

- **Residual variance:** s^2 is an estimate of σ^2 , the variance of the error, a measure of the “residual variability” unexplained by the model.

$$\hat{\sigma}^2 = s^2 = \frac{Q(\hat{\beta})}{n - (p + 1)} = \frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{n - (p + 1)} = \frac{\mathbf{e}^T \mathbf{e}}{n - (p + 1)}$$

Note: $p + 1$ is the total number of β -parameters in the model.

Properties of parameter estimates

Note: For a constant matrix \mathbf{A} and a random matrix \mathbf{Y} we have $E(\mathbf{AY}) = \mathbf{A}E(\mathbf{Y})$ and $\text{Var}(\mathbf{AY}) = \mathbf{A}\text{Var}(\mathbf{Y})\mathbf{A}^\top$.

$$E(\hat{\beta}) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \underbrace{E(\mathbf{Y})}_{\mathbf{X}\beta} = \underbrace{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}}_{\mathbf{I}} \beta = \beta$$

$$\begin{aligned} \text{Var}(\hat{\beta}) &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \text{Var}(\mathbf{Y}) ((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top)^\top \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \cdot \sigma^2 \mathbf{I} \cdot \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} \\ &= \begin{bmatrix} V(\hat{\beta}_0) & \cdots & C(\hat{\beta}_0, \hat{\beta}_p) \\ \vdots & \ddots & \vdots \\ C(\hat{\beta}_p, \hat{\beta}_0) & \cdots & V(\hat{\beta}_p) \end{bmatrix} \end{aligned}$$

When $p = 1$:

$$\begin{aligned} V(\hat{\beta}_0) &= \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right), \quad V(\hat{\beta}_1) = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \\ C(\hat{\beta}_0, \hat{\beta}_1) &= C(\hat{\beta}_1, \hat{\beta}_0) = -\bar{x} V(\hat{\beta}_1) \end{aligned}$$

For a specific set of x -values, $\mathbf{x}_0 = [1 \ x_{01} \ \dots \ x_{0p}]$, we have

- ▶ on average (the fitted line) $\hat{Y}_0 = \mathbf{x}_0 \hat{\boldsymbol{\beta}}$:

$$E(\hat{Y}_0) = \mathbf{x}_0 E(\hat{\boldsymbol{\beta}}) = \mathbf{x}_0 \boldsymbol{\beta} = \beta_0 + \beta_1 x_{01} + \dots + \beta_p x_{0p} = \mu_0,$$

$$V(\hat{Y}_0) = \mathbf{x}_0 \text{Var}(\hat{\boldsymbol{\beta}}) \mathbf{x}_0^T = \sigma^2 \mathbf{x}_0 (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0^T$$

$$[p = 1] = \sigma^2 \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)$$

- ▶ for a new observation $\hat{Y}_{\text{pred}_0} = \mathbf{x}_0 \hat{\boldsymbol{\beta}} + \epsilon_0$:

$$E(\hat{Y}_{\text{pred}_0}) = \mathbf{x}_0 E(\hat{\boldsymbol{\beta}}) + E(\epsilon_0) = \mathbf{x}_0 \boldsymbol{\beta} = \mu_0,$$

$$V(\hat{Y}_{\text{pred}_0}) = \mathbf{x}_0 \text{Var}(\hat{\boldsymbol{\beta}}) \mathbf{x}_0^T + V(\epsilon_0) = \sigma^2 (1 + \mathbf{x}_0 (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0^T)$$

$$[p = 1] = \sigma^2 \left(1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)$$

Standard error

The **standard error** (s.e. = medelfel), $d(\hat{\theta}) = \sqrt{\hat{V}(\hat{\theta})}$, of an estimate $\hat{\theta}$, replaces any unknown parameters in $V(\hat{\theta})$ by their estimates. Here, replace σ by $\hat{\sigma} = s$.

Comments

- ▶ $\hat{\beta}$ exists only when $\mathbf{X}^T \mathbf{X}$ is non-singular. Near-singularity makes all β -estimates very uncertain.
- ▶ Never calculate $(\mathbf{X}^T \mathbf{X})^{-1}$ “by hand”. Take a course in Numerical analysis (or rely on R).
- ▶ The uncertainty of the β -estimates is mostly due to the sample size, n , and the structure of the x -variables, as expressed in $(\mathbf{X}^T \mathbf{X})^{-1}$.
- ▶ The uncertainty of the predictions \hat{Y}_0 is **also** due to how far \mathbf{x}_0 is from the midpoint of the x -variables.
- ▶ $V(\hat{Y}_{\text{pred}_0}) \rightarrow \sigma^2 > 0$ when $V(\hat{Y}_0) \rightarrow 0$.

Ice cream: estimates

Y = weight loss (g), x = storage time (weeks).

Model $Y = \beta_0 + \beta_1 x + \epsilon$

Variable	parameter	estimate	s.e.	unit
intercept (time = 0)	β_0	-5.7	0.81	g
storage time	β_1	1.33	0.03	g/week
resid.std.dev	σ	0.80		g

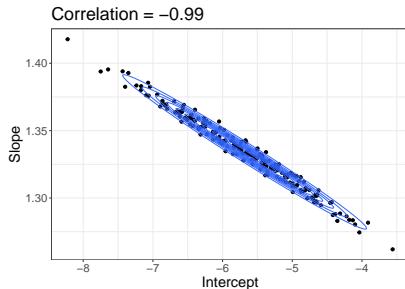
Fitted line: $\hat{Y} = -5.7 + 1.33x$.

The estimates $\hat{\beta}_0$ and $\hat{\beta}_1$ are correlated

$$\hat{E}(\hat{\beta}) = \hat{\beta} = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{bmatrix} = \begin{bmatrix} -5.7 \\ 0.03 \end{bmatrix}$$

$$\text{Var}(\hat{\beta}) = \hat{\sigma}^2(\mathbf{X}^T \mathbf{X})^{-1} = \begin{bmatrix} 0.65 & -0.02 \\ -0.02 & 0.0007 \end{bmatrix}$$

$$\begin{aligned} \hat{\rho}(\hat{\beta}_0, \hat{\beta}_1) &= \frac{-0.02}{\sqrt{0.65 \cdot 0.0007}} = \\ &= \frac{-0.02}{0.81 \cdot 0.03} \\ &\approx -0.99 \end{aligned}$$



Predictions

If we store the ice cream for $x_0 = 34$ weeks, how much weight loss can we expect on average? in a single package?

	estimate	s.e.	unit
on average	$\hat{Y}_0 = -5.7 + 1.33 \cdot 34 = 39.7$	0.15	g
single package	$\hat{Y}_{\text{pred}_0} = 39.7 + \epsilon_0$	0.82	g

Note: $0.82 = \sqrt{0.80^2 + 0.15^2}$.