

# Linear Regression Models

We consider the following simple linear model for the data

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i \quad \text{with } \epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2), i = 1, \dots, n$$

or equivalently

$$y_i | x_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$$

using vector notation, we have

$$\mathbf{y} = \mathbf{x}\beta + \epsilon, \quad \mathbf{y} | \mathbf{x}\beta \sim N(\mathbf{x}\beta, \sigma^2)$$

Note: The parameters are  $\beta_0, \beta_1$  and the variance  $\sigma^2$  assumed to be constant accross all observations. They are assumed to be fixed and unknown. The  $x_i$  are assumed to be fixed and the  $y_i$  are assumed to be random, because of the random error terms  $e_i$ .

# Classical approach to Bayesian Linear Regression Models

As usual in Bayesian statistics, the posterior distribution is proportional to the likelihood times the prior distributions of the parameters, so that we have  $\pi(\beta, \sigma^2 \mid \mathbf{y}) \propto L(\beta, \sigma^2) \pi(\beta, \sigma^2)$ .

Many choices for priors are available. For example, if the following noninformative reference priors are chosen  $\pi(\beta \mid \sigma^2) \propto 1$  and  $\pi(\sigma^2) \propto \frac{1}{\sigma^2}$ , then, the posterior, up to proportionality, becomes

$$\begin{aligned}\pi(\beta, \sigma^2 \mid \mathbf{y}) &\propto \prod_{i=1}^n L(\beta, \sigma^2) \pi(\beta, \sigma^2) \\ &\propto (\sigma^2)^{-n/2} \exp \left[ -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \right] \left( \frac{1}{\sigma^2} \right) \\ &\propto (\sigma^2)^{-(n+2)/2} \exp \left[ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) \right]\end{aligned}$$

# Inference for the Bayesian Linear Regression Models

As usual in Bayesian statistics, the posterior distribution is proportional to the likelihood times the prior distributions of the parameters, so that we have  $\pi(\beta, \sigma^2 \mid \mathbf{y}) \propto L(\beta, \sigma^2) \pi(\beta, \sigma^2)$ .

Many choices for priors are available. For example, if the following noninformative reference priors are chosen  $\pi(\beta \mid \sigma^2) \propto 1$  and  $\pi(\sigma^2) \propto \frac{1}{\sigma^2}$ , then, the posterior, up to proportionality, becomes

$$\begin{aligned}\pi(\beta, \sigma^2 \mid \mathbf{y}) &\propto \prod_{i=1}^n L(\beta, \sigma^2) \pi(\beta, \sigma^2) \\ &\propto (\sigma^2)^{-n/2} \exp \left[ -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \right] \left( \frac{1}{\sigma^2} \right) \\ &\propto (\sigma^2)^{-(n+2)/2} \exp \left[ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) \right]\end{aligned}$$

# Cars dataset in R

We work with the dataset 'cars' to illustrate linear and log-linear models, both in the traditional frequentist and bayesian framework. Two variables are considered: 'speed', which corresponds to speed (mph) and 'dist' which corresponds to the stopping distance (ft). The dataset, which consists in 50 observations for each variable is then split into a train set (70%) and a test set. Below the R code to access the initial data set, partitioning it and creating linear models.

```
1 library(tidyverse)
2 data(cars)
3
4 # 0.1 Split the dataset into train and test sets
5 set.seed(2024)
6 sample = sample(c(TRUE, FALSE), nrow(cars), replace=T, prob=c(0.7, 0.3))
7 #Split your data into training (70%) and test (30%) sets
8 train = cars[sample, ]
9 test = cars[!sample, ]
10
11 # 1.1 Create frequentist linear model, log linear model and plot
12 freq.lin.mod = lm(dist~speed, data=train)
13 freq.log.lin.mod = lm(log(dist)~speed, data = train)
14
15 summary(freq.lin.mod)
16 summary(freq.log.lin.mod)
```

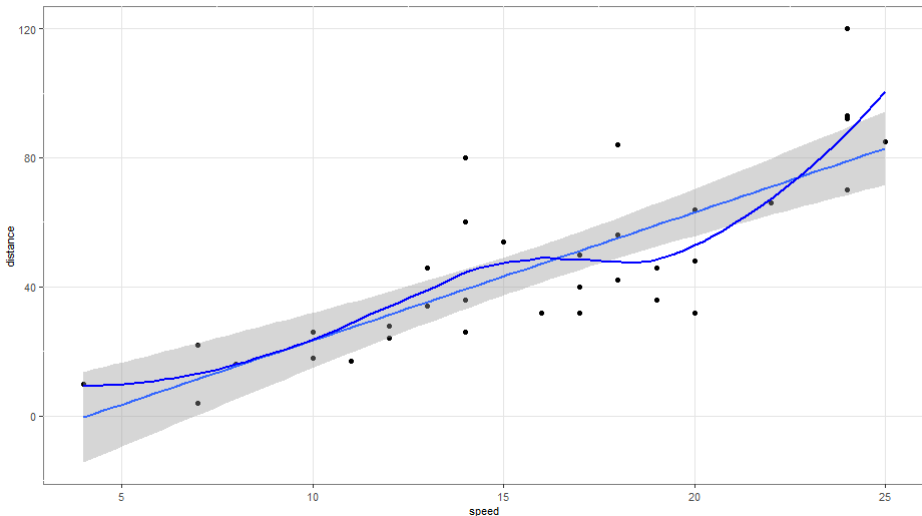
# Frequentist LM outputs

<i>Dependent variable:</i>	
	dist
speed	3.968*** (0.527)
Constant	-16.293* (8.843)
Observations	34
R <sup>2</sup>	0.639
Adjusted R <sup>2</sup>	0.628
Residual Std. Error	16.569 (df = 32)
F Statistic	56.652*** (df = 1; 32)
<i>Note:</i> *p<0.1; **p<0.05; ***p<0.01	

# Frequentist LM, observations

Plot of observations and the fitted Linear model

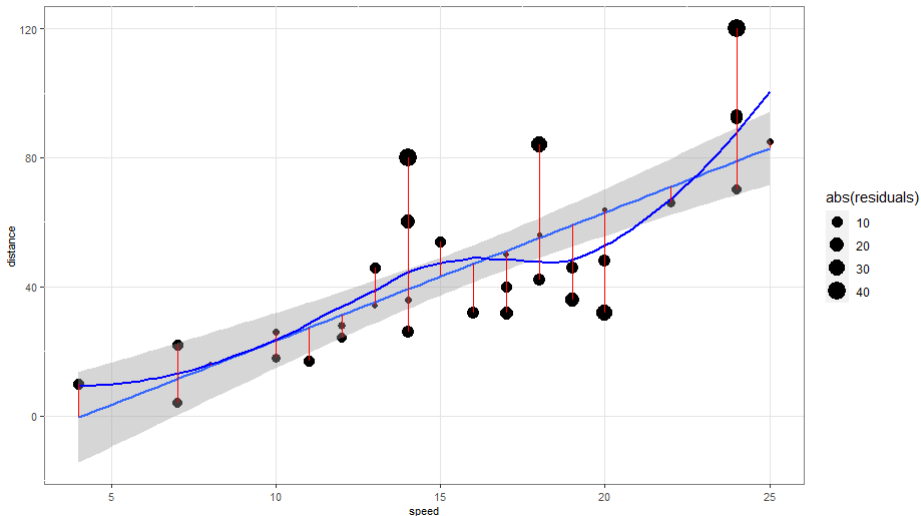
*Cars train data set*



# Frequentist LM, residuals

Plot of the residuals and their distance to the fitted Linear model

*Cars train data set*



# Frequentist LM outputs

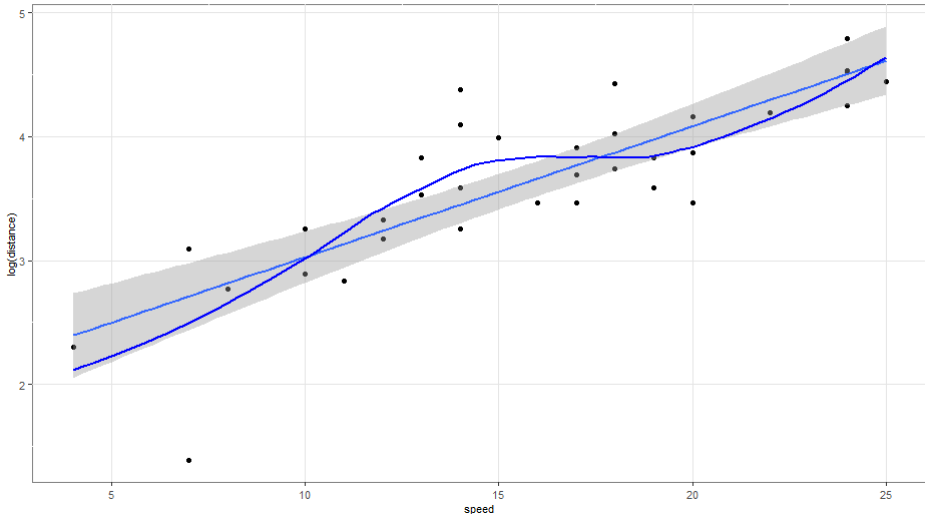
	<i>Dependent variable:</i>
	dist
speed	3.968*** (0.527)
Constant	-16.293* (8.843)
Observations	34
R <sup>2</sup>	0.639
Adjusted R <sup>2</sup>	0.628
Residual Std. Error	16.569 (df = 32)
F Statistic	56.652*** (df = 1; 32)
Note:	*p<0.1; **p<0.05; ***p<0.01



# Frequentist LLM, observations

Plot of observations and the fitted Log-Linear model

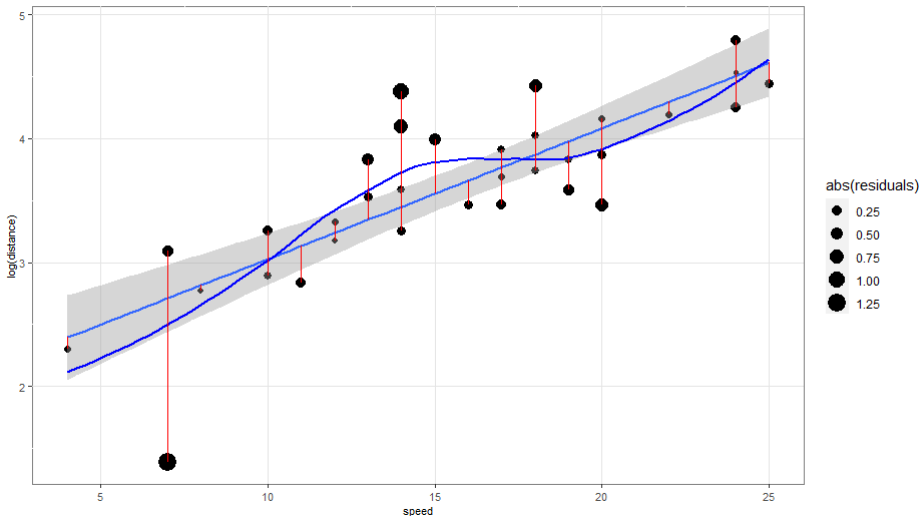
*Cars train data set*



# Frequentist LLM, residuals

Plot of the residuals and their distance to the fitted Log-Linear model

*Cars train data set*

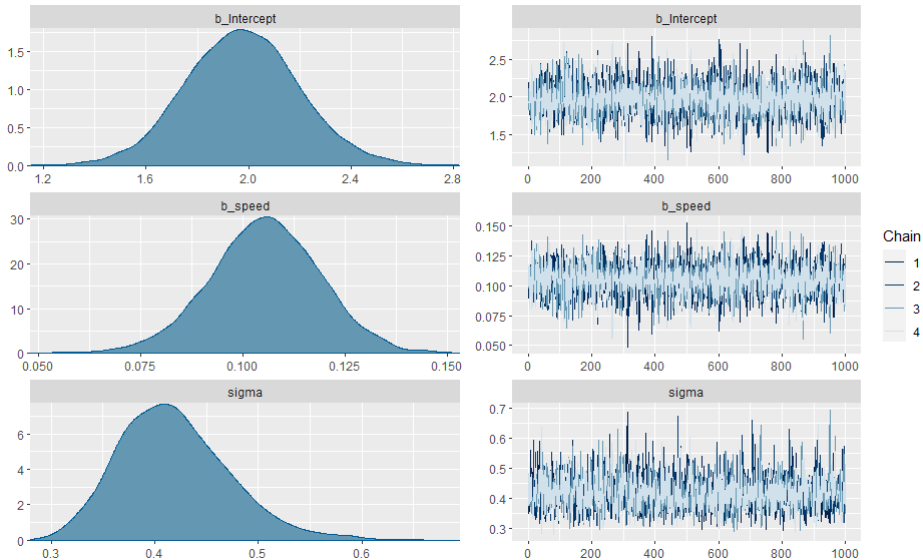


# Bayesian LM outputs

<i>Dependent variable:</i>	
dist	
speed	3.96 (0.54)
Constant	-16.27 (9.15)
Observations	34

# Bayesian LLM outputs

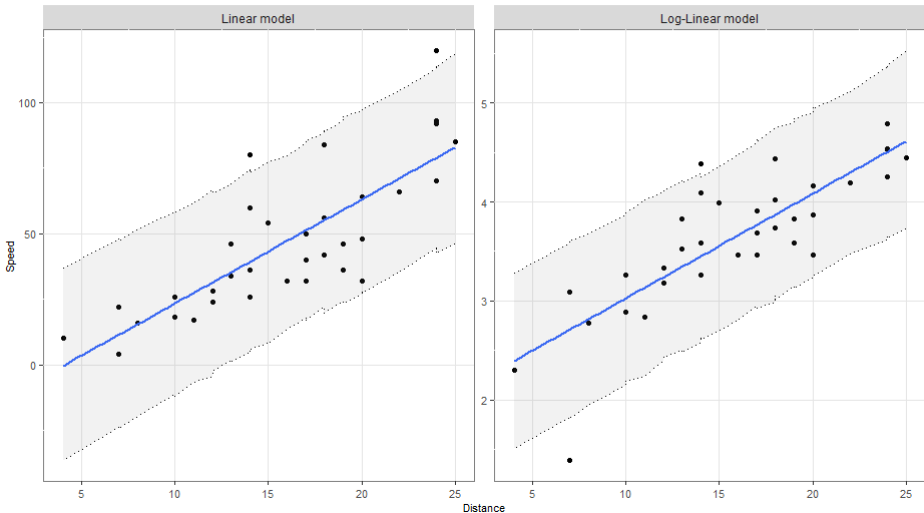
# Bayesian LLM, posterior densities



# Bayesian LLM, fitted models

## Scatterplots, Regression lines and Prediction intervals

*Cars train data set*



# References

The R Project for Statistical Computing:

<https://www.r-project.org/>

Python:

<https://www.python.org/>

course notes