

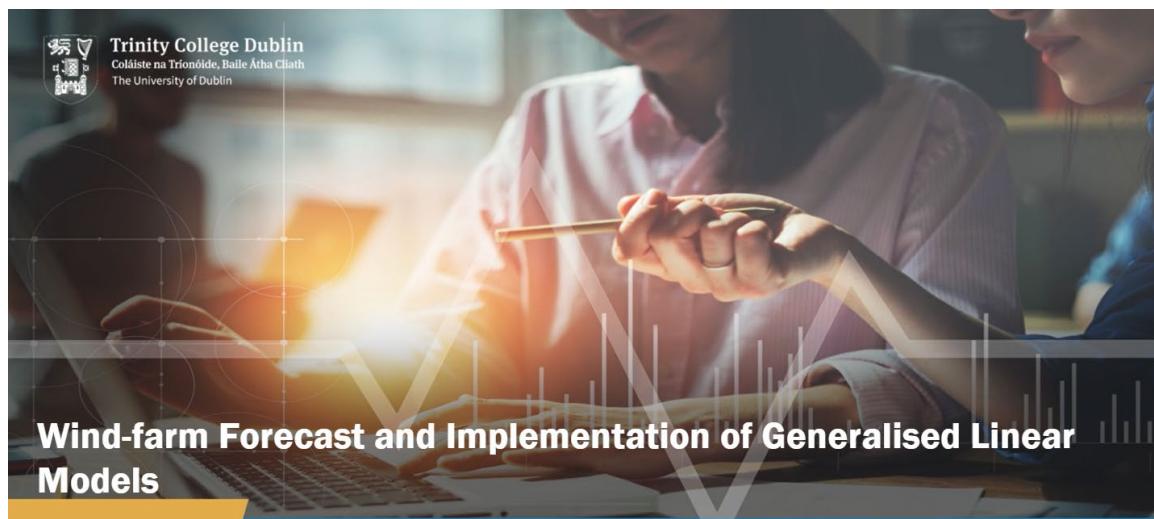
Wind-Farm Forecast and Implementation of Generalised Linear Models

Slide 1:	Introduction	3
Slide 2:	Wind Farm Forecast Example: Dataset and Considerations	3
Tab 1:	The Wind Turbine Dataset and Considerations	4
Slide 3:	Further Predictors	10
Slide 4:	Response Variable	11
Slide 5:	Logistic Function and Likelihood	12
Slide 6:	Logistic Regression	12
Slide 7:	Slide 7: Deriving the Model	13
Slide 8:	What is the Appropriate Likelihood Model?	14
Tab 1:	Likelihood	14
Slide 9:	Linear Regression and Matrix Formulation	19
Tab 1:	Maximising the Log-Likelihood.....	19
Tab 2:	Minimisation.....	20
Tab 3:	Solution of Parameter Estimation for Linear Regression.....	20
Tab 4:	Matrix Form	21
Tab 5:	Linear Regression in Matrix Form.....	21
Slide 10:	Turbine Problem: Matrix Notation	22
Slide 11:	Derivative of the Log-likelihood with Respect to Beta	22
Tab 1:	Derivative of mu With Respect to Beta	23
Tab 2:	Derivative of the Log-likelihood With Respect to Rho	24
Slide 12:	Finding the Maximum	24
Slide 13:	R Code: Log-likelihood	25
Slide 14:	R Code: Gradient of Log-likelihood	26
Slide 15:	R Code: Optimisation (1)	27
Slide 16:	R code: Optimisation (2).....	28
Slide 17:	R Code: Model Performance	29
Slide 18:	Akaike Information Criterion (AIC).....	30
Slide 19:	R Code: What About Time?.....	30
Slide 20:	R Code: Hours	32
Slide 21:	R Code: Second Model	33
Slide 22:	R Code: Model Performance	34
Slide 23:	R Code: Can We Do Better?	35



Slide 24: R Code: Model Performance with Additional Predictors	36
Slide 25: Statistics	37
Slide 26: Summary.....	37

Slide 1: **Introduction**



Trinity College Dublin
Coláiste na Tríonóide, Baile Átha Cliath
The University of Dublin

Wind-farm Forecast and Implementation of Generalised Linear Models

Presenter: Alessio Benavoli
Duration: 22:37
School: Computer Science and Statistics

My name is **Alessio Benavoli** and I am **the instructor for this session**. In this session, we apply generalised linear models to the problem of wind power forecasting. This problem cannot be formulated as linear regression, or as logistic regression, or as any other generalised linear model implemented in the R function “`glm`”. For this reason, we will derive the model from scratch both theoretically and practically in R. This will allow us to address the wind power forecasting problem but also to understand (indirectly) what is inside the “`glm`” function in R.

Slide 2: **Wind Farm Forecast Example: Dataset and Considerations**



Wind-farm Forecast Example: Dataset and Considerations

2 of 26

- Our goal is to predict the power output of a wind turbine using weather forecast data as covariates.
- These data are:
 - Wind velocity
 - Direction



Click the tab to learn more. Then, click Next to continue.



Wind Turbine Dataset and Considerations

In this session, we will face the problem of predicting the power generated by a wind farm. More precisely, our Goal is to predict the power output of a wind turbine using weather-forecast data (wind velocity and direction) as covariates.

Click the tab to learn more about this dataset and other considerations.

Tab 1: The Wind Turbine Dataset and Considerations

Wind Turbine Dataset and Considerations (1/9)

Wind Turbine Data Set



- The dataset includes:
 - Historical data from a wind-farm
 - Weather forecast data
- For each turbine and each time stamp, the dataset includes:
 - The wind forecast along zonal and meridional directions at heights of 10 metres and 100 metres above ground

Wind forecast	
Zonal	West-east direction
Meridional flow	North-south direction

We will use a dataset which includes historical data from a wind-farm, together with weather-forecast data.

The data also includes the wind forecast (delivered by a Weather-forecast provider) along the zonal and meridional directions at heights of 10 and 100 meters above ground (four measurements) for each turbine and for each time stamp.

Note that Zonal denotes the west-east direction and Meridional flow denotes the north-south direction.

Tab 1.1: Dataset Size

Wind Turbine Dataset and Considerations (2/9)

Dataset Size



- The dataset contains 168,000 hourly measurements made between:
 - 1 January 2012 and 1 December 2013 (100 weeks)
- Each measurement consists of:
 - One turbine power output
 - The four associated wind forecasts
 - A time stamp

! **Turbine power outputs are normalised to their nominal capacities, making TARGETVAR a number between 0 and 1.**

Variable	Definition
ZONEID	Zone (turbine) ID
TIMESTAMP	Date and time of observation
U10	Zonal wind velocity at 10m above ground
V10	Meridional wind velocity at 10 m above ground
U100	Zonal wind velocity at 100m above ground
V100	Meridional wind velocity at 100m above ground
TARGETVAR	Output of wind turbine, as a fraction of maximum capacity

The data set contains 168,000 hourly measurements made between January 1, 2012 and December 1, 2013 (100 weeks). Each measurement consisted of one turbine power output, plus the four associated wind forecasts, plus a time stamp. You can see the variables in the table u10 and u100 correspond to the zonal component of wind at 10m and, respectively, 100m. v10 and v100 correspond to the meridional component of wind at 10m and, respectively, 100m.

The TARGET variable is the produced power. Note that, the turbine power outputs are normalized to their nominal capacities, making TARGETVAR a number between 0 and 1.

Tab 1.2: Train and Test

Wind Turbine Dataset and Considerations (3/9)

Train and Test



- The dataset is divided into training and testing components.
- The goal is:
 - To learn a regression model using the training dataset
 - To predict the power for the instances in the test dataset

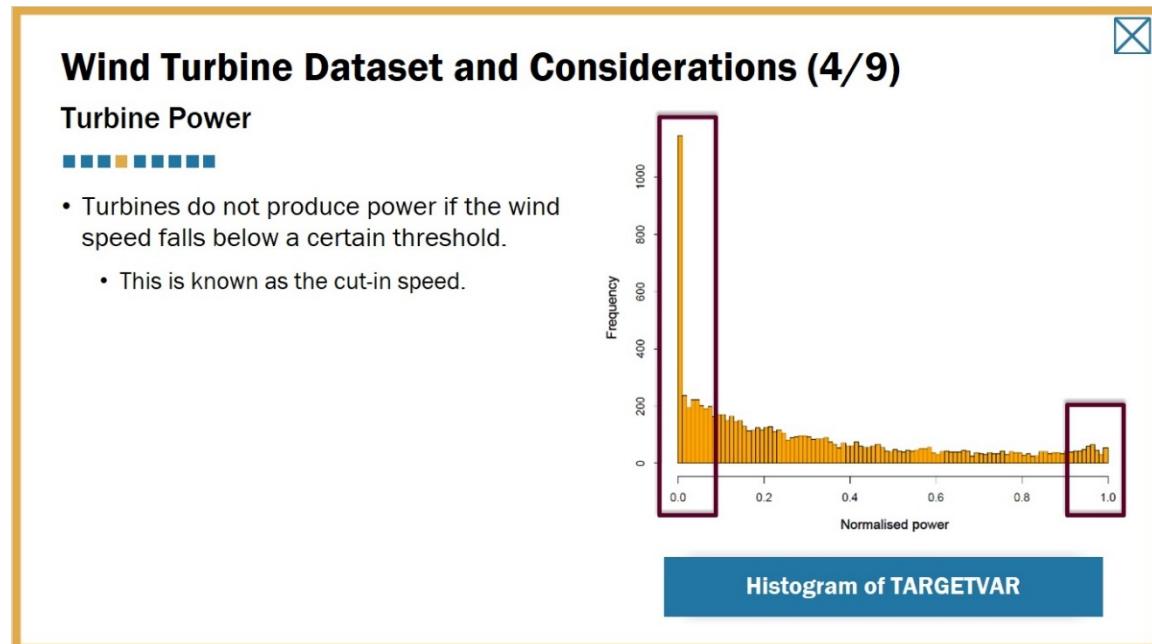
- This division allows us to assess the predicting performance of our regression model using unseen data.



The dataset is divided into training and testing components. Our goal is to learn a regression model using the data in the training dataset and, then, predict the power for the instances in the test dataset.

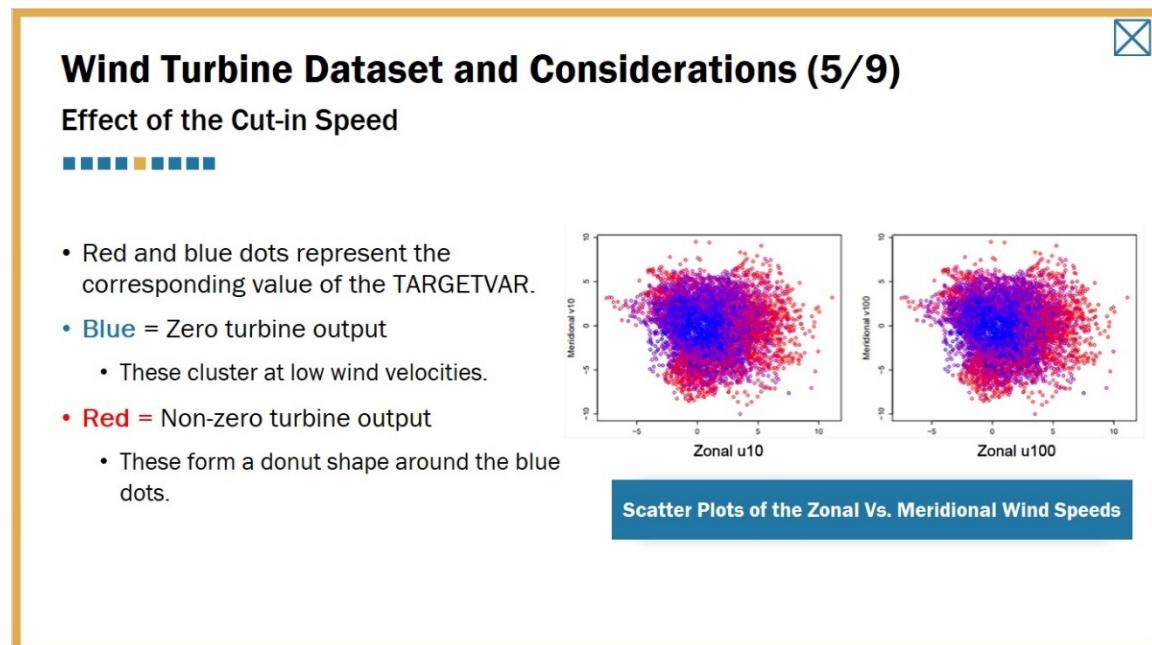
In this way we can assess the predicting performance of our regression model using unseen data, that is data were not used during fitting (that is computing MLE).

Tab 1.3: **Turbine Power**



The histogram shows two characteristics of this dataset: a large spike near 0, and a bump near 1. The bump corresponds to turbines working near full capacity. The large spike corresponds to low wind velocities. Turbines won't produce any power if the wind speed falls below a threshold known as the *cut-in speed*.

Tab 1.4: **Effect of the Cut-in Speed**



These are the scatter plots of the zonal versus meridional wind speed at 10meters (left) and 100meters (right). The colours represent the corresponding value of the TARGETVAR: blue dots correspond to zero turbine output and cluster at low wind velocities. Red dots correspond to non-zero turbine output and form a donut shape around the blue ones.

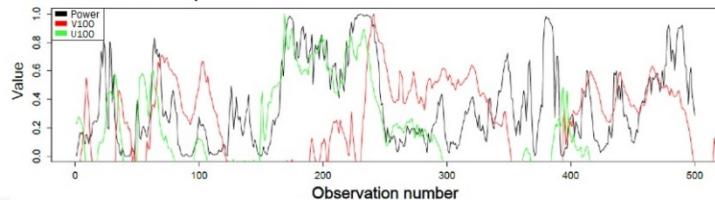
Tab 1.5: Time Series

Wind Turbine Dataset and Considerations (6/9)

Time Series



- Peaks in wind speed correspond to peaks in turbine power.
- The effect of the cut-in speed is visible.
 - There is zero turbine output at low velocity.
 - There is non-zero turbine output at high velocity.
- Wind direction affects turbine power.



Time Series

Peaks in wind speed correspond to peaks in turbine power, and the effect of the cut-in speed is also visible (periods of zero turbine output associated with low, but non-zero wind speed). Note also the effect of the wind direction, that is the effect of the zonal (u_{100}) or meridional wind (v_{100}) on the turbine power. As can be seen, a strong wind corresponds to a large produced power.

Tab 1.6: Predictors

Wind Turbine Dataset and Considerations (7/9)

Predictors



- The wind turbine:
 - Depends on the magnitude of the wind to function
 - Can change orientation, depending on the wind
- We will consider the following predictors instead of working directly with v10, u10, v100, u100:

$$ws10 = \sqrt{v10^2 + u10^2}, \quad ws100 = \sqrt{v100^2 + u100^2}, \quad dir100 = \tan^{-1}\left(\frac{v100}{u100}\right)$$

- We have verified that the variables dir100 and dir10 are highly correlated.
 - Therefore, we will only include one variable, dir100, in the model.

From knowledge about how a wind-turbine works, it is clear that the power mainly depends on the magnitude of the wind. The turbine can change orientation depending on the wind. So, we will consider the following predictors:

ws10 which is magnitude of the wind at 10m, ws100 which is magnitude of the wind at 100m, and dir100 which is the direction of the wind at 100m, instead of working directly with v10, u10, v100, u100.

We consider only dir100 and not dir10, because we have verified that these two variables are highly correlated. So, we should only include one of them in the model.

Tab 1.7: Forecasting Power Problem

Wind Turbine Dataset and Considerations (8/9)

Forecasting Power Problem



- We need to use wind measurements ws10, ws200 and dir100, to predict each turbine output.
 - Both wind magnitude and direction are important.



How do we handle the peak at zero turbine output?



How do we take into account the daily and seasonal effects of wind speed and direction?

- The dataset seems unbalanced towards zero.



We need to use wind measurements (ws10, ws100, dir100) for predicting each turbine output (because both magnitude and direction are important).

But this may cause difficulties because of the following aspects:

How should we handle the peak at zero turbine output? It seems the dataset is unbalanced towards zero. This may be a problem for our regression model.

Wind speed/direction has daily and seasonal effects (in winter wind is on average stronger), how do we take into account of these effects?

Tab 1.8: Model

Wind Turbine Dataset and Considerations (9/9)

Model



- We would like to use a linear regression model to predict power as a function of the inputs.

$$power = f(ws10, ws100, dir100, \boxed{daypart, season})$$

- We include “daypart” and “season” for the daily and seasonal effects.

- The linear regression model will not work here.
 - The response variable is bounded between 0 and 1.

Ideally, we would like to use a linear regression model to predict power as a function of the inputs

$\text{power} = f(\text{ws10}; \text{ws100}; \text{dir100}; \text{daypart}; \text{season})$

We have included time (“day part” and “season” to account for the daily and seasonality effects).

The linear regression model will not work here since the response variable is bounded between 0 and 1, we have seen a similar issue in the lesson about logistic regression. Linear regression predictions can be below 0 or above 1 and these predictions do not make sense for this problem (where the response variable is bounded between 0 and 1).

Slide 3: Further Predictors

 Further Predictors 3 of 26

- Here we focus on the predictors and on f .

$$\text{power} = f(\text{ws10}, \text{ws100}, \text{dir100}, \text{daypart}, \text{season})$$

Wind Magnitude = Numeric Variable

- Wind magnitude can be included in the linear model.

$$\beta_0 + \beta_1 \text{ws10} + \beta_2 \text{ws100} + \beta_3 \text{dir100}$$

$$\begin{aligned} &\beta_0 + \beta_1 \text{ws10} + \beta_2 \text{ws100} + \beta_3 \text{dir100} \\ &+ \beta_4 I_{\text{Spring}}(s) + \beta_5 I_{\text{Summer}}(s) + \beta_6 I_{\text{Fall}}(s) + \beta_7 I_{\text{Winter}}(s) \\ &+ \beta_8 I_{\text{Morning}}(d) + \beta_9 I_{\text{Afternoon}}(d) + \beta_{10} I_{\text{Evening}}(d) + \beta_{11} I_{\text{Night}}(d) \end{aligned}$$

Season and Daypart = Categorical Variables

$$\begin{aligned} \text{season} \in \{\text{Spring}, \text{Summer}, \text{Fall}, \text{Winter}\}, \\ \text{daypart} \in \{\text{Morning}, \text{Afternoon}, \text{Evening}, \text{Night}\} \end{aligned}$$

Before discussing the issue with the response variable further, let us focus on the predictors and on f .

$\text{power} = f(\text{ws10}; \text{ws100}; \text{dir100}; \text{daypart}; \text{season})$

The wind magnitude is a numeric variable and so we can include it in the linear model as shown in the top formula $\beta_0 + \beta_1 \text{ws10} + \beta_2 \text{ws100} + \beta_3 \text{dir100}$

while season can be Spring; Summer; Fall; Winter

$\text{season} \in \{\text{Spring}; \text{Summer}; \text{Fall}; \text{Winter}\};$

$\text{daypart} \in \{\text{Morning}; \text{Afternoon}; \text{Evening}; \text{Night}\}$

daypart can be Morning; Afternoon; Evening; Night

So, they are categorical (discrete) variables. Our goal here is to have a different linear model depending on the time of the day and on the season of the year and so we can use indicator functions, as seen in this last formula:

$$\begin{aligned} & \beta_0 + \beta_1 ws10 + \beta_2 ws100 + \beta_3 dir100 \\ & + \beta_4 / \text{Spring}^{(s)} + \beta_5 / \text{Summer}^{(s)} + \beta_6 / \text{Fall}^{(s)} + \beta_7 / \text{Winter}^{(s)} \\ & + \beta_8 / \text{Morning}^{(d)} + \beta_9 / \text{Afternoon}^{(d)} + \beta_{10} / \text{Evening}^{(d)} + \beta_{11} / \text{Night}^{(d)} \end{aligned}$$

where d stands for daypart and s for season.

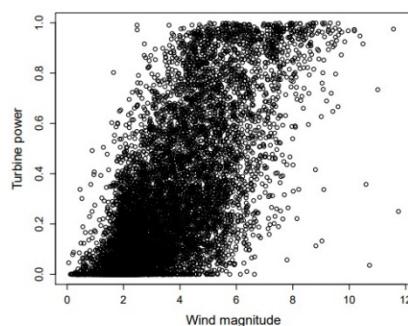
Slide 4: Response Variable



Response Variable

4 of 26

- Knowing the wind magnitude allows us to calculate the wind intensity.



Scatter Plot Between Turbine Power and Wind Magnitude



We see an S-shaped pattern. Can we squeeze the linear model through an S-shaped function to fit the data?

Looking at the scatter plot between turbine power (x-axis) and wind magnitude at 10m (y-axis), where the wind magnitude is defined as the square-root of the sum of the zonal and meridional wind components squared.

We can see an S-shape pattern, so can we use an approach similar to logistic regression to deal with the normalised turbine power? That is, we can try to squeeze the linear model through a S-shaped function to fit the data.

Slide 5: **Logistic Function and Likelihood**



Logistic Function and Likelihood

5 of 26

- Recall the logistic function $\phi(\beta_0 + \beta_1 x_1 + \cdots + \beta_m x_m) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \cdots + \beta_m x_m)}}$

and the likelihood we used for logistic regression (a product of Bernoulli distributions)

$$\begin{aligned} p(\mathcal{D} | \beta_0, \beta_1, \dots, \beta_m) &= \prod_{(x_1^{(i)}, \dots, x_m^{(i)}, y_i) \in \mathcal{D}} \phi(\beta_0 + \beta_1 x_1^{(i)} + \cdots + \beta_m x_m^{(i)})^{y_i} \\ &\quad \cdot (1 - \phi(\beta_0 + \beta_1 x_1^{(i)} + \cdots + \beta_m x_m^{(i)}))^{1-y_i} \end{aligned}$$

where ⁱ denotes the i-th observation in dataset D and x denotes m predictors.

We recall here from the session on logistic regression the logistic function and the likelihood we used for logistic regression, which is a product of Bernoulli distributions, where I used the superscript i to denote the i-th observation in the dataset.

$$\Phi(\beta_0 + \beta_1 x_1 + \cdots + \beta_m x_m) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \cdots + \beta_m x_m)}}$$

I have assumed we have m predictors denoted with x. In the turbine power prediction problem defined previously, m=11 and the predictors are ws10, ws100, dir100 and the indicators.

Slide 6: **Logistic Regression**



Logistic Regression

6 of 26

- The logistic regression model is not appropriate for the turbine problem, because:

$$\phi(\beta_0 + \beta_1 x_1 + \cdots + \beta_m x_m)$$

is the probability of y=1

$$1 - \phi(\beta_0 + \beta_1 x_1 + \cdots + \beta_m x_m)$$

is the probability of y=0

In the wind example:

- y can be any real number between 0 and 1
 - Therefore, we lose the interpretation of ϕ as a probability of y being equal to 1.

Is a Logistic regression model appropriate for the turbine power prediction problem? No, because in logistic regression, we interpreted Phi as the probability that y is equal to one and 1-Phi as the probability that y is zero.

But in the wind example y can be any real number between 0 and 1, and so we lose this interpretation of Phi as a probability of y being equal to one.

Slide 7: Slide 7: Deriving the Model



Deriving the Model

7 of 26

- We need to derive the appropriate model for the turbine power prediction.
- We can directly interpret $\phi(\beta_0 + \beta_1 x_1 + \cdots + \beta_m x_m)$ as the value of y.
 - This is correct, as ϕ is a real number between 0 and 1.
 - We can assume that $y = \phi$.

$$y \stackrel{?}{=} \phi(\beta_0 + \beta_1 x_1 + \cdots + \beta_m x_m)$$



We need to derive the appropriate model for the turbine power prediction. We could directly interpret the output of the logistic function, that is the value Phi

as the value of y. This is correct because Phi is a real number between 0 and 1, so we could assume that y is equal to Phi.

Slide 8:

What is the Appropriate Likelihood Model?



What is the Appropriate Likelihood Model?

8 of 26

- In practice, we can find betas that satisfy
 $y_i = \phi(\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_m x_m^{(i)})$
 for all data points, if all the y's lie exactly on the logistic function.

We need to define an appropriate likelihood function.

- We cannot use normal distribution.
- We need another probability density function to use as likelihood.

- In practice, linear regression will never find betas that satisfy

$$y_i = \phi(\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_m x_m^{(i)})$$

due to observations errors.

- Instead, we ask that the above equality only has to hold on average.

$$E[y_i] = \phi(\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_m x_m^{(i)})$$



Finding the Appropriate Likelihood Model



Click the tab to learn more. Then click Next to continue.

In practice, we can find betas that satisfy the set of equations in the slide. Note that we have an equation for each observation "i" in the dataset. These N equations can all be satisfied only if all the ys lie exactly on the logistic function. As for linear regression, this will never happen in practice due to observation errors. So, to account for these errors, we won't ask our model to satisfy the above set of equations. Similar to what was done for linear regression, we will ask the above equality only to hold on average, so that Phi is considered only be the average value of the variable y, as you can see from the formula.

Click the tab to learn more about how we do this. Then click next to continue.

Tab 1: Likelihood

Finding the Appropriate Likelihood Model (1/6)

Likelihood



- Assuming the likelihood:

$$p(y|\phi(\beta_0 + \beta_1 x_1 + \dots + \beta_m x_m)) \\ = \frac{1}{B(\mu\rho, (1-\mu)\rho)} y^{\mu\rho-1} (1-y)^{(1-\mu)\rho-1}$$

with

$$\mu = \phi(\beta_0 + \beta_1 x_1 + \dots + \beta_m x_m)$$

and rho is a positive scalar (parameter).

- This implies that the expectation of y = mu:

$$E[y] = \mu = \phi(\beta_0 + \beta_1 x_1 + \dots + \beta_m x_m)$$

This is exactly what we wanted.



We will do that by defining an appropriate likelihood function. We cannot use a Normal distribution because both y and Φ are numbers between 0 and 1. We need another probability density function to use as likelihood. We will consider the likelihood you see in the first equation on this slide, where μ is equal to the logistic function Φ . This is a probability density function of y , and one can show that the mean of this probability density function is equal to μ , that is the logistic function. This is exactly what we want.

Tab 1.1: Beta Distribution

Finding the Appropriate Likelihood Model (2/6)

Beta Distribution



Beta Distribution:

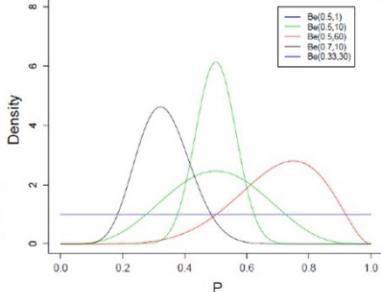
The probability density function of the variable y parametrised by μ and ρ

$$\text{Beta}(y; \mu, \rho) = \frac{1}{B(\mu\rho, (1-\mu)\rho)} y^{\mu\rho-1} (1-y)^{(1-\mu)\rho-1}$$

- μ is the mean of the distribution.
- ρ is the strength parameter.

Beta Function

$0 < y < 1, \quad \mu \in (0, 1), \quad \rho > 0$



Beta Distribution

This probability density function of the variable y parametrized by μ and ρ is called Beta distribution. You can see its probability density function in the formula. It has two parameters μ , which is the mean of the distribution, and ρ which is the strength parameter, whose domain is defined below. At the increase of ρ , the Beta density concentrates more around the mean. As you can see from the plots, comparing for instance the two green lines corresponding to a Beta with the same $\mu=0.5$ and different ρ .

This will be our likelihood. What is the denominator in the fraction in the top formula? It is called Beta function.

Tab 1.2: Beta Function

Finding the Appropriate Likelihood Model (3/6)

Beta Function



- The two-arguments function

$$B(\mu\rho, (1 - \mu)\rho)$$

is the normalisation constant of the Beta distribution.

In R, compute as follows:

```
1 beta(1,1) # result 1.0
2 beta(1,5) # result 0.2
```

- It is defined as the integral of unnormalised probability density function (PDF)

$$B(\mu\rho, (1 - \mu)\rho) = \int_0^1 y^{\mu\rho-1} (1 - y)^{(1-\mu)\rho-1} dy$$

The two-argument function you saw in the previous slide in the denominator of the fraction appearing in the formula of the Beta density and that you see in the first equation here

is the normalization constant of the Beta distribution. It is defined as the integral of the unnormalised PDF.

In R, you can compute calling the function “beta” as you can see in the slide.

```
1 beta(1,1) # result 1.0
```

```
2 beta(1,5) # result 0.2
```

Tab 1.3: Likelihood of the Dataset

Finding the Appropriate Likelihood Model (4/6)

Likelihood of the Dataset



- Therefore, assuming the observations are independent and identically distributed, our likelihood is:

$$p(\mathcal{D}|\beta_0, \beta_1, \dots, \beta_m) = \prod_{(x_1^{(i)}, \dots, x_m^{(i)}, y_i) \in \mathcal{D}} \frac{1}{B(\mu_i \rho, (1 - \mu_i) \rho)} y_i^{\mu_i \rho - 1} (1 - y_i)^{(1 - \mu_i) \rho - 1}$$

where

$$\mu_i = \phi(\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_m x_m^{(i)})$$

- These are the unknown parameters:

$$\beta_0, \beta_1, \dots, \beta_m, \rho$$

- The betas are real-valued, while rho is positive real-valued.

Therefore, assuming the observations are independent and identically distributed, our likelihood is a product of Beta density, with one term for each observation in the dataset. The mean mu depends on i and it is equal to the logistic function Phi we use to squeeze the linear model into 0 and 1.

The unknown parameters for this regression model are the Betas and rho.

The betas are real-valued, while rho is positive real-valued.

Tab 1.4: MLE

Finding the Appropriate Likelihood Model (5/6)

MLE



- Parameter estimation is performed by maximum likelihood:

$$\begin{aligned} & \arg \max_{\beta_0, \beta_1, \dots, \beta_m, \rho} p(\mathcal{D}|\beta_0, \beta_1, \dots, \beta_m, \rho) \\ &= \arg \max_{\beta_0, \beta_1, \dots, \beta_m, \rho} \log p(\mathcal{D}|\beta_0, \beta_1, \dots, \beta_m, \rho) \end{aligned}$$

where

$$\begin{aligned} \log p(\mathcal{D}|\beta_0, \beta_1, \dots, \beta_m, \rho) &= - \sum_{i=1}^N \log B(\mu_i \rho, (1 - \mu_i) \rho) \\ &\quad + (\mu_i \rho - 1) \log(y_i) \\ &\quad + ((1 - \mu_i) \rho - 1) \log(1 - y_i) \end{aligned}$$

with

$$\mu_i = \phi(\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_m x_m^{(i)})$$

We work with log-likelihood to:

- Avoid numerical numbers

The log likelihood can become very small.

- Numerical precision errors may occur when implemented in R.

Parameter estimation is performed by maximum likelihood, that is we want to find the values (called argument) of the betas and rho which maximize the likelihood. This is equivalent to maximizing the log-likelihood. We work with the log-likelihood to avoid numerical problems. The likelihood can become very very small and so we can have numerical precision errors when we implement it in R.

The expression of the log-likelihood is given in the second formula.

Tab 1.5: **R glm**

Finding the Appropriate Likelihood Model (6/6)

R glm



This generalised linear model has been defined by:

- Squeezing the linear model through the logistic function
- Defining an appropriate likelihood for the turbine power problem

! As **glm** from R does NOT provide the implementation of this generalised linear model, we must implement everything from scratch.

We have just defined another generalised linear model, by squeezing the linear model through the logistic function and by defining an appropriate likelihood for our turbine power problem.

Does **glm** from R provide the implementation of this generalised linear model? No, it doesn't.

What do we do then? We can implement everything from scratch.

Slide 9:

Linear Regression and Matrix Formulation



Linear Regression and Matrix Formulation

9 of 26

Maximising the Log-likelihood

Minimisation

Solution of Parameter Estimation for Linear Regression

Matrix Form

Linear Regression in Matrix Form

Introduction

- We will now see how we can apply matrix formulation to concisely write the linear regression parameter estimation.

$$y = \beta_1 x + \beta_0$$



Click each tab to learn more. Then, click Next to continue.

In a previous session we looked at linear regression. Let's recall what we did for standard linear regression and how we can apply the matrix formulation to compactly write the solution of linear regression parameter estimation.

Click each tab to learn more. Then click next to continue.

Tab 1: Maximising the Log-Likelihood



Linear Regression and Matrix Formulation

9 of 26

Maximising the Log-likelihood

Minimisation

Solution of Parameter Estimation for Linear Regression

Matrix Form

Linear Regression in Matrix Form

Maximising the Log-likelihood



- Given the dataset $D = \{(x_i, y_i) : i = 1, \dots, N\}$
- The parameters β_0, β_1 (intercept and slope) are determined by squared-error loss function

$$\min_{\beta_0, \beta_1} \sum_{i=1}^N (y_i - x_i \beta_1 - \beta_0)^2$$

- This is equivalent to maximising the log-likelihood (a Normal distribution).

In linear regression, given the dataset D, the parameters beta (intercept and slope) are determined by minimising the squared-error loss function which is equivalent to maximising the log-likelihood (which is a Normal distribution).

Tab 2: Minimisation

Linear Regression and Matrix Formulation				
9 of 26				
Maximising the Log-likelihood	Minimisation	Solution of Parameter Estimation for Linear Regression	Matrix Form	Linear Regression in Matrix Form
Minimisation <div style="float: right;"><input checked="" type="checkbox"/></div> <ul style="list-style-type: none"> To calculate the minimum: <ul style="list-style-type: none"> We compute the first derivative with respect to each parameter beta $\frac{\partial(\cdot)}{\partial \beta_1} = -2 \sum_{i=1}^N (y_i - x_i \beta_1 - \beta_0) x_i$ and $\frac{\partial(\cdot)}{\partial \beta_0} = -2 \sum_{i=1}^N (y_i - x_i \beta_1 - \beta_0)$ We find the values of beta which make these derivatives = zero by solving the system of equations below. $0 = -2 \sum_{i=1}^N (y_i - x_i \beta_1 - \beta_0) x_i, \quad 0 = -2 \sum_{i=1}^N (y_i - x_i \beta_1 - \beta_0)$ 				

To compute the minimum we compute the first derivative with respect to each parameter alpha,beta. The derivatives are shown in the first two equations in the slide.

To find the minimum, we look for the values of alpha,beta which make the above derivatives equal to zero. In linear regression, we can find these values analytically by solving the system of equations:

Tab 3: Solution of Parameter Estimation for Linear Regression

Linear Regression and Matrix Formulation				
9 of 26				
Maximising the Log-likelihood	Minimisation	Solution of Parameter Estimation for Linear Regression	Matrix Form	Linear Regression in Matrix Form
Solution of Parameter Estimation for Linear Regression <div style="float: right;"><input checked="" type="checkbox"/></div> <ul style="list-style-type: none"> The optimal values for beta are shown below: The solution is $\hat{\beta}_1 = \frac{\sum_{i=1}^N x_i y_i - \bar{y} \bar{x}}{\sum_{i=1}^N x_i^2 - \bar{x} \bar{x}}, \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x},$ <p>where</p> $\bar{y} = \left(\sum_{i=1}^N y_i \right) / N, \quad \bar{x} = \left(\sum_{i=1}^N x_i \right) / N$				

The optimal values for beta are shown in the slide.

Tab 4: Matrix Form


Linear Regression and Matrix Formulation
9 of 26

Maximising the Log-likelihood	Minimisation	Solution of Parameter Estimation for Linear Regression	Matrix Form	Linear Regression in Matrix Form				
Matrix Form <ul style="list-style-type: none"> If there are m predictors, we can use the matrix form to find the optimal values of the parameters. $X = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \dots & x_m^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} & \dots & x_m^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(N)} & x_2^{(N)} & \dots & x_m^{(N)} \end{bmatrix}, \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}, \quad \vec{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_m \end{bmatrix}$ <ul style="list-style-type: none"> The subscript in the x = covariate index The superscript in the x = index of the observation in the dataset N = observations m = covariates (or predictors) 								
<p>Arrow-beta = vector which includes all unknown betas</p> <ul style="list-style-type: none"> We can write the sum of the squared errors in matrix form: $\sum_{i=1}^N (y_i - \beta_0 - \beta_1 x_1^{(i)} - \beta_2 x_2^{(i)} - \dots - \beta_m x_m^{(i)})^2 = (\mathbf{y} - \mathbf{X} \vec{\beta})^T (\mathbf{y} - \mathbf{X} \vec{\beta})$ <p>where arrow-beta is the vector which includes all unknown betas</p> <table border="1" style="margin-left: auto; margin-right: auto;"> <tr> <td>Subscript in x</td> <td>Covariate index</td> </tr> <tr> <td>Superscript in the x</td> <td>Index of the observation in the dataset</td> </tr> </table> <p>$\cdot {}^T$ = transpose</p>					Subscript in x	Covariate index	Superscript in the x	Index of the observation in the dataset
Subscript in x	Covariate index							
Superscript in the x	Index of the observation in the dataset							

If there are m predictors (instead of only one as in the previous slide), we can use the matrix form to find the optimal values of the parameters. We define the matrix capital X and the vector y as in the slide,

where the subscript in the x denotes the index of the covariate and the superscript in the x the index of the observation in the dataset. I am assuming there are N observations and m covariates (also called predictors).

We can write the sum of the squared errors in this matrix form, where the arrow above the beta is used to denote the vector which includes all the unknown betas.

Tab 5: Linear Regression in Matrix Form


Linear Regression and Matrix Formulation
9 of 26

Maximising the Log-likelihood	Minimisation	Solution of Parameter Estimation for Linear Regression	Matrix Form	Linear Regression in Matrix Form
Linear Regression in Matrix Form <ul style="list-style-type: none"> By minimising the sum of the squared errors with respect to the vector arrow-beta, we find the familiar solution of linear regression in matrix form. 				
<ul style="list-style-type: none"> We need to solve $\arg \min_{\vec{\beta}} (\mathbf{y} - \mathbf{X} \vec{\beta})^T (\mathbf{y} - \mathbf{X} \vec{\beta})$				
<ul style="list-style-type: none"> The optimal solution is $\vec{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$				

By minimizing the sum of the squared errors with respect to the vector arrow-beta, we find the familiar solution of linear regression in matrix form

Slide 10: Turbine Problem: Matrix Notation


Turbine Problem: Matrix Notation
10 of 26

- We re-write our problem using matrix notation.

$$\arg \max_{\vec{\beta}, \rho} p(\mathcal{D} | \vec{\beta}, \rho) = \arg \max_{\vec{\beta}, \rho} \log p(\mathcal{D} | \vec{\beta}, \rho)$$

- where

$$\log p(\mathcal{D} | \vec{\beta}, \rho) = - \sum_{i=1}^N \log B(\mu_i \rho, (1 - \mu_i) \rho) + (\mu_i \rho - 1) \log(y_i) + ((1 - \mu_i) \rho - 1) \log(1 - y_i)$$

- with

$$\mu_i = \phi(X \vec{\beta})$$

Now that we have explained the matrix form (which we need to solve the turbine problem), let's go back to our problem, which I have rewritten using the matrix notation. Also, in this case the unknown is the vector arrow-beta.

Take time to view the information on this slide. Then click next to continue.

Slide 11: Derivative of the Log-likelihood with Respect to Beta


Derivative of the Log-likelihood with Respect to Beta
11 of 26

Derivative of mu with Respect to Beta

Derivative of Log-likelihood with Respect to Rho

Introduction

- We use the chain rule to compute the derivatives of each beta parameter by computing:
 - The derivative of log-likelihood with respect to μ_i
 - The derivative of μ_i with respect to beta j

There is a link to the digamma function wikipedia page in the Extend section of your session homepage.

The μ_i derivative depends on the digamma (ψ) function.

- This arises from differentiating the beta function.

$$\frac{d}{d\mu_i} \left(\log B(\mu_i \rho, (1 - \mu_i) \rho) + (\mu_i \rho - 1) \log(y_i) + ((1 - \mu_i) \rho - 1) \log(1 - y_i) \right) = -\psi(\mu_i \rho) + \psi((1 - \mu_i) \rho) + \rho \log \left(\frac{y_i}{1 - y_i} \right)$$

$$\frac{d}{d\beta_j} \log p(\mathcal{D} | \beta_0, \beta_1, \dots, \beta_m, \rho) = - \sum_{i=1}^N \frac{d}{d\mu_i} \left(\log B(\mu_i \rho, (1 - \mu_i) \rho) + (\mu_i \rho - 1) \log(y_i) + ((1 - \mu_i) \rho - 1) \log(1 - y_i) \right) \frac{d\mu_i}{d\beta_j}$$

Click each tab to learn more. Then, click Next to continue.

We use the chain rule to compute the derivatives of each beta parameter, by first computing the derivative of the log-likelihood with respect to μ_i and then computing the derivative of μ_i with respect to β_j , as you see in the formula below.

The derivative of the log-likelihood with respect to μ_i shown on the right depends on the digamma function, which arises from differentiating the Beta function (see wikipedia for details).

Click each tab to learn more about calculating these derivatives. Then, click next to continue.

Tab 1: Derivative of mu With Respect to Beta


Derivative of the Log-likelihood with Respect to Beta
11 of 26

Derivative of mu with Respect to Beta <ul style="list-style-type: none"> • The expression of μ_i is: $\mu_i = \frac{1}{1 + e^{-X_i\vec{\beta}}}$	Derivative of Log-likelihood with Respect to Rho
✖	

Derivative of mu with Respect to Beta

- Its derivative is:

$$\frac{d\mu_i}{d\beta_j} \frac{1}{1 + e^{-X_i\vec{\beta}}} = \frac{X_{ij}e^{-X_i\vec{\beta}}}{(1 + e^{-X_i\vec{\beta}})^2}$$

where

$$X_{ij}$$

The derivative of μ_i with respect to β_j can be easily calculated from the expression of the logistic function.

Tab 2: Derivative of the Log-likelihood With Respect to Rho

Derivative of the Log-likelihood with Respect to Beta
11 of 26

Derivative of mu with Respect to Beta	Derivative of Log-likelihood with Respect to Rho
---------------------------------------	--

Derivative of Log-likelihood with Respect to Rho

- Rho is a positive scalar, so the derivative is:

$$\begin{aligned} \frac{d}{d\rho} & \left(\log B(\mu_i\rho, (1 - \mu_i)\rho) + (\mu_i\rho - 1) \log(y_i) \right. \\ & \left. + ((1 - \mu_i)\rho - 1) \log(1 - y_i) \right) \\ & = \psi(\rho)(1 - \mu_i - (1 - \mu_i)) + \mu_i \log(y_i) + (1 - \mu_i) \log(1 - y_i) \end{aligned}$$

where ψ appears in the expression.

Rho is a positive scalar, so the derivative can be computed easily and again the digamma function appears in the expression of the derivative.

Slide 12: Finding the Maximum

Finding the Maximum
12 of 26

How do we find the values of betas and rho which maximise the log-likelihood?	Gradient-based Optimisation Method
---	------------------------------------

How do we find the values of betas and rho which maximise the log-likelihood?

- We could try to equate the first derivative to zero

$$0 = \frac{X_{ij}e^{-X_i\beta}}{(1+e^{-X_i\beta})^2} \left(-\psi(\mu_i\rho) + \psi((1 - \mu_i)\rho) + \rho \log\left(\frac{y_i}{1 - y_i}\right) \right)$$

- for all j from 1 to m, and

$$0 = \psi(\rho)(1 - \mu_i - (1 - \mu_i)) + \mu_i \log(y_i) + (1 - \mu_i) \log(1 - y_i)$$

- This system of equations can only be solved using a numerical method.

Gradient-based Optimisation Method

This optimisation method:

- Searches for the optimal value of parameters along the directions defined by the gradient of the function
- Stops when the gradient is close to zero

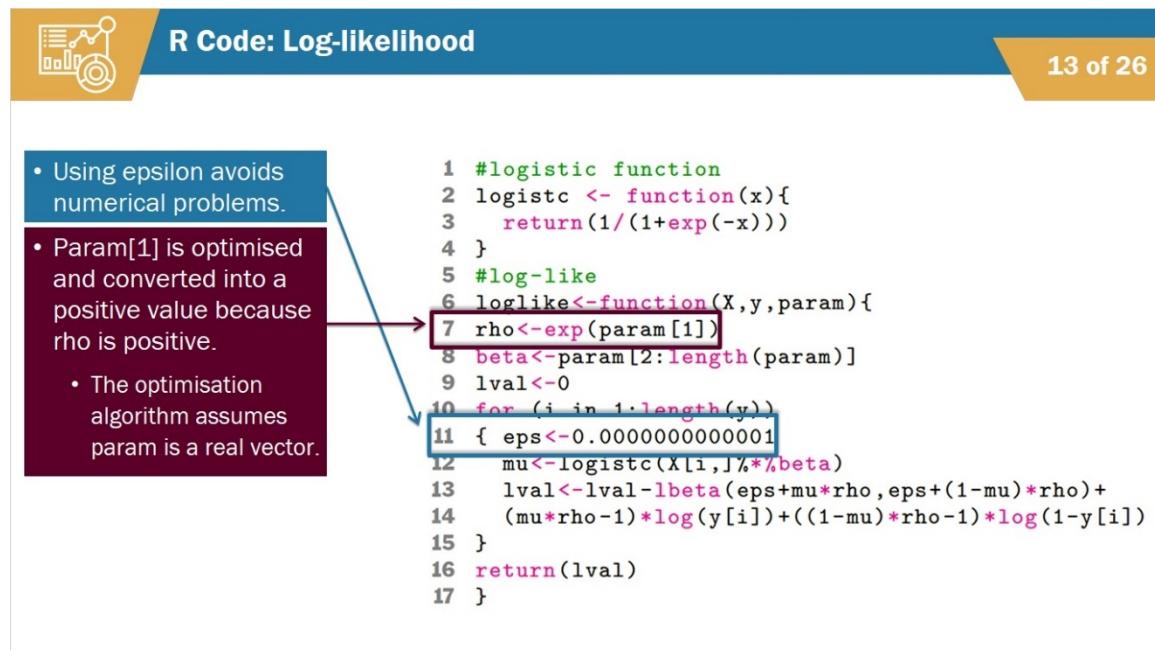
We want to find the values of betas and rho which maximize the log-likelihood. We could try to equate the first derivative to zero to do that and solve the system of equations you see in the slide.

Unfortunately, we cannot solve this system of equations analytically. Therefore, we use another strategy, a numerical method. We use a gradient-based optimisation method, that is an optimisation method which searches for the optimal value of the parameters

by iteratively changing the parameters along the directions defined by the gradient of the function (at the current value of the parameter). It stops when the gradient is close to zero.

We will now implement it in R.

Slide 13: R Code: Log-likelihood



R Code: Log-likelihood

13 of 26

- Using epsilon avoids numerical problems.
- Param[1] is optimised and converted into a positive value because rho is positive.
 - The optimisation algorithm assumes param is a real vector.

```

1 #logistic function
2 logistc <- function(x){
3   return(1/(1+exp(-x)))
4 }
5 #log-like
6 loglike <-function(X,y,param){
7   rho<-exp(param[1])
8   beta<-param[2:length(param)]
9   lval<-0
10  for (i in 1:length(y))
11    { eps<-0.00000000000001
12      mu<-logistc(X[i,]*%*%beta)
13      lval<-lval-lbeta(eps+mu*rho ,eps+(1-mu)*rho)+
14        (mu*rho-1)*log(y[i])+((1-mu)*rho-1)*log(1-y[i])
15    }
16  return(lval)
17 }
```

The following code implements the log-likelihood in R. AT the top you can see the implementation of the logistic function. The second R function is the log-likelihood.

```

1 #logistic function
2 logistc <- function(x){
3   return(1/(1+exp(-x)))
4 }
5 #log -like
6 loglike <-function(X,y,param){
7   rho<-exp(param [1])
8   beta<-param[2: length(param)]
9   lval<-0
10  for (i in 1:length(y))
11    { eps<-0.00000000000001
12      mu<-logistc(X[i,]*%*%beta)
13      lval<-lval-lbeta(eps+mu*rho ,eps+(1-mu)*rho)+
```

```

14 (mu*rho -1)*log(y[i])+((1-mu)*rho -1)*log(1-y[i])
15 }
16 return(lval)
17 }
```

We have introduced the epsilon to avoid numerical problems. Note the term $\exp(\text{param}[1])$, we optimize $\text{param}[1]$ which is real-variable and we use the \exp transform to convert it to a positive value because rho is positive. This is necessary because the optimization algorithm, we will use later, assumes param is a real vector.

Slide 14: R Code: Gradient of Log-likelihood

R Code: Gradient of Log-likelihood
14 of 26

- This is the implementation of the vector of the derivatives of the log-likelihood with respect to all the betas and rho.

```

1 dloglike<-function(X,y,param){
2   rho<-exp(param[1])
3   beta<-param[2:length(param)]
4   lval<-0
5   eps<-0.000000000001
6   for (i in 1:length(y))
7     { mui<-logistic(X[i,] %*% beta)
8      #derivative wrt betas
9      dbeta<-X[i,]*as.vector((-rho*digamma(eps+mui*rho) +
10 rho*digamma(eps+(1-mui)*rho) +
11 rho*log(y[i]/(1-y[i])))*exp(-X[i,] %*% beta)/((1+exp(-X[i,] %*% beta))^2))
12      #derivative wrt rho
13      drho<-(digamma(eps+rho)-mui*digamma(eps+rho*mui)-(1-mui)*digamma(eps+rho*(1-mui))+
14        mui*log(y[i])+(1-mui)*log(1-y[i]))*rho# this *rho comes from exp(param[1])
15   lval<-lval+c(drho,dbeta)
16 }
16 return(lval)
```

Implementation of Gradient of Log-Likelihood

This is the implementation of the gradient of the log-likelihood, that is the vector that includes the derivatives of the log-likelihood with respect to all the betas and rho.

```

1 dloglike <-function(X,y,param){
2   rho<-exp(param [1])
3   beta<-param[2: length(param)]
4   lval<-0
5   eps<-0.000000000001
6   for (i in 1:length(y))
7     { mui<-logistic(X[i,] %*% beta)
8      #derivative wrt betas
9      dbeta <-X[i,]*as.vector((-rho*digamma(eps+mui*rho) +
10 rho*digamma(eps+(1-mui)*rho) +
11 rho*log(y[i]/(1-y[i])))*exp(-X[i,] %*% beta)/((1+exp(-X[i,] %*% beta))^2))
12      #derivative wrt rho
13      drho<-(digamma(eps+rho)-mui*digamma(eps+rho*mui)-(1-mui)*digamma(eps+rho*(1-mui))+
14        mui*log(y[i])+(1-mui)*log(1-y[i]))*rho# this *rho comes from exp(param[1])
15   lval<-lval+c(drho,dbeta)
16 }
16 return(lval)
```

```

11 rho*log(y[i]/(1-y[i]))*exp(-X[i,]%^%^beta)/((1+exp(-X[i,]%^%^beta))^2))
12 #derivative wrt rho
13 drho<-digamma(eps+rho)-mui*digamma(eps+rho*mui)-(1-
    mui)*digamma(eps+rho*(1-mui))+
    mui*log(y[i])+(1-mui)*log(1-y[i]))*rho# this *rho comes from exp(param [1])
14 lval<-lval+c(drho ,dbeta)
15 }
16 return(lval)

```

Slide 15: R Code: Optimisation (1)


R Code: Optimisation (1)
15 of 26

- We will now apply the code to the wind power forecasting problem.
- We first use ws10 as predictor.
- We add rep(1,n) to insert the intercept.
- The Beta regression model assumes y is in (0,1).
 - Transformation $(y * (n - 1) + 0.5) / n$ is used to avoid y equalling 0 or 1, thereby dealing with the zero issues in TARGETVAR (turbine power).

```

1 #define the data
2 n<-length(train$p)
3 X<-cbind(rep(1,n),train$ws10)
4 y<-(train$p*(n-1)+0.5)/n
5 fr <- function(param){
6     return(loglike(X,y,param))
7 }
8 grr <- function(param){
9     return(dloglike(X,y,param))
10 }
```

- We define two auxiliary functions for the log-likelihood and its derivative, fr and grr.
 - These only depend on the parameter vector param.

We are now going to apply the code to the wind power forecasting problem.

```

1 #define the data
2 n<-length(train$p)
3 X<-cbind(rep(1,n),train$ws10)
4 y<-(train$p*(n-1)+0.5)/n
5 fr <- function(param){
6     return(loglike(X,y,param))
7 }
8 grr <- function(param){
9     return(dloglike(X,y,param))
10 }
```

For the moment we only use ws10 as predictor. We have added rep(1,n) which is a column of ones to insert the intercept term. The Beta regression model we defined assume y is in (0,1), so we use the transformation $(y^*(n-1)+0.5)/n$

where n is the number of observations, to avoid that y is equal to either 0 and 1. In this way, we deal with the issue of all the zeroes in TARGETVAR (turbine-power). We then define two auxiliary functions for the log-likelihood and its derivative, fr,grr, which only depend on the parameter vector param. This because the only variable we aim to optimize is param.

Slide 16: R code: Optimisation (2)


R Code: Optimisation (2)
16 of 26

- We run optim from R to numerically solve the optimisation problem.

? How do we get an initial value for the betas and rho?

```
1 res<-optim(rep(1,ncol(X)+1), fr,grr, method = "L-BFGS-B",control=list(fnscale=-1))
```

- We use optimiser “L-BFGS-B”.
- control=list(fnscale=-1) tells optim we are solving a maximisation problem.

- The solution to the optimisation problem is:

```
res$par
  0.95, -3.00,  0.56
```

rho	exp(0.95)
beta0	-3
beta1	0.56

We then run optim from R which numerically solves the optimization problem. We need an initial value for the betas and rho: a vector of ones in the code seen here (this is arbitrary).

```
1 res<-optim(rep(1,ncol(X)+1), fr,grr , method = "L-BFGS -B",control=list(fnscale=-1))
```

"L-BFGS-B" is the name of the optimizer we use and control=list(fnscale=-1) tells optim we are solving a maximization problem.

The solution to the optimisation problem is rho equal to exp(0.95), beta0 is -3 and beta1 0.56

Slide 17: R Code: Model Performance



R Code: Model Performance

17 of 26

- Using the fitted model, we can now compute:
 - Predictions on the test data
 - The mean absolute error (MAE) between the true value of the turbine power and the predicted one
 - The AIC for model comparison

MAE	0.160125
AIC	-15358.43

```

1 param<-res$par
2 param[1]<-exp(param[1])
3 Xtest<-cbind( rep(1,length(test$ws10)), test$ws10)
4 pred <-logistc(Xtest%*%param[2:length(param)])#prediction
5 MAE<-mean(abs(pred-test$p))
6 MAE
7 AIC<-2*length(res$par)-2*loglike(X,y,res$par)
8 AIC

```

We can now use the fitted model to compute predictions on the test data. We then use these predictions to compute the mean absolute error (MAE) between the true value of the turbine power and the predicted one. This allows us to quantify the accuracy of the model. Finally we compute the AIC for model comparison. You can see the result on this slide.

```

1 param <-res$par
2 param[1]<-exp(param [1])
3 Xtest <-cbind( rep(1,length(test$ws10)), test$ws10)
4 pred <-logistc(Xtest%*%param[2: length(param)])#prediction
5 MAE<-mean(abs(pred -test$p))
6 MAE
7 AIC<-2*length(res$par)-2*loglike(X,y,res$par)
8 AIC

```

$$MAE = 0.1601025$$

$$AIC = 15358.43$$

Slide 18:

Akaike Information Criterion (AIC)



Akaike Information Criterion (AIC)

18 of 26

- Let k be the number of parameters in the model.
- Let \hat{L} be the maximum value of the likelihood function for the model.
- The AIC value of the model is:

$$AIC = 2k - 2 \ln(\hat{L})$$



Given a set of candidate models, the preferred model is the one with the minimum AIC value.

In an earlier presentation, we describe the AIC to perform model selection, in particular to select the most important covariates for predicting the response variable. Recall: let k be the number of parameters in the model. Let

\hat{L} be the maximum value of the likelihood function for the model. Then the AIC value of the model is $2*k - 2 \log(\hat{L})$

Given a set of candidate models, the preferred model is the one with the minimum AIC value.

Slide 19:

R Code: What About Time?



R Code: What About Time?

19 of 26



Should we include “time” to account for daily and seasonality effects?

- We can add indicators and interaction terms.
- We make four indicators for the four seasons by using the library “lubridate” from R.

Season	Month Number
Winter	12, 1, 2
Spring	3, 4, 5
Summer	6, 7, 8
Autumn	9, 10, 11

```

1 library(lubridate)
2 new_date=ymd_hm(train$timestamp)
3 months<-month(new_date)
4 Iwinter<-rep(0.0,length(new_date))
5 ind<-which(months %in% c(12,1,2) )
6 Iwinter[ind]=1.0
7 Ispring<-rep(0.0,length(new_date))
8 ind<-which(months %in% c(3,4,5) )
9 Ispring[ind]=1.0
10 Isummer<-rep(0.0,length(new_date))
11 ind<-which(months %in% c(6,7,8) )
12 Isummer[ind]=1.0
13 Ifall<-rep(0.0,length(new_date))
14 ind<-which(months %in% c(9,10,11) )
15 Ifall[ind]=1.0

```

Should we include “time” to account for daily and seasonality effects? We may try to add indicators and interaction terms.

First, we need to make the indicators. We make 4 indicators for the 4 seasons, for Winter (month number 12,1,2), for Spring (month number 3,4,5) etc. To do that, we need to parse the dates in the timestamps, we use the library ``lubridate'' from R.

```

1 library(lubridate)
2 new_date=ymd_hm(train$timestamp)
3 months <-month(new_date)
4 Iwinter <-rep(0.0,length(new_date))
5 ind<-which(months %in% c(12,1,2) )
6 Iwinter[ind]=1.0
7 Ispring <-rep(0.0,length(new_date))
8 ind<-which(months %in% c(3,4,5) )
9 Ispring[ind]=1.0
10 Isummer <-rep(0.0,length(new_date))
11 ind<-which(months %in% c(6,7,8) )
12 Isummer[ind]=1.0
13 Ifall <-rep(0.0,length(new_date))
14 ind<-which(months %in% c(9,10,11) )
15 Ifall[ind]=1.0

```

Slide 20: R Code: Hours



R Code: Hours

20 of 26

- We use a similar code for parts of the day by creating indicators for:
 - Morning
 - Afternoon
 - Evening
 - Night

```

1 library(lubridate)
2 new_date=ymd_hm(train$timestamp)
3 hours<-hour(new_date)
4 Imorn<-rep(0.0,length(new_date))
5 ind<-which(hours %in% c(6,7,8,9,10,11) )
6 Imorn[ind]=1.0
7 Iafter<-rep(0.0,length(new_date))
8 ind<-which(months %in% c(12,13,14,15,16) )
9 Iafter[ind]=1.0
10 Ieven<-rep(0.0,length(new_date))
11 ind<-which(months %in% c(17,18,19,20,21,22) )
12 Ieven[ind]=1.0
13 Inight<-rep(0.0,length(new_date))
14 ind<-which(months %in% c(23,1,2,3,4,5) )
15 Inight[ind]=1.0
  
```

We use a similar code for dayparts.

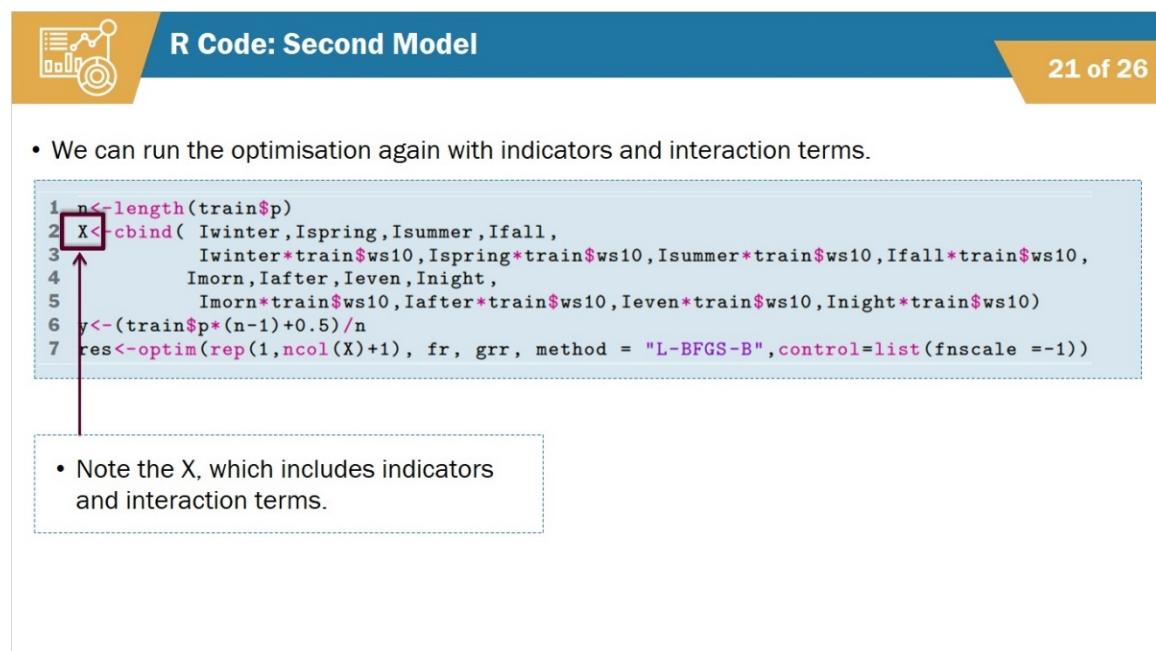
```

1 library(lubridate)
2 new_date=ymd_hm(train$timestamp)
3 hours <-hour(new_date)
4 Imorn <-rep(0.0,length(new_date))
5 ind<-which(hours %in% c(6,7,8,9,10,11) )
6 Imorn[ind]=1.0
7 Iafter <-rep(0.0,length(new_date))
8 ind<-which(months %in% c(12,13,14,15,16) )
9 Iafter[ind]=1.0
10 Ieven <-rep(0.0,length(new_date))
11 ind<-which(months %in% c(17,18,19,20,21,22) )
12 Ieven[ind]=1.0
13 Inight <-rep(0.0,length(new_date))
14 ind<-which(months %in% c(23,1,2,3,4,5) )
15 Inight[ind]=1.0
  
```

We make 4 indicators for morning, afternoon, evening and night.

Slide 21:

R Code: Second Model



• We can run the optimisation again with indicators and interaction terms.

```

1 n<-length(train$p)
2 X<-cbind( Iwinter,Ispring,Isummer,Ifall,
3           Iwinter*train$ws10,Ispring*train$ws10,Isummer*train$ws10,Ifall*train$ws10,
4           Imorn,Iafter,Ieven,Inight,
5           Imorn*train$ws10,Iafter*train$ws10,Ieven*train$ws10,Inight*train$ws10)
6 y<-(train$p*(n-1)+0.5)/n
7 res<-optim(rep(1,ncol(X)+1), fr, grr, method = "L-BFGS-B",control=list(fnscale =-1))

```

• Note the X, which includes indicators and interaction terms.

This is our model which has indicators and interaction terms.

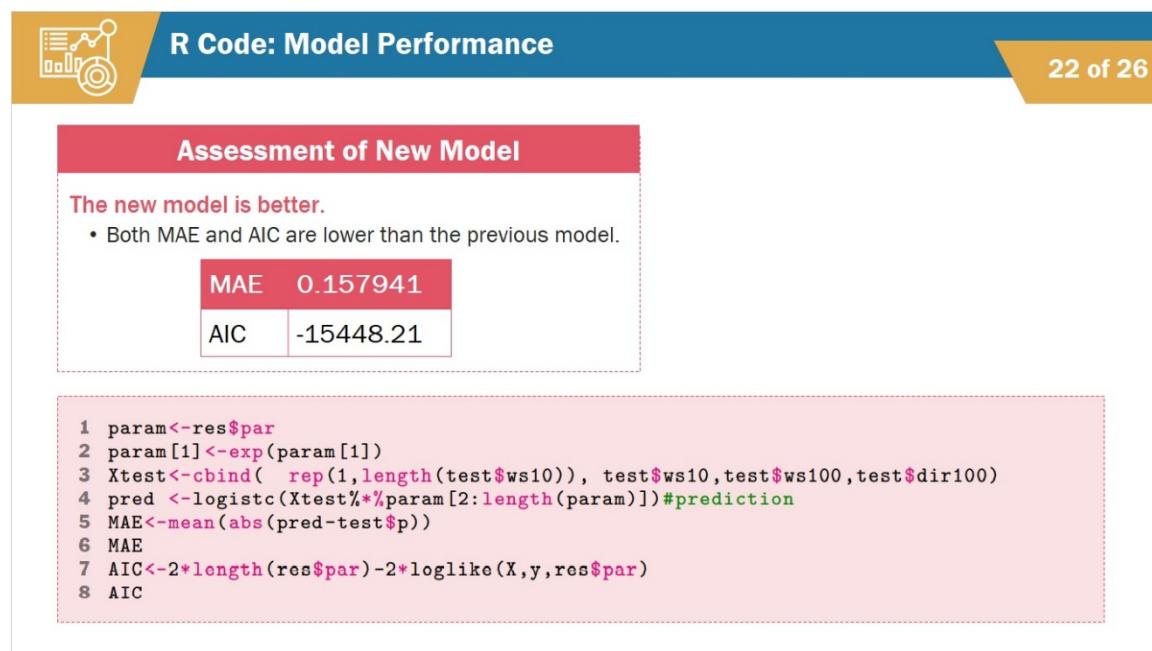
```

1 n<-length(train$p)
2 X<-cbind( Iwinter,Ispring,Isummer,Ifall,
3           Iwinter*train$ws10,Ispring*train$ws10,Isummer*train$ws10,Ifall*train$ws10,
4           Imorn,Iafter,Ieven,Inight,
5           Imorn*train$ws10,Iafter*train$ws10,Ieven*train$ws10,Inight*train$ws10)
6 y<-(train$p*(n-1)+0.5)/n
7 res<-optim(rep(1,ncol(X)+1), fr, grr, method = "L-BFGS -B",control=list(fnscale =-1))

```

Note in the code that, we only need to change X, so our code is very general (similar to `glm` in R).

Slide 22: R Code: Model Performance



R Code: Model Performance

Assessment of New Model

The new model is better.

- Both MAE and AIC are lower than the previous model.

MAE	0.157941
AIC	-15448.21

```

1 param<-res$par
2 param[1]<-exp(param[1])
3 Xtest<-cbind( rep(1,length(test$ws10)), test$ws10,test$ws100,test$dir100)
4 pred <-logistic(Xtest%*%param[2:length(param)])#prediction
5 MAE<-mean(abs(pred-test$p))
6 MAE
7 AIC<-2*length(res$par)-2*loglike(X,y,res$par)
8 AIC

```

We can assess the performance of the new model.

```

1 param <-res$par
2 param[1]<-exp(param [1])
3 Xtest <-cbind( rep(1,length(test$ws10)), test$ws10 ,test$ws100 ,test$dir100)
4 pred <-logistic(Xtest%*%param[2: length(param)])#prediction
5 MAE<-mean(abs(pred -test$p))
6 MAE
7 AIC<-2*length(res$par)-2*loglike(X,y,res$par)
8 AIC

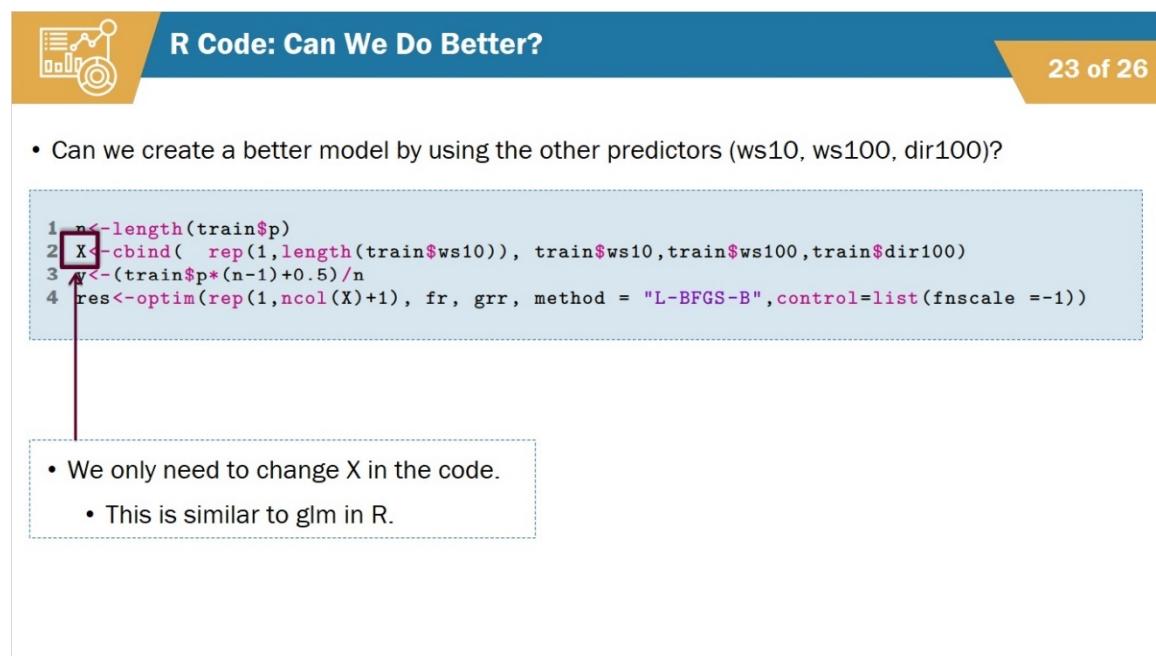
```

MAE = 0.1579741

AIC = -15448.21

The new model has a lower MAE and a lower AIC, so it is definitely better.

Slide 23: R Code: Can We Do Better?



R Code: Can We Do Better?

23 of 26

- Can we create a better model by using the other predictors (ws10, ws100, dir100)?

```

1 n<-length(train$p)
2 X<-cbind( rep(1,length(train$ws10)), train$ws10,train$ws100,train$dir100)
3 y<-(train$p*(n-1)+0.5)/n
4 res<-optim(rep(1,ncol(X)+1), fr, grr, method = "L-BFGS-B",control=list(fnscale =-1))

```

- We only need to change X in the code.
 - This is similar to glm in R.

Can we create a better model for instance, by using the other predictors (ws10, ws100, dir100)?

```

1 n<-length(train$p)
2 X<-cbind( rep(1,length(train$ws10)), train$ws10 ,train$ws100 ,train$dir100)
3 y<-(train$p*(n-1)+0.5)/n
4 res<-optim(rep(1,ncol(X)+1), fr, grr , method = "L-BFGS -B",control=list(fnscale =-1))

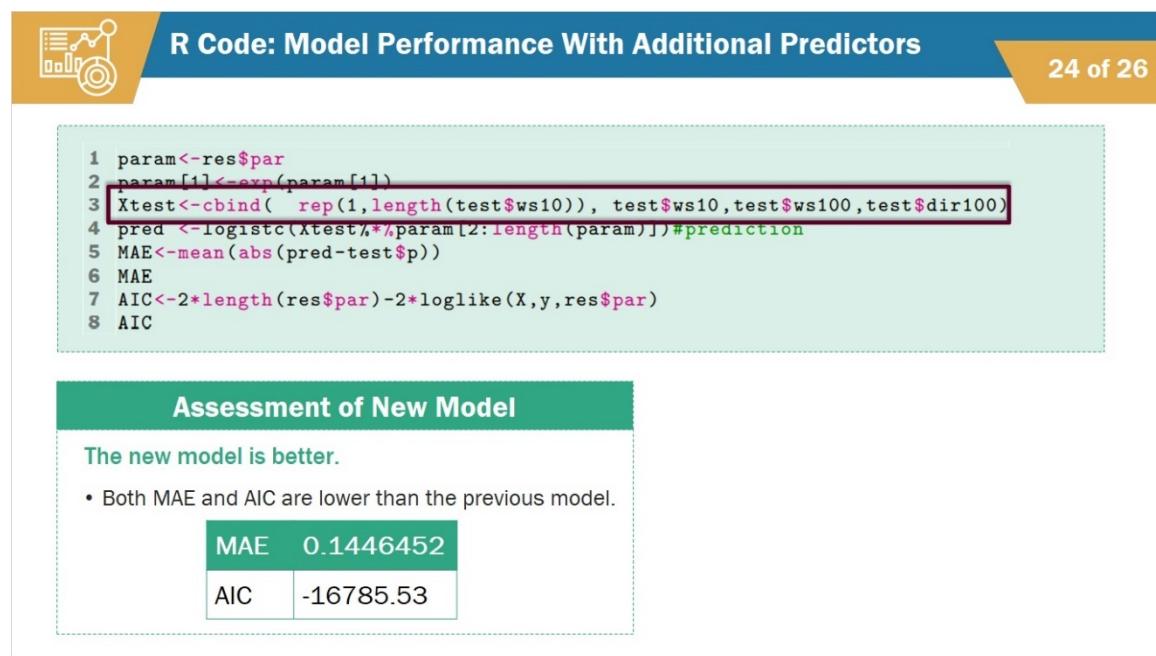
```

Note that we only need to change x, so our code is very general. It is similar to glm in R.

Take time to view the information on this slide. When you are ready, click next to continue.

Slide 24:

R Code: Model Performance with Additional Predictors



The screenshot shows an R script for model performance assessment. The code includes a section for Xtest definition and a prediction step, which is highlighted with a red box. Below the code, a green box titled 'Assessment of New Model' contains the message 'The new model is better.' and a bulleted list stating 'Both MAE and AIC are lower than the previous model.' To the right of the list is a table showing MAE and AIC values.

MAE	0.1446452
AIC	-16785.53

We can assess the performance of the new model that includes the predictors (ws10, ws100, dir100), note we changed the way Xtest is defined because the covariates of the new model include ws10, ws100 and dir100.

```

1 param <-res$par
2 param[1]<-exp(param [1])
3 Xtest <-cbind( rep(1,length(test$ws10)), test$ws10 ,test$ws100 ,test$dir100)
4 pred <-logistic(Xtest%*%param[2: length(param)])#prediction
5 MAE<-mean(abs(pred -test$p))
6 MAE
7 AIC<-2*length(res$par)-2*loglike(X,y,res$par)
8 AIC

```

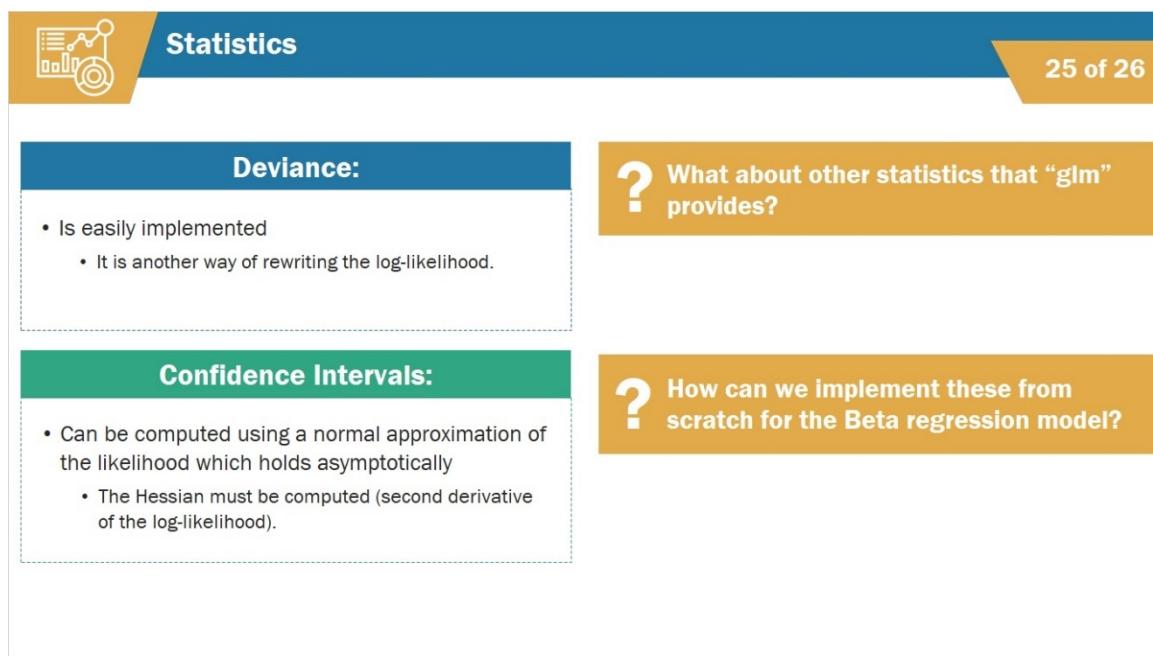
$$MAE = 0.1446452$$

$$AIC = -16785.53$$

The new model has a lower MAE and a lower AIC, so it is definitely better than both the previous ones.

Slide 25:

Statistics



Statistics

Deviance:

- Is easily implemented
 - It is another way of rewriting the log-likelihood.

Confidence Intervals:

- Can be computed using a normal approximation of the likelihood which holds asymptotically
 - The Hessian must be computed (second derivative of the log-likelihood).

What about other statistics that “glm” provides?

How can we implement these from scratch for the Beta regression model?

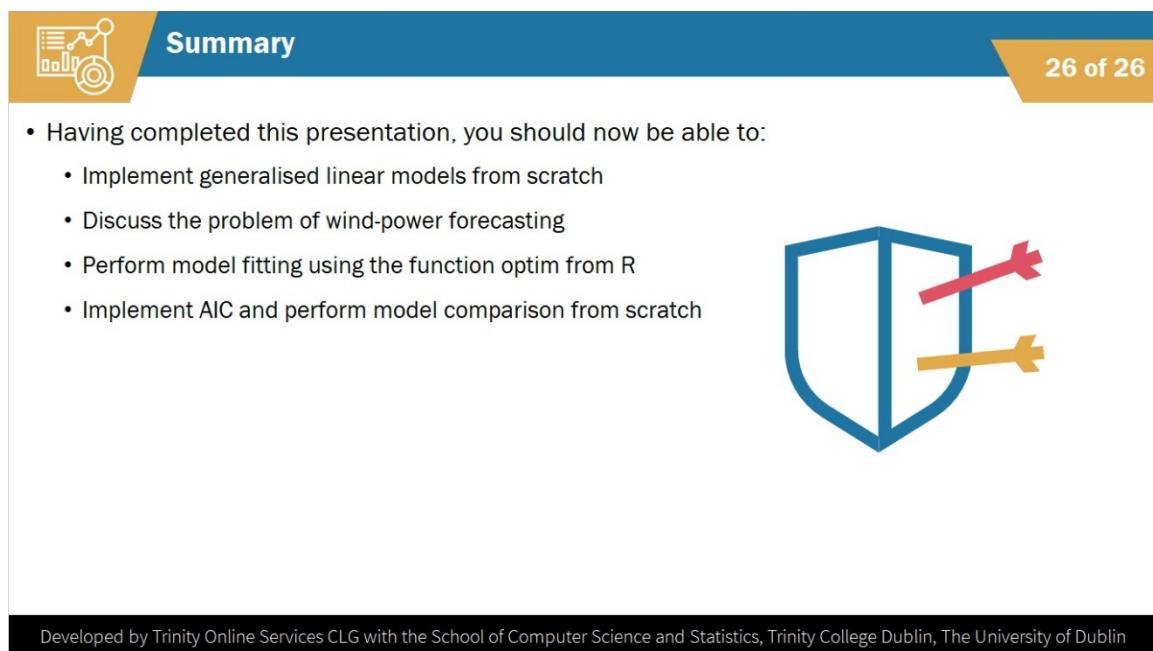
What about the other statistics “glm” provides us, for instance, deviance, confidence intervals etc., **how can we implement them from scratch for the Beta regression model?**

Deviance can be easily implemented because, as we saw for logistic regression, deviance is just a way of rewriting the log-likelihood.

Confidence intervals can be computed using a Normal approximation of the likelihood which holds asymptotically. In order to do that, we need to compute the Hessian of the log-likelihood, that is the second derivatives of the log-likelihood. It is not difficult, but it requires lengthy calculations and so we won't do it in this lesson.

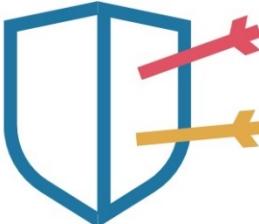
Slide 26:

Summary



Summary

- Having completed this presentation, you should now be able to:
 - Implement generalised linear models from scratch
 - Discuss the problem of wind-power forecasting
 - Perform model fitting using the function optim from R
 - Implement AIC and perform model comparison from scratch



In this presentation, we have developed and implemented from scratch in R a new generalised linear model to solve a -wind-power forecasting problem.

Having completed this presentation, you should now be able to:

Implement generalised linear models from scratch

Discuss the problem of wind-power forecasting

Perform model fitting using the function optim from R and

Implement AIC and perform model comparison from scratch