

**Inf2 - Foundations of Data Science:
Regression and inference -
From the maximum likelihood principle
to linear regression**



THE UNIVERSITY *of* EDINBURGH
informatics

F O U N D A T I O N S
O F
D A T A
S C I E N C E

Announcements

- Badges! (and Stickers) Please let me know if I owe you one



We want to investigate the relationship between the number of bikes hired in an hour and the mean temperature during that hour

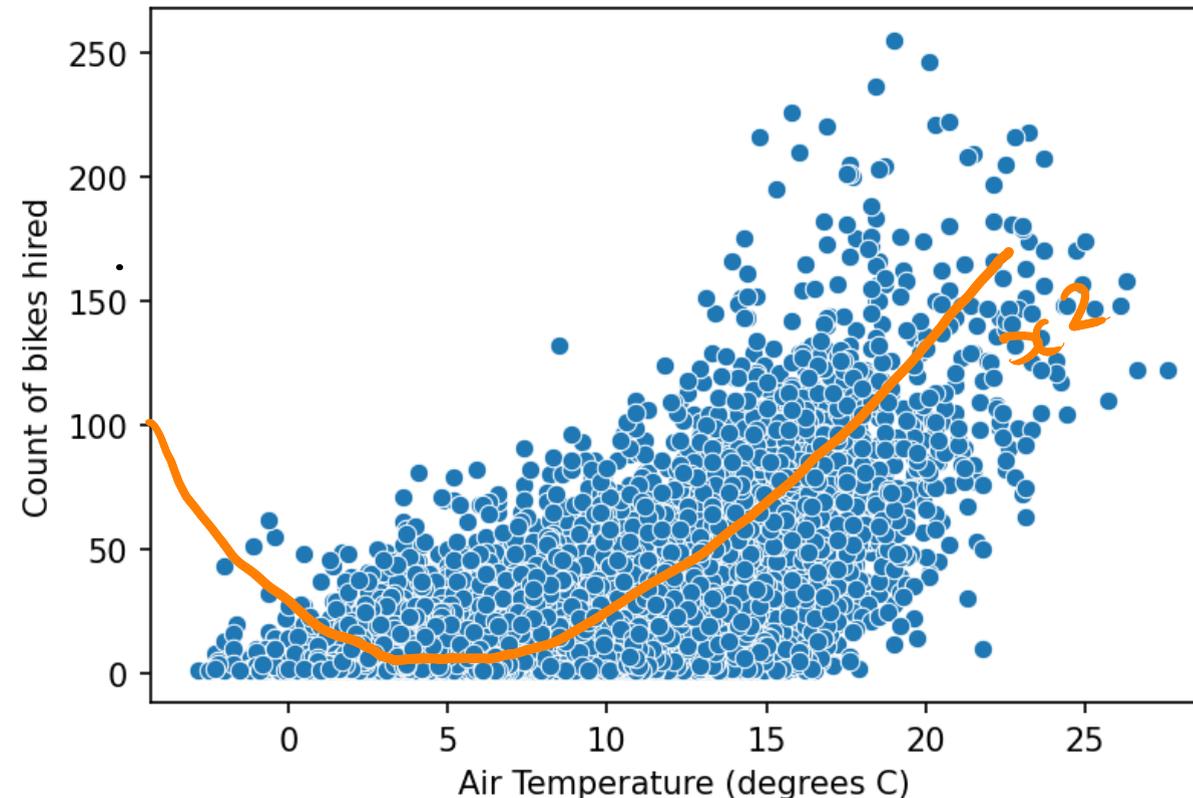
Is there a problem with using ordinary least squares linear regression to do this?



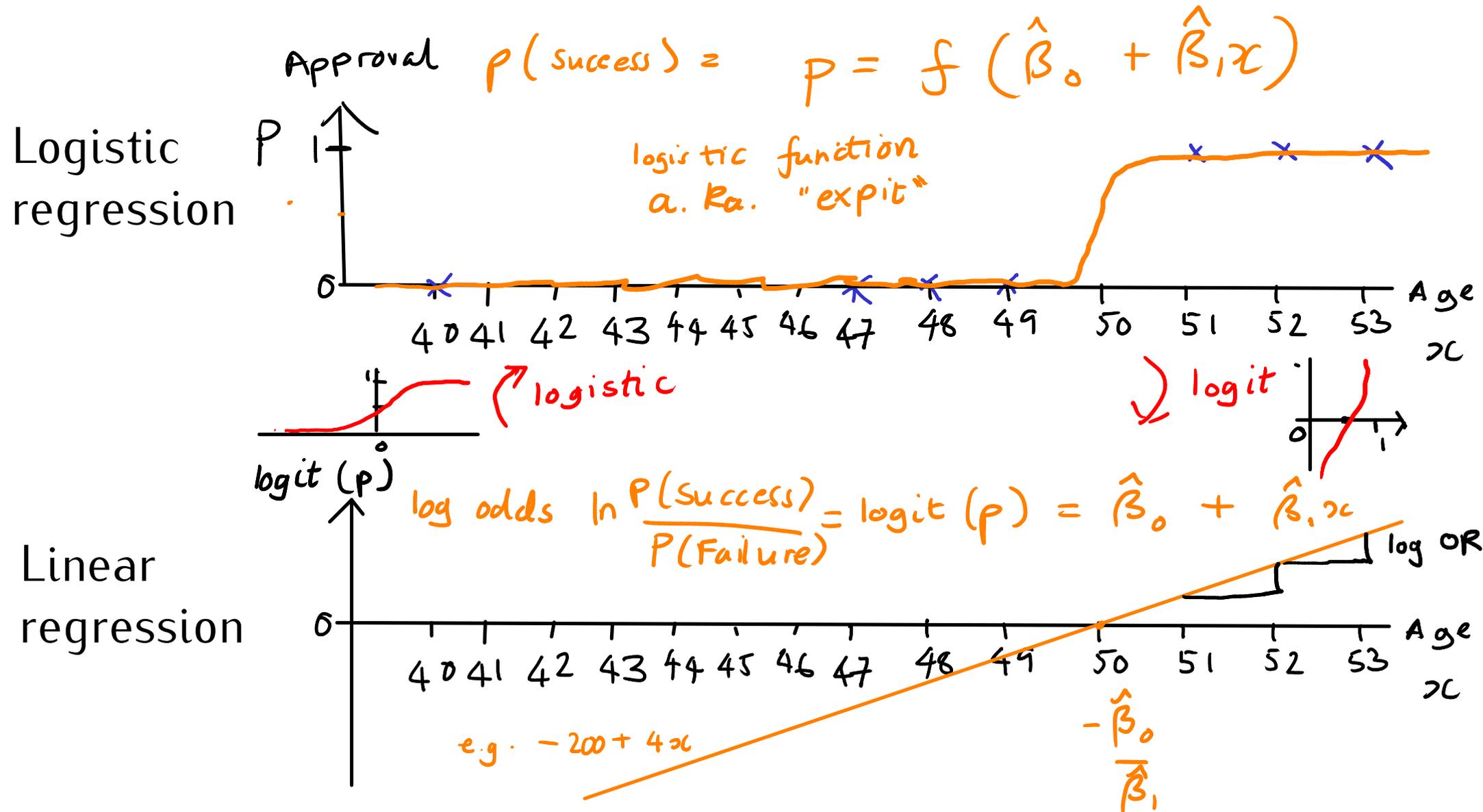
Image copyright Pashley Cycles

Data sources:

- Edinburgh Just Eat Bikes data 2020, on on OpenData Scotland
- Edinburgh temperature observations, Met Office via MIDAS



Where we're at in the Maximum Likelihood Principle and Regression



Overview

Today

1. The maximum likelihood principle
2. Application of maximum likelihood to a simple example
3. Application of maximum likelihood to linear regression

Wednesday

- Max likelihood with non-normal distributions
- Generalised linear regression, including logistic regression

**Inf2 - Foundations of Data Science:
Regression and inference -
The maximum likelihood principle**



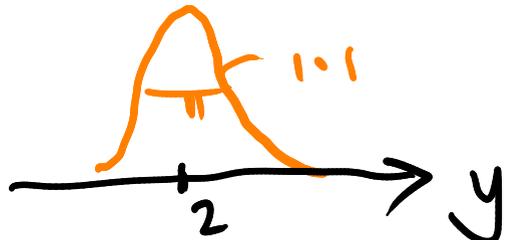
THE UNIVERSITY *of* EDINBURGH
informatics

F O U N D A T I O N S
O F
D A T A
S C I E N C E

Intuition for maximum likelihood principle

Statistical model, e.g. normal dist.

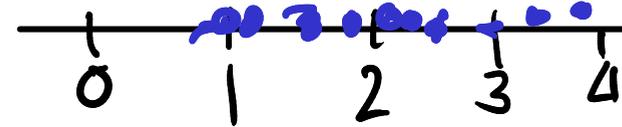
Data



A hand-drawn orange normal distribution curve on a horizontal axis labeled 'y'. The mean is marked as 2. The standard deviation is indicated as 1.1. The curve is centered at 2 and has a peak at that point.

$$p(y | 2, 1.1^2) = \frac{1}{\sqrt{2\pi} \cdot 1.1} e^{-\frac{1}{2} \left(\frac{y-2}{1.1}\right)^2}$$

sample



Alternatively

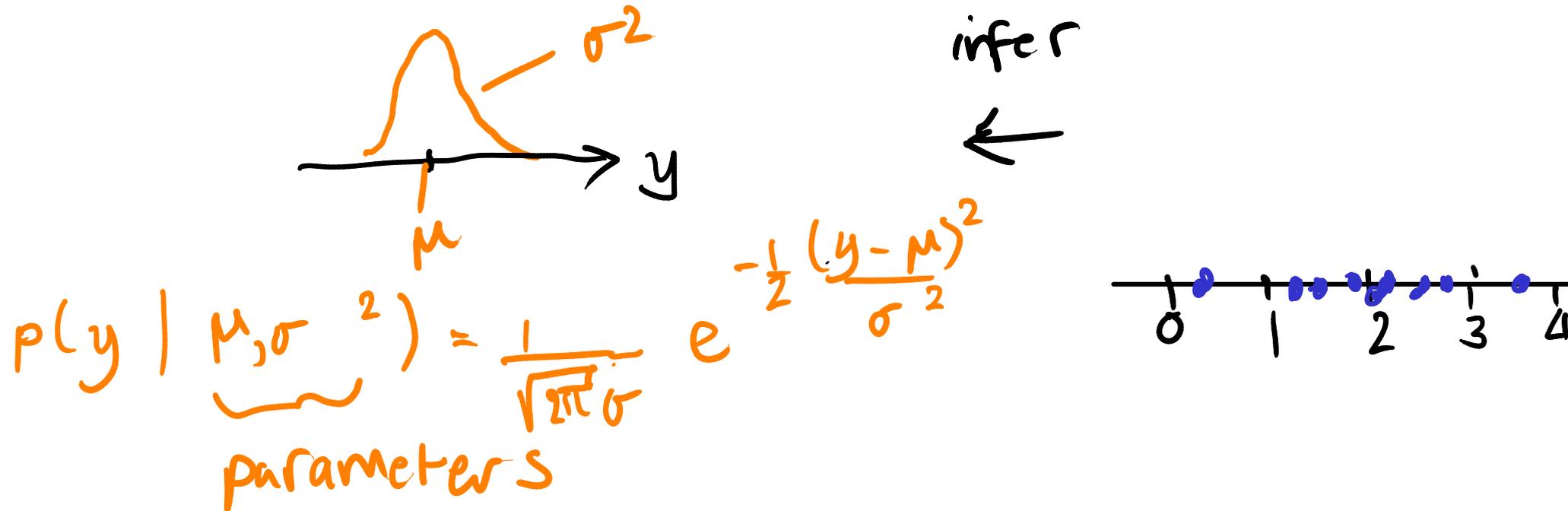
$$y \sim \mathcal{N} \left(\underset{\uparrow \mu}{2}, \underset{\uparrow \sigma^2}{1.1^2} \right)$$

y is drawn from a normal dist.
with $\mu = 2$ and $\sigma^2 = 1.1^2$

Intuition for maximum likelihood principle

Statistical model, e.g. normal dist.

Data



Given a normal distribution, what parameters are most likely to have generated data?

Exercise

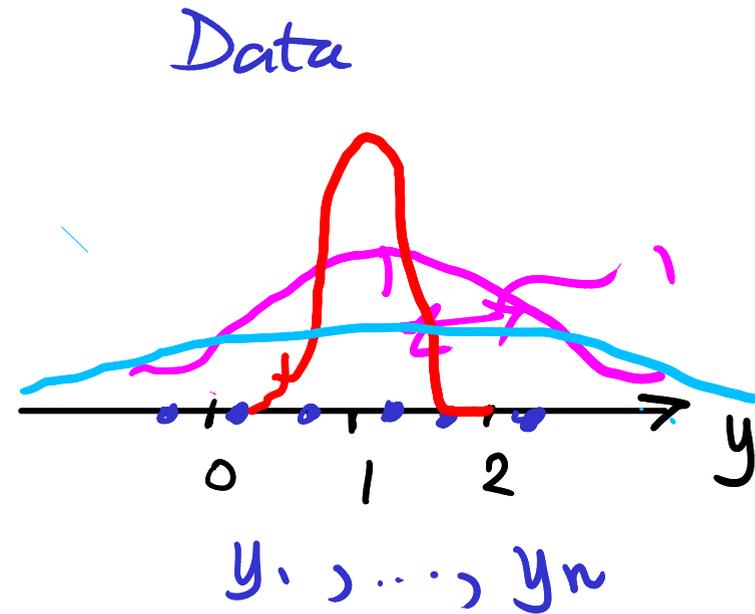
Which model is most likely to have generated this data?

① $y_i \sim \mathcal{N}(\overset{\mu}{0}, \overset{\sigma^2}{1})$

② $y_i \sim \mathcal{N}(1, 0.1)$

③ $y_i \sim \mathcal{N}(1, 5)$

④ $y_i \sim \mathcal{N}(1, 1)$



Definition of Principle of Maximum Likelihood

For a set of observed data and a given statistical model the principle of maximum likelihood states that the parameters of the model are adjusted so as to maximise the likelihood that the model generated the observed data.

$$\text{Data : } \underline{Y} = \{y_1, \dots, y_n\}$$

$$\text{Likelihood } P(\underline{Y} = y_1, \dots, y_n | \vartheta_1, \dots, \vartheta_m)$$

$$\text{Normal : } \vartheta_1 = \mu$$
$$\vartheta_2 = \sigma^2$$

**Inf2 - Foundations of Data Science:
Regression and inference -
Application of the maximum likelihood
principle to a simple example**



THE UNIVERSITY *of* EDINBURGH
informatics

F O U N D A T I O N S
O F
D A T A
S C I E N C E

Application to 1-variable example

1. Assume samples are drawn independently
2. Assume each sample is drawn from a normal distribution

$$P(Y = y_i | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2} \left(\frac{y_i - \mu}{\sigma}\right)^2} \quad (2)$$

Ass. ① :

$$P(\underline{Y} = y_1, \dots, y_n | \mu, \sigma^2) = P(Y = y_1 | \mu, \sigma^2) \\ \times P(Y = y_2 | \mu, \sigma^2) \\ \times \dots \\ \times P(Y = y_n | \mu, \sigma^2)$$

More compact notation...

Likelihood

$$P(\underline{y} = y_1, \dots, y_n) = P(Y = y_1 | \mu, \sigma^2) \\ \times P(Y = y_2 | \mu, \sigma^2) \\ \times \dots \\ \times P(Y = y_n | \mu, \sigma^2)$$

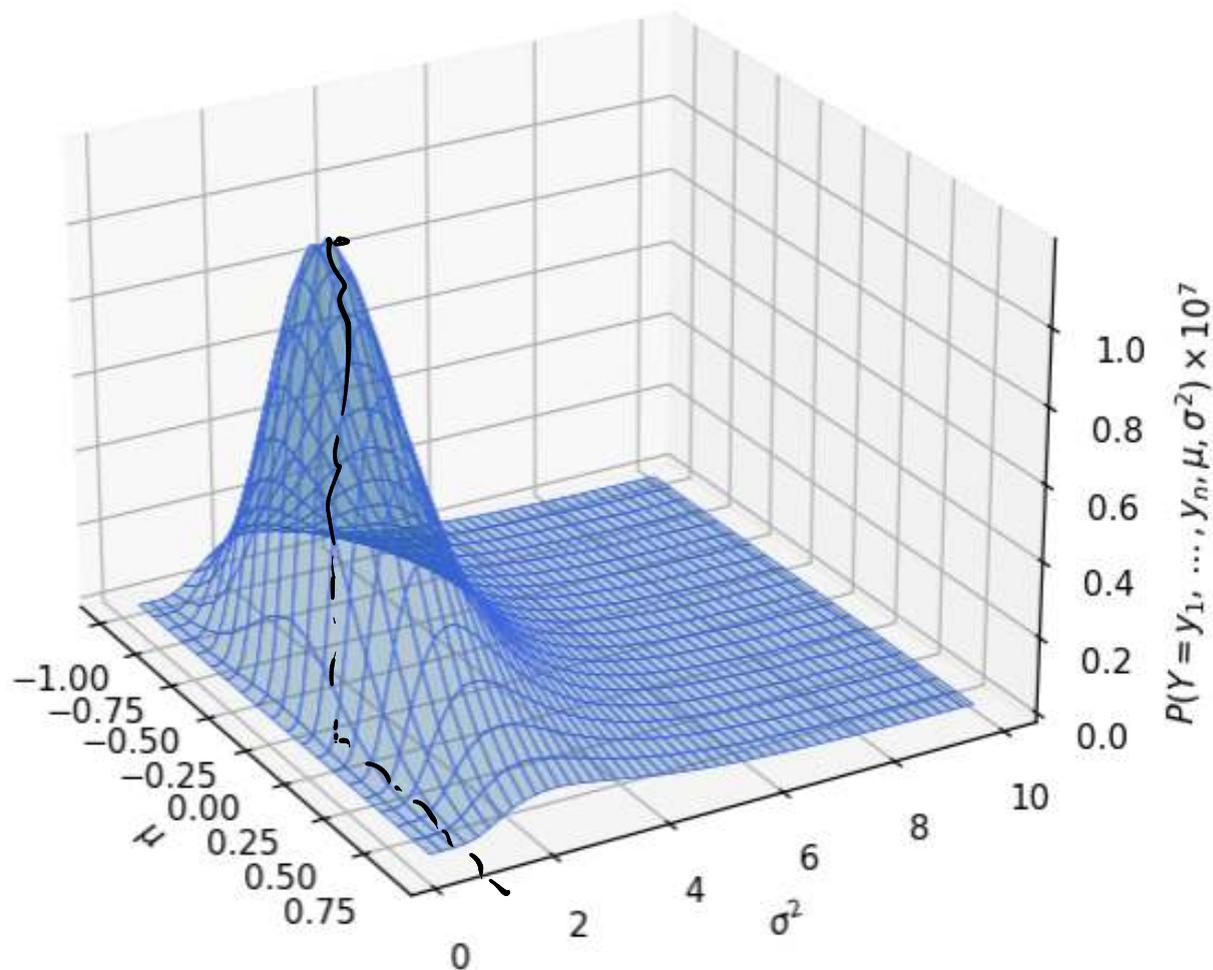
$$= \prod_{i=1}^n P(Y = y_i | \mu, \sigma^2) \\ = \prod_{i=1}^n \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{1}{2} \left(\frac{y_i - \mu}{\sigma} \right)^2}$$

Likelihood as a function of parameters

Data:

y_1, \dots, y_{10} drawn from
 $\mathcal{N}(0, 1)$

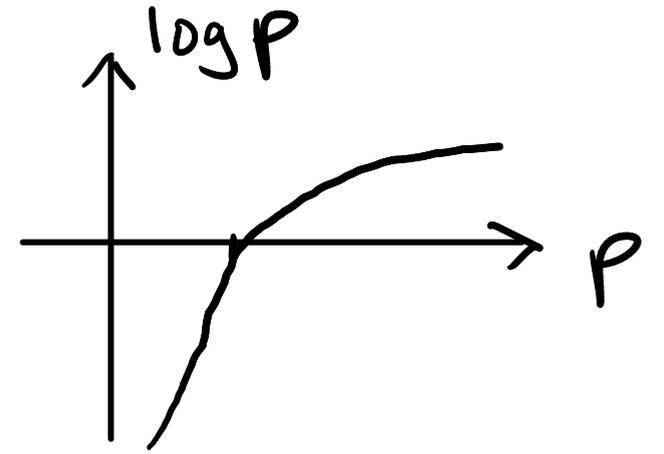
| | Data |
|----------|-----------|
| y_1 | 1.624345 |
| y_2 | -0.611756 |
| y_3 | -0.528172 |
| y_4 | -1.072969 |
| y_5 | 0.865408 |
| y_6 | -2.301539 |
| y_7 | 1.744812 |
| y_8 | -0.761207 |
| y_9 | 0.319039 |
| y_{10} | -0.249370 |



[code]

Log-likelihood equations: products to sums

Remember: $\ln a b = \ln a + \ln b$



$$\ln(p_1 \times \dots \times p_n) = \ln p_1 + \dots + \ln p_n$$

$$\ln \left(\prod_{i=1}^n p_i \right) = \sum_{i=1}^n \ln p_i$$

Log likelihood

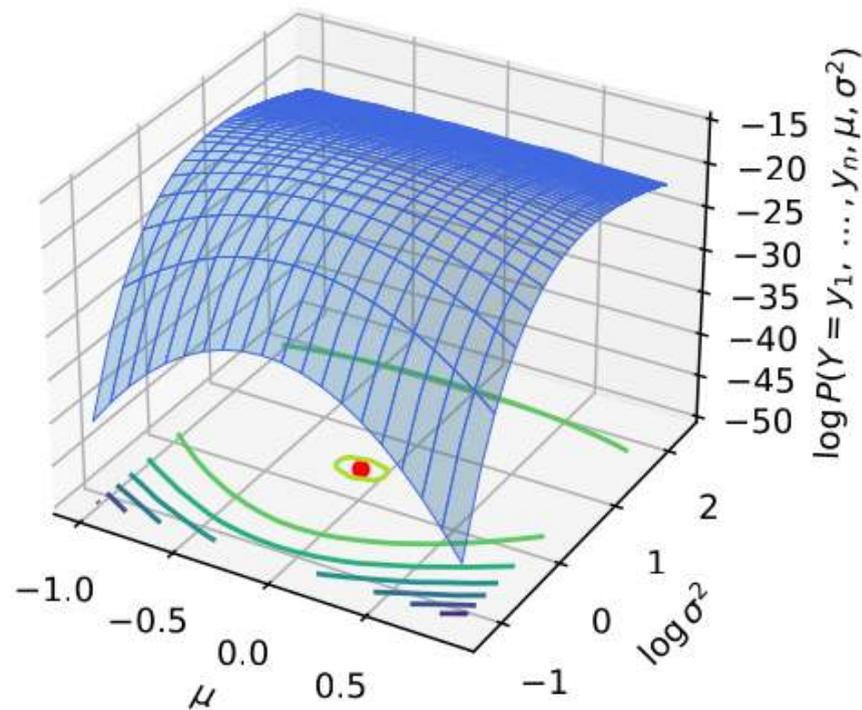
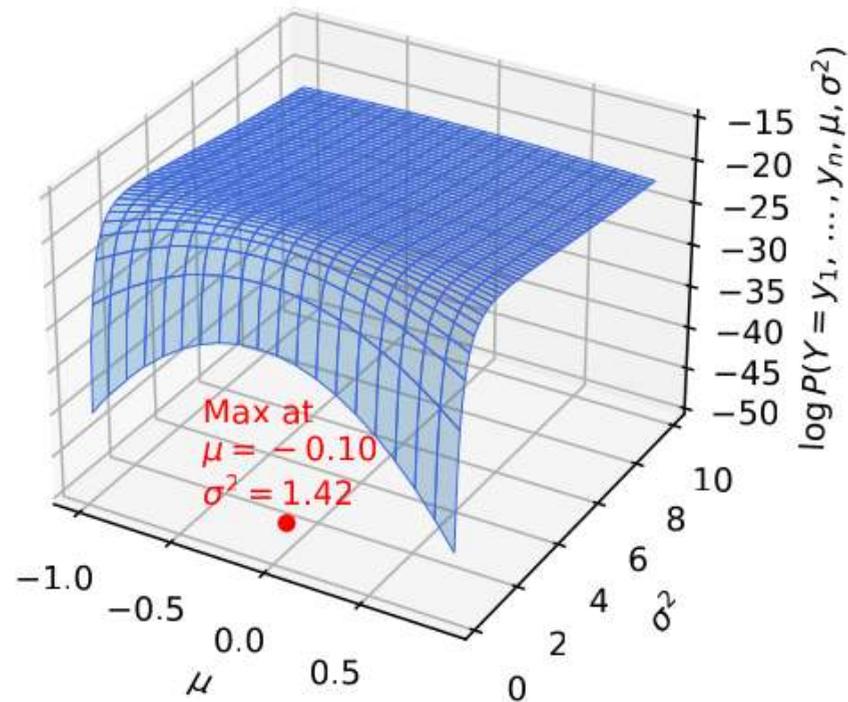
$$\ln \left(\prod_{i=1}^n P(y = y_i | \mu, \sigma^2) \right) = \sum_{i=1}^n \ln P(y = y_i | \mu, \sigma^2)$$

The log of the normal distribution

$$\begin{aligned}\ln P(X=y_i | \mu, \sigma^2) &= \ln \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{y_i-\mu}{\sigma}\right)^2} = \ln \left(\frac{1}{\sqrt{2\pi}\sigma} \right) + \ln \left(e^{-\frac{1}{2}\left(\frac{y_i-\mu}{\sigma}\right)^2} \right) \\ &= -\frac{1}{2} \ln 2\pi\sigma^2 - \frac{1}{2}\left(\frac{y_i-\mu}{\sigma}\right)^2\end{aligned}$$

Log-likelihood of data drawn from a normal distribution as a function of parameters

$$\ln P(\underline{Y} = y_1, \dots, y_n \mid \mu, \sigma^2)$$
$$= \sum_{i=1}^n \left(-\frac{1}{2} \ln 2\pi\sigma^2 - \frac{1}{2} \frac{(y_i - \mu)^2}{\sigma^2} \right)$$



Maximum likelihood estimates of parameters

$$\begin{aligned} \ln P(\underline{Y} = y_1, \dots, y_n \mid \mu, \sigma^2) \\ = \sum_{i=1}^n \left(-\frac{1}{2} \ln 2\pi\sigma^2 - \frac{1}{2} \frac{(y_i - \mu)^2}{\sigma^2} \right) \end{aligned}$$

Maximise w.r.t μ and σ^2 to give max likelihood estimates (MLE)

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n y_i$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{\mu})^2$$

Exercise: prove these statements by differentiating w.r.t. μ and σ^2

The beauty of logs and sums

- Sum of logs is easy to represent within limits of floating point arithmetic
- Log likelihood function is smoother than likelihood function
- Sums are easy to differentiate; products are not

.

.

**Inf2 - Foundations of Data Science:
Regression and inference -
Application of the maximum likelihood
principle to linear regression**



THE UNIVERSITY *of* EDINBURGH
informatics

F O U N D A T I O N S
O F
D A T A
S C I E N C E

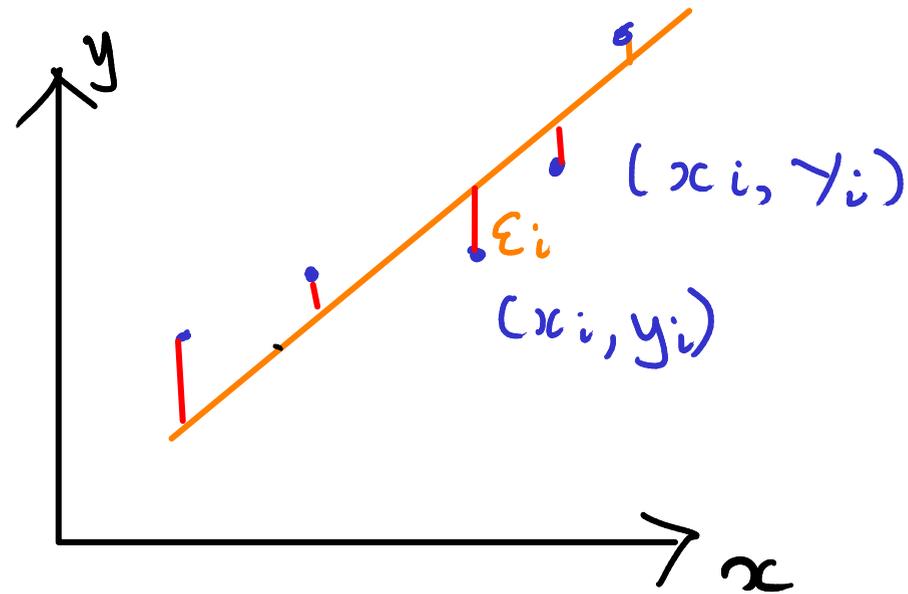
Application of max likelihood to linear regression

$$\hat{y}_i = \beta_0 + \beta_1 x_i$$

$$y_i = \beta_0 + \beta_1 x_i + \underbrace{\varepsilon_i}_{\text{error term}}$$

$$\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$$

$$y_i \sim \mathcal{N}(\beta_0 + \beta_1 x_i, \sigma^2)$$



$$\begin{aligned} \ln P(Y = y_1, \dots, y_n; x_1, \dots, x_n | \beta_0, \beta_1, \sigma^2) \\ = \prod_{i=1}^n \left(-\frac{1}{2} \ln 2\pi\sigma^2 - \frac{1}{2} \left(\frac{y_i - \beta_0 - \beta_1 x_i}{\sigma} \right)^2 \right) \end{aligned}$$

Relationship to ordinary least squares

$$\begin{aligned} \ln p(y_1, \dots, y_n; x_1, \dots, x_n | \beta_0, \beta_1, \sigma^2) \\ &= \sum_{i=1}^n \left(-\frac{1}{2} \ln 2\pi\sigma^2 - \frac{1}{2} \frac{(y_i - \beta_0 - \beta_1 x_i)^2}{\sigma^2} \right) \\ &= -\frac{n}{2} \ln 2\pi\sigma^2 - \underbrace{\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2}_{\text{SSE}} \end{aligned}$$

Estimate of coefficients

- Analytical solutions for $\hat{\beta}_0$, $\hat{\beta}_1$ and $\hat{\sigma}^2$ that maximise like lihood
- $\hat{\beta}_0$ and $\hat{\beta}_1$: as per ordinary least squares
- Variance of residuals:

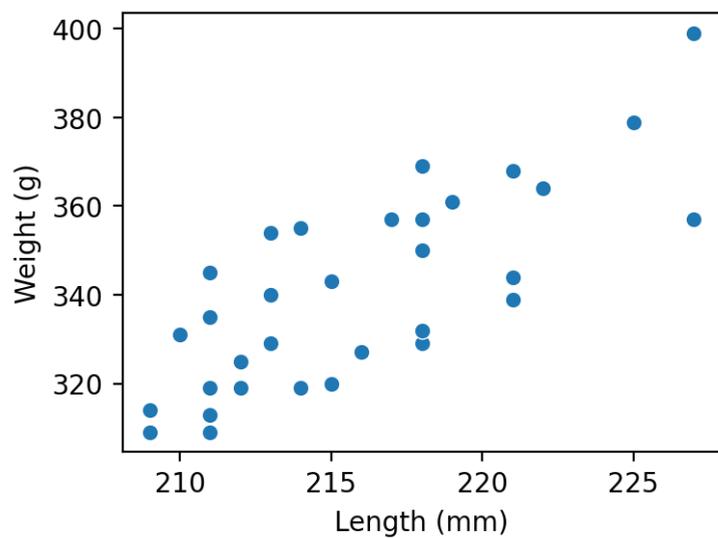
$$\begin{aligned}\hat{\sigma}^2 &= \frac{1}{n} \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 \\ &= \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \frac{SSE}{n} \leftarrow \text{Biased}\end{aligned}$$

Sampling theory $\hat{\sigma}^2 = \frac{SSE}{n-2} \leftarrow \text{Unbiased}$

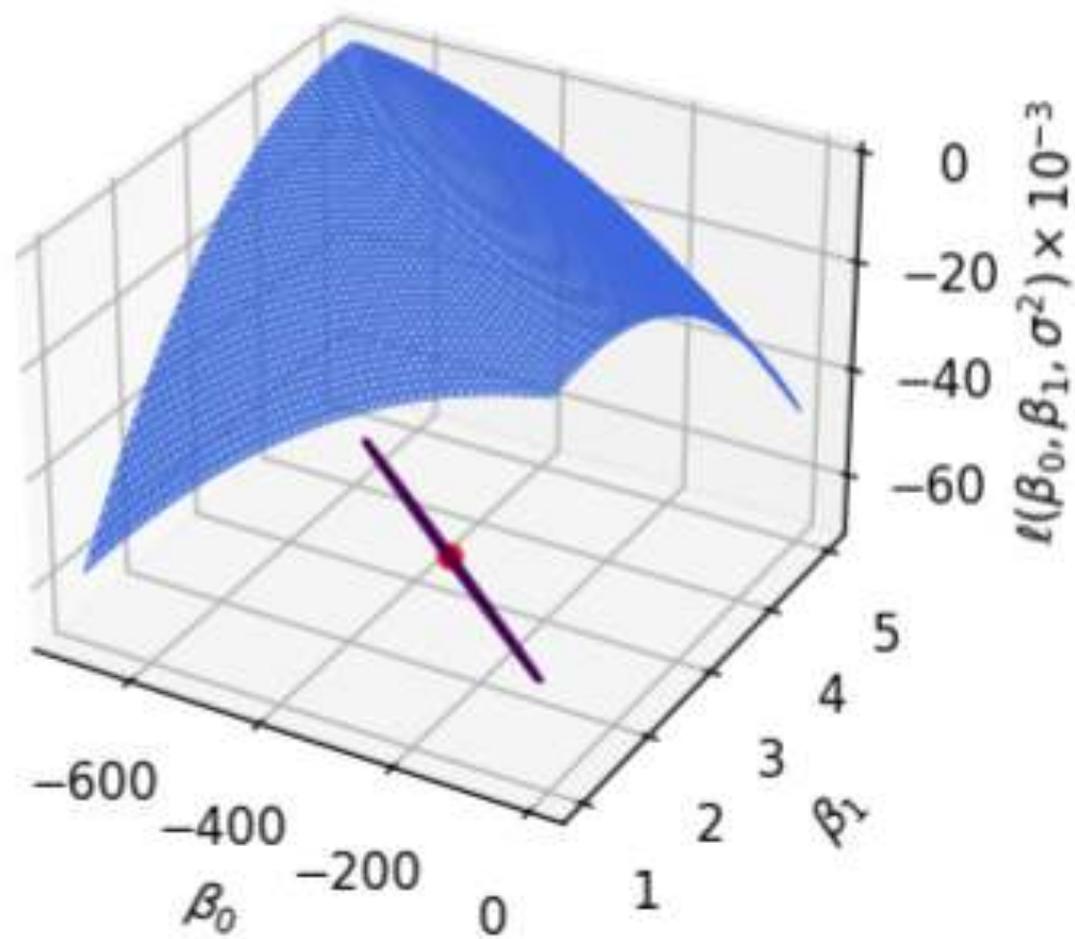
Log likelihood of coefficients



Peter Trimming, Wikimedia Commons, CC BY 2.0

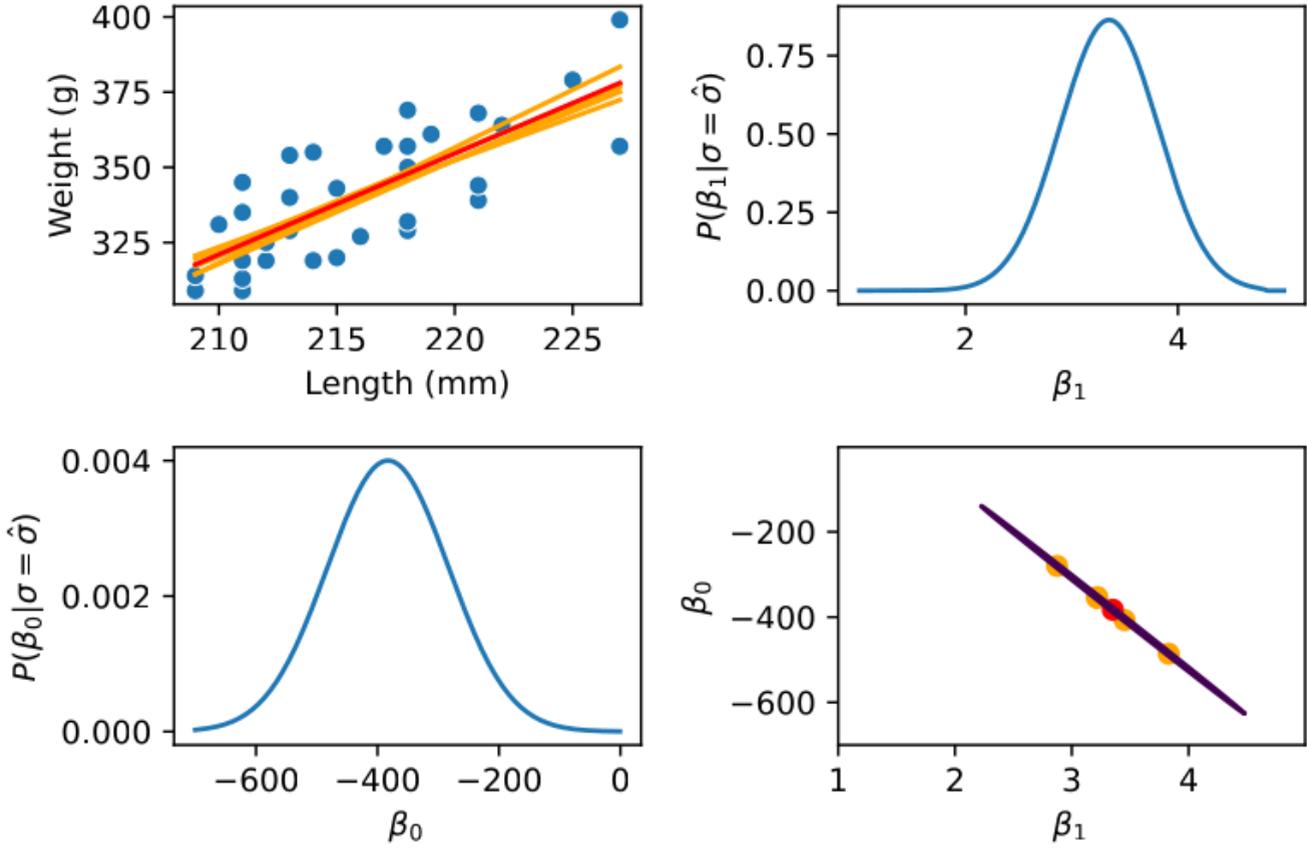


Data from Wauters and Dhondt 1989

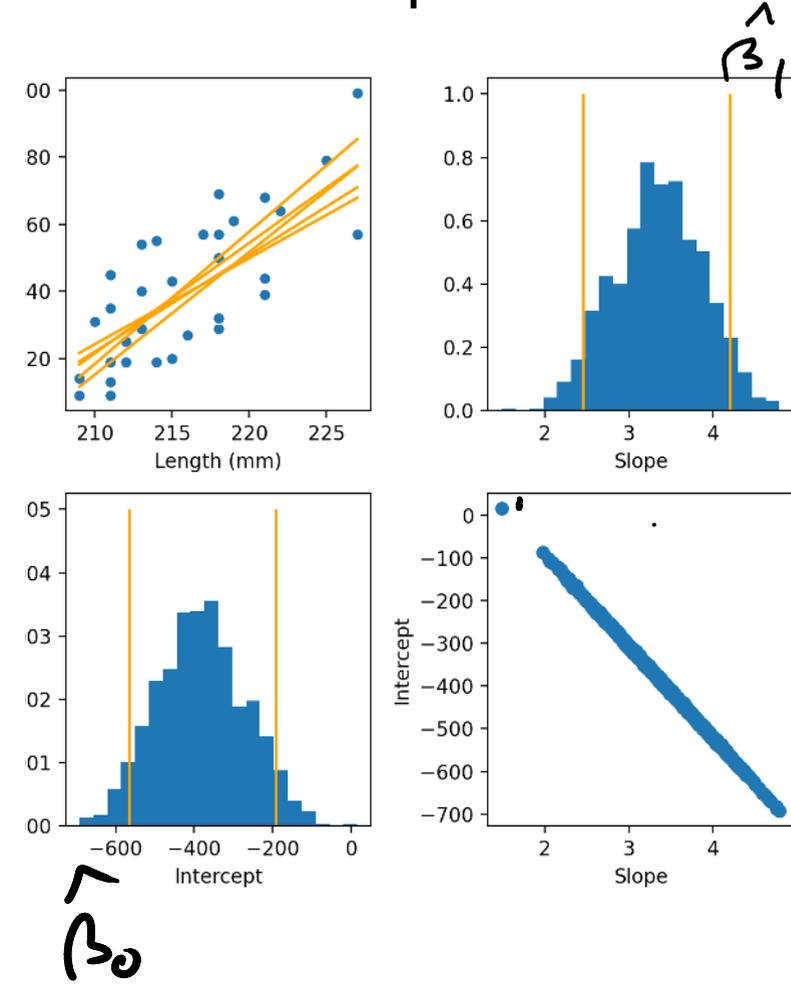


Parameter uncertainty

Max likelihood



cf Bootstrap



Overview

1. Maximum likelihood principle
 - What model was most likely to have generated the data
2. Maximum likelihood principle applied to simple example
 - Log likelihood turns out to be useful
 - Gives rise to familiar estimates for mean and variance
3. Maximum likelihood principle applied to linear regression
 - Turns out to give ordinary least squares
 - Link with coefficient uncertainty and the bootstrap estimates of parameter uncertainty

We want to investigate the relationship between the number of bikes hired in an hour and the mean temperature during that hour

Is there a problem with using ordinary least squares linear regression to do this?

Data sources:

- Edinburgh Just Eat Bikes data 2020
- Edinburgh temperature observations, Met Office via MIDAS

