# Inf2 – Foundations of Data Science: Regression and inference – Generalised linear models

THE UNIVERSITY of EDINBURGH

**informatics**

FOUNDATIONS OF DATA SCIENCE

We want to investigate the relationship between the number of bikes hired in an hour and the mean temperature during that hour
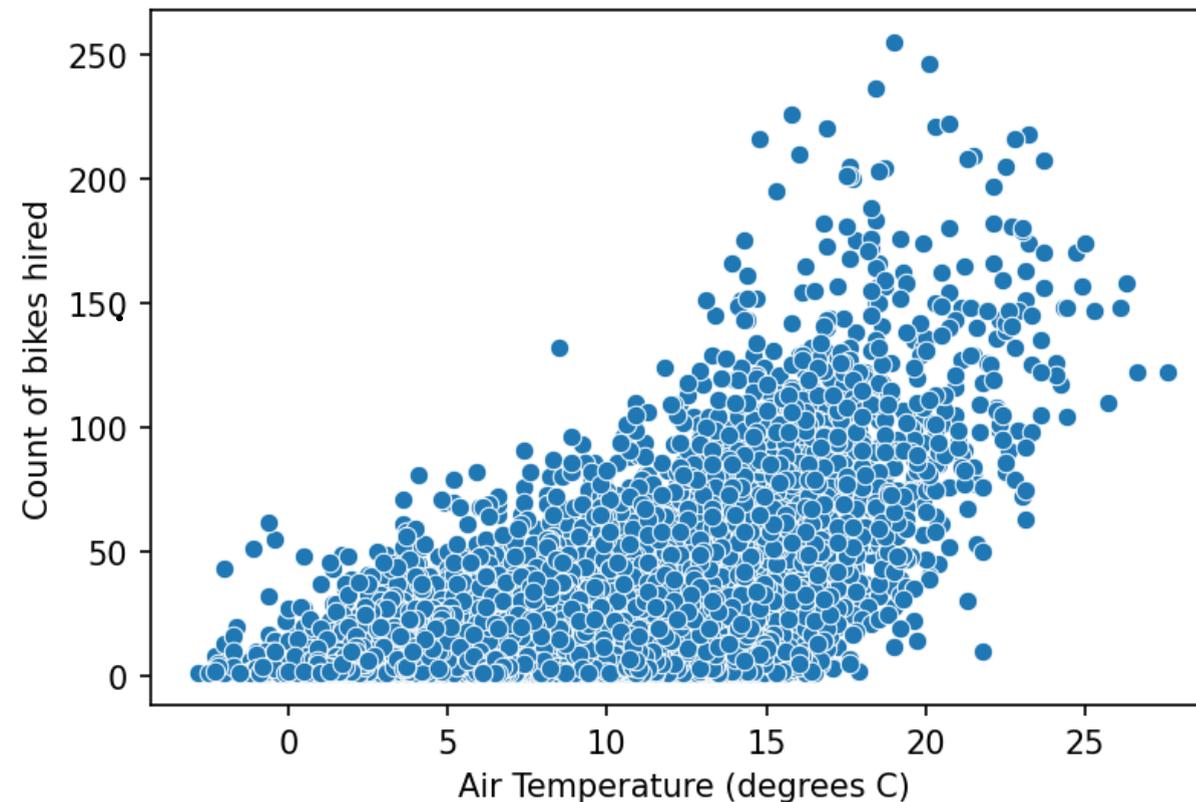
Is there a problem with using ordinary least squares linear regression to do this?

Data sources:
– Edinburgh Just Eat Bikes data 2020
– Edinburgh temperature observations, Met Office via MIDAS

# Overview

Monday

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

Today

0. Recap
1. Max likelihood with non-normal distributions
2. Poisson regression
3. Logistic regression and generalised linear models

# Inf2 – Foundations of Data Science: Regression and inference – Recap of max likelihood applied to linear regression

THE UNIVERSITY of EDINBURGH

informatics

FOUNDATIONS OF DATA SCIENCE

# Likelihood and log likelihood as a function of parameters

$$P(Y = y_1, \ldots, y_n \mid \mu, \sigma^2)$$

Data:

$y_1, \ldots, y_{10}$ drawn from

$$\mathcal{N}(0, 1)$$

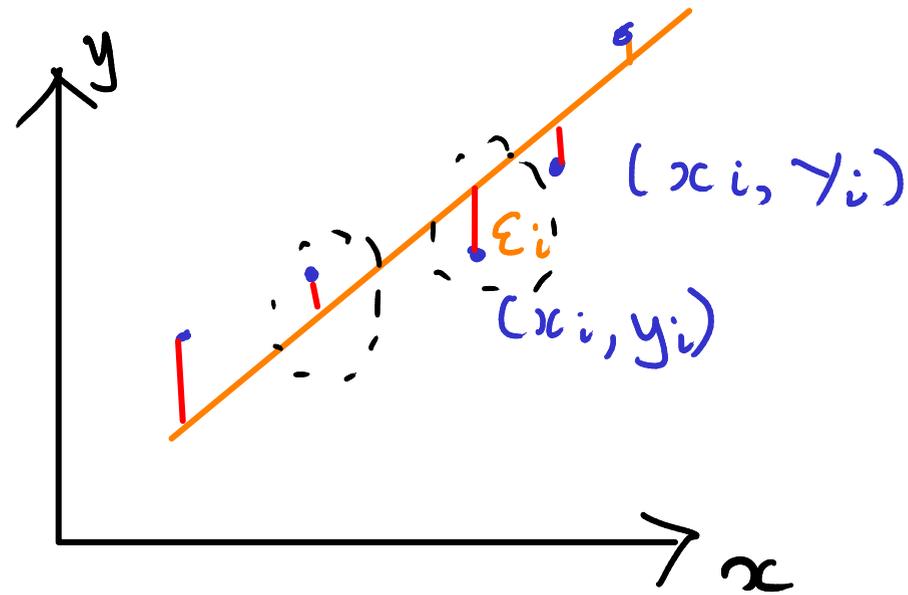| | Data |
|---|---|
| $y_1$ | 1.624345 |
| $y_2$ | -0.611756 |
| $y_3$ | -0.528172 |
| $y_4$ | -1.072969 |
| $y_5$ | 0.865408 |
| $y_6$ | -2.301539 |
| $y_7$ | 1.744812 |
| $y_8$ | -0.761207 |
| $y_9$ | 0.319039 |
| $y_{10}$ | -0.249370 |

Optimisation.



[Code]

# Application of max likelihood to linear regression

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

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

↑ residual

OR

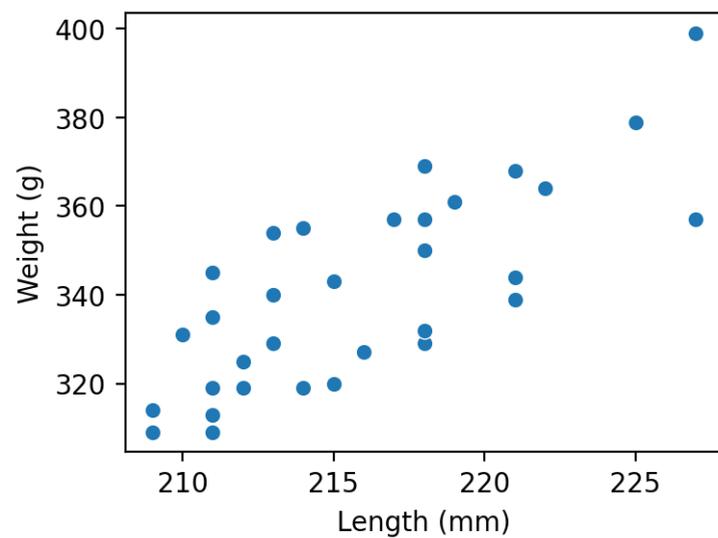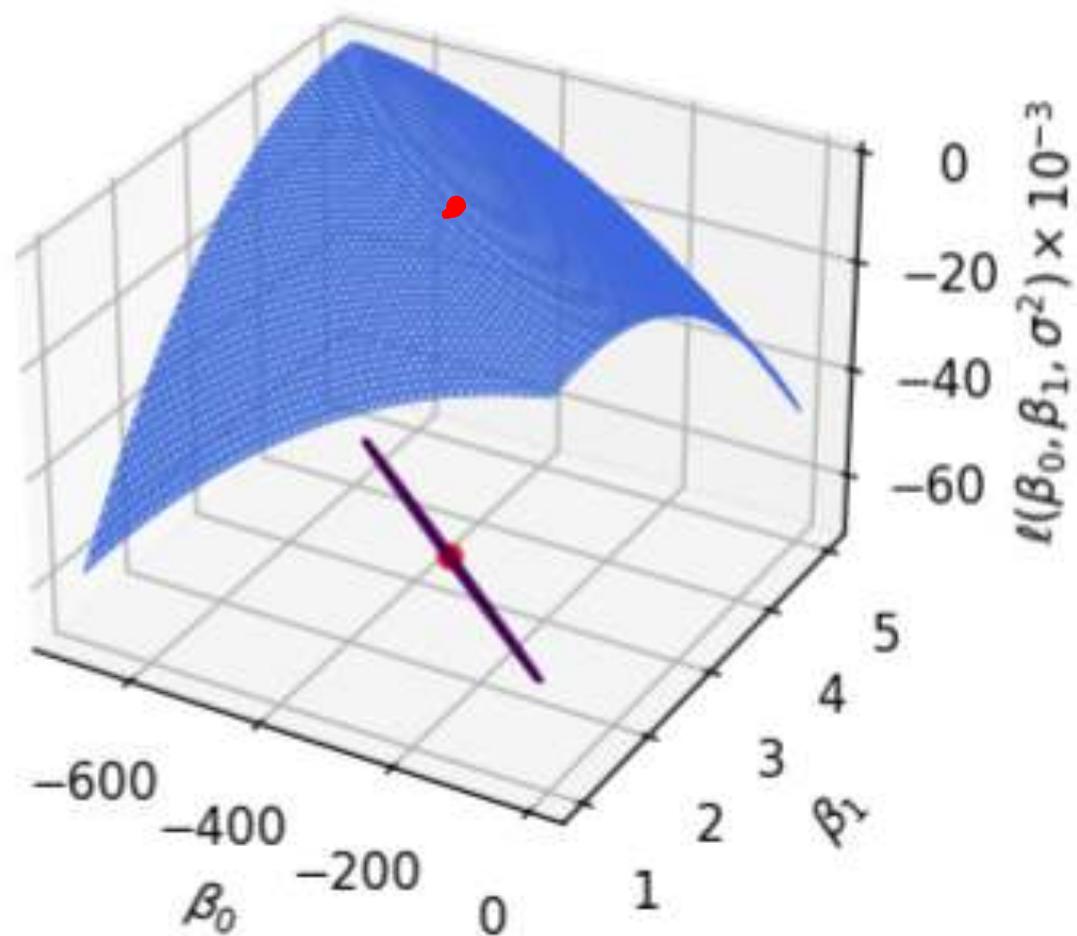$$Y_i \sim N(\overbrace{\beta_0 + \beta_1 x_i}^{\mu}, \sigma^2)$$

$$\ln p\left(\underline{Y} = y_1, \ldots, y_n \; ; \; x_1, \ldots, x_n \; \middle|\; \overbrace{\beta_0, \beta_1}^{\mu}, \sigma^2\right)$$

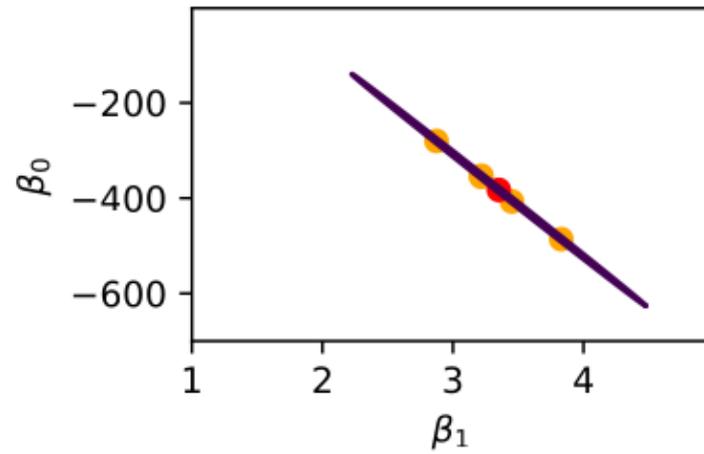$$= \sum_{i=1}^{n}\left(-\frac{1}{2}\ln 2\pi\sigma^2 - \frac{1}{2}\left(\frac{y_i - \beta_0 - \beta_1 x_i}{\sigma^2}\right)^2\right)$$
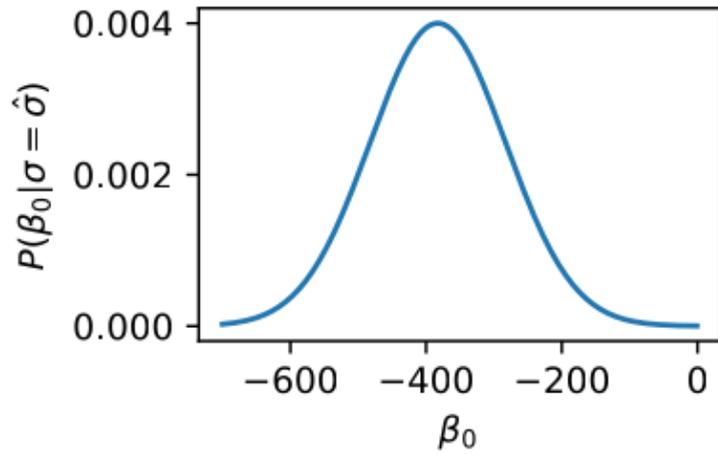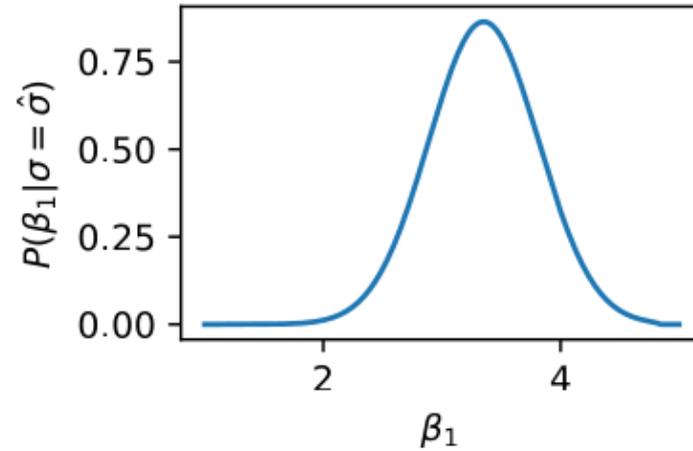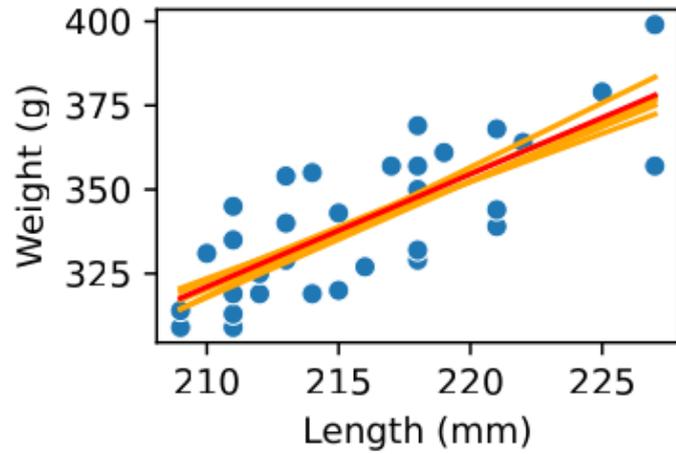
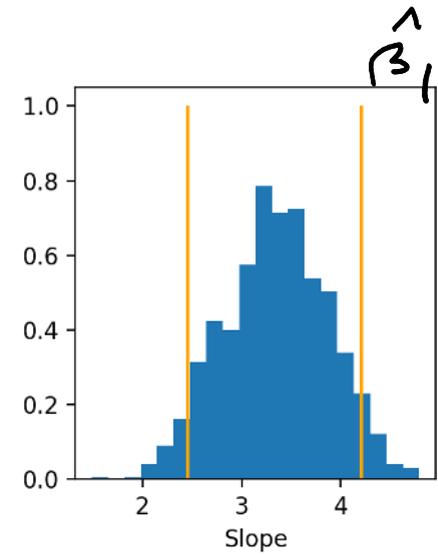# Log likelihood of coefficients
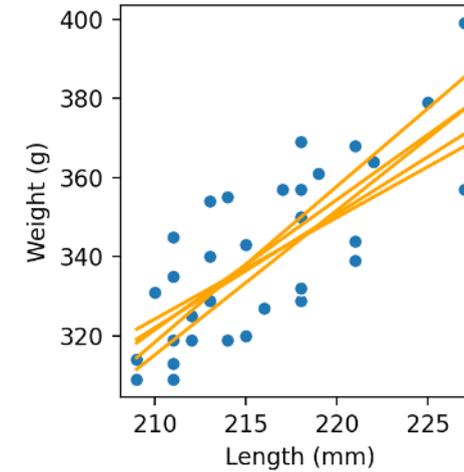


Peter Trimming, Wikimedia Commons, CC BY 2.0
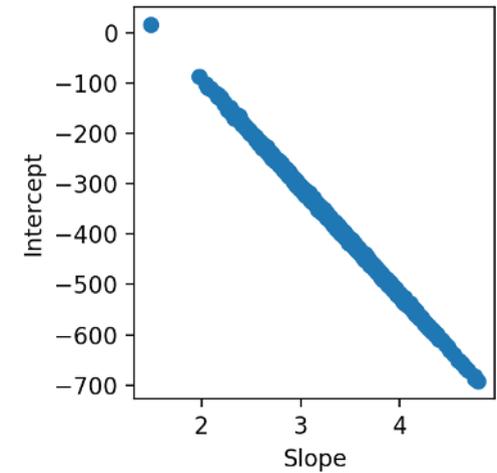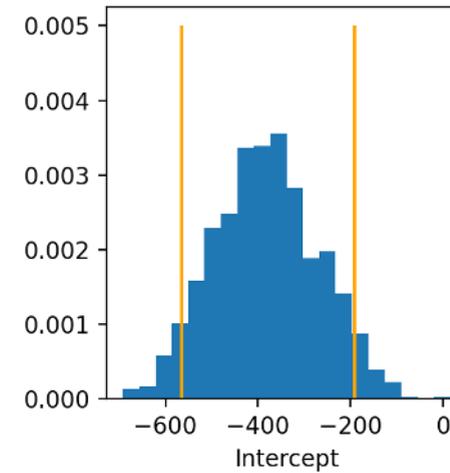


Data from Wauters and Dhondt 1989

# Parameter uncertainty

## Max likelihood

cf Bootstrap on data points + fitting

# Understanding more regression output

```
results = smf.ols('Weight ~ Length', data=datf).fit()
results.summary()
```

OLS Regression Results

| | | | |
|---|---|---|---|
| Dep. Variable: | Weight | R-squared: | 0.597 |
| Model: | OLS | Adj. R-squared: | 0.583 |
| Method: | Least Squares | F-statistic: | 44.37 |
| Date: | Sun, 10 Jan 2021 | Prob (F-statistic): | 2.24e-07 |
| Time: | 21:08:04 | Log-Likelihood: | -129.18 |
| No. Observations: | 32 | AIC: | 262.4 |
| Df Residuals: | 30 | BIC: | 265.3 |
| Df Model: | 1 | | |
| Covariance Type: | nonrobust | | |

*Max log likelihood* $\ln \hat{L}$

*Akaike Information criterion*

| | coef | std err | t | P>|t| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| Intercept | -382.7372 | 108.680 | -3.522 | 0.001 | -604.692 | -160.783 |
| Length | 3.3515 | 0.503 | 6.661 | 0.000 | 2.324 | 4.379 |

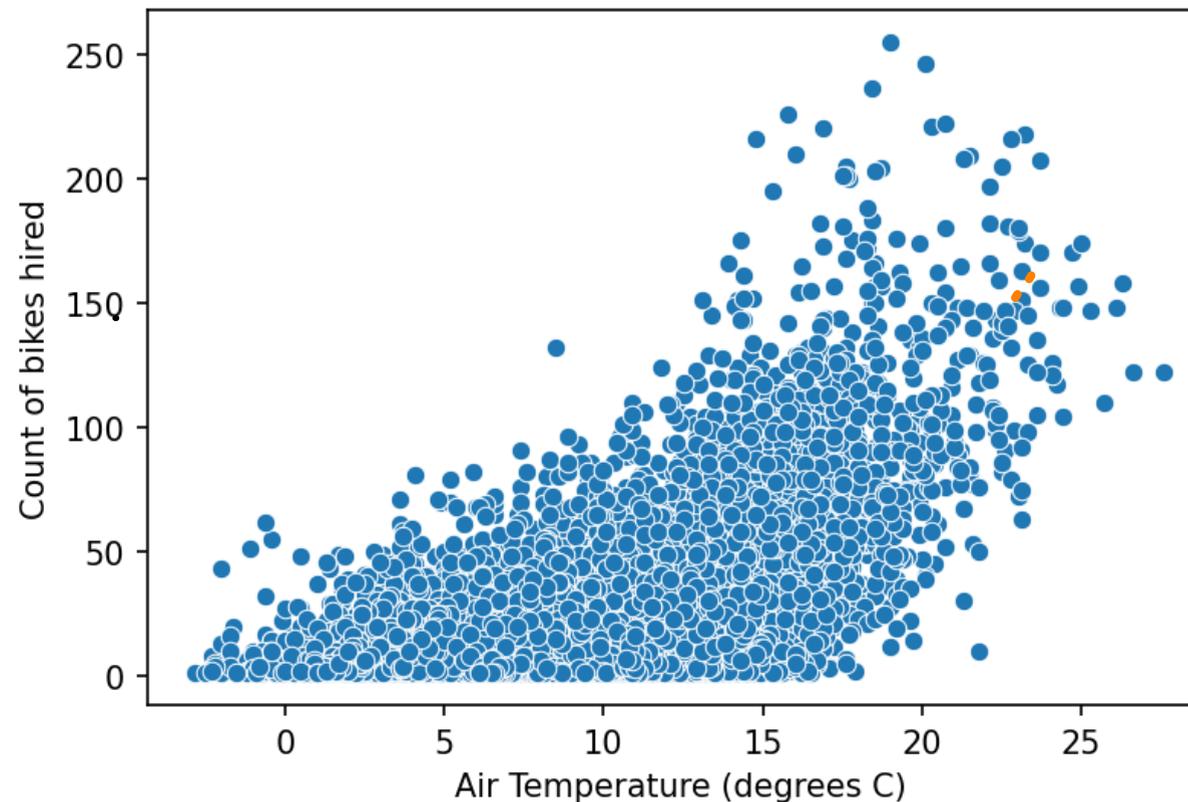| | | | |
|---|---|---|---|
| Omnibus: | 8.046 | Durbin-Watson: | 2.337 |
| Prob(Omnibus): | 0.018 | Jarque-Bera (JB): | 2.231 |
| Skew: | 0.092 | Prob(JB): | 0.328 |
| Kurtosis: | 1.720 | Cond. No. | 9.38e+03 |

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

Are there any techniques described in the course so far that could fit the data?

# Inf2 – Foundations of Data Science:
# Regression and inference –
# Max likelihood of univariate non-normal distributions
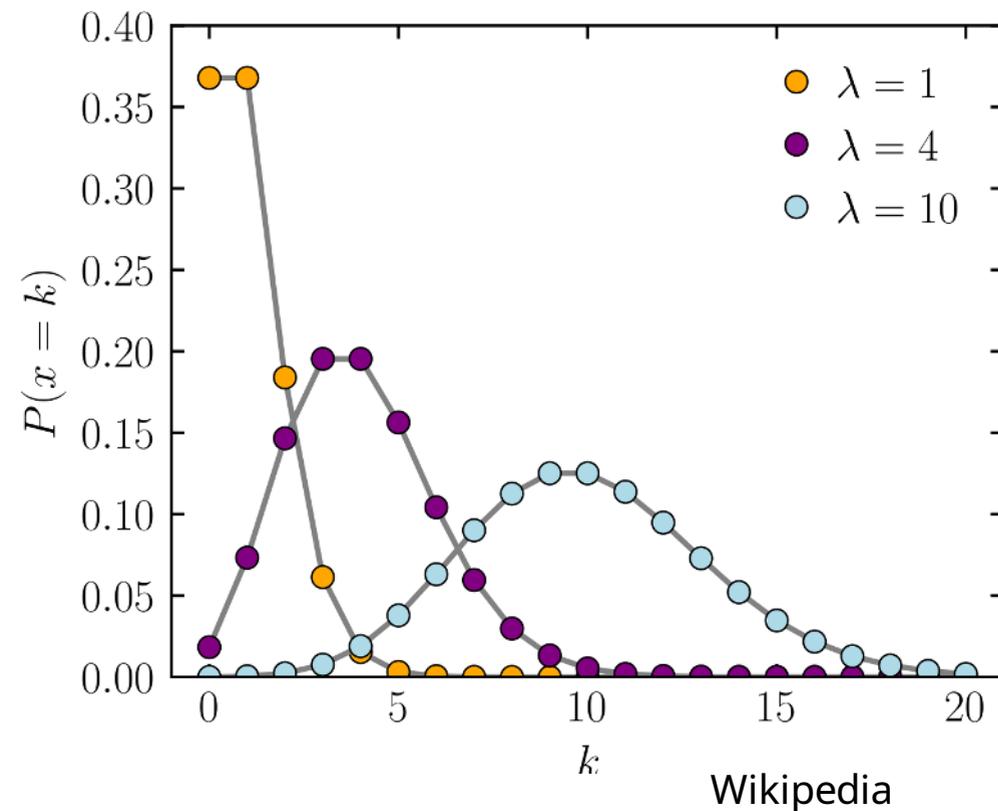
THE UNIVERSITY *of* EDINBURGH

**informatics**

FOUNDATIONS
OF
DATA
SCIENCE

# Max likelihood for models other than the normal

We don't have to assume the data is normally distributed.

E.g. Poisson distribution



Wikipedia

# E.g. Number of goals in World Cup football matches

Assumptions: Discrete events, uniform probability over time

Goals

0      90   )

t (minutes

Expected number of goals in a match $\lambda = 2.5$

# Number of deaths by horse kicks in the Prussian army


Wikpedia, CC-BY 2.0


Bortkewitsch 1898

$$\underline{y} = (0, 2, 2, 0, \cdots \qquad )$$
$$y_1 \quad y_2 \quad \cdots \qquad \cdots \qquad y_{280}$$

$$n_k = \sum_{i=1}^{n} I(y_i = k)$$

Can we infer the parameter from the data?

| $k$ | $n_k$ |
|---|---|
| 0 | 144 |
| 1 | 91 |
| 2 | 32 |
| 3 | 11 |
| 4 | 2 |

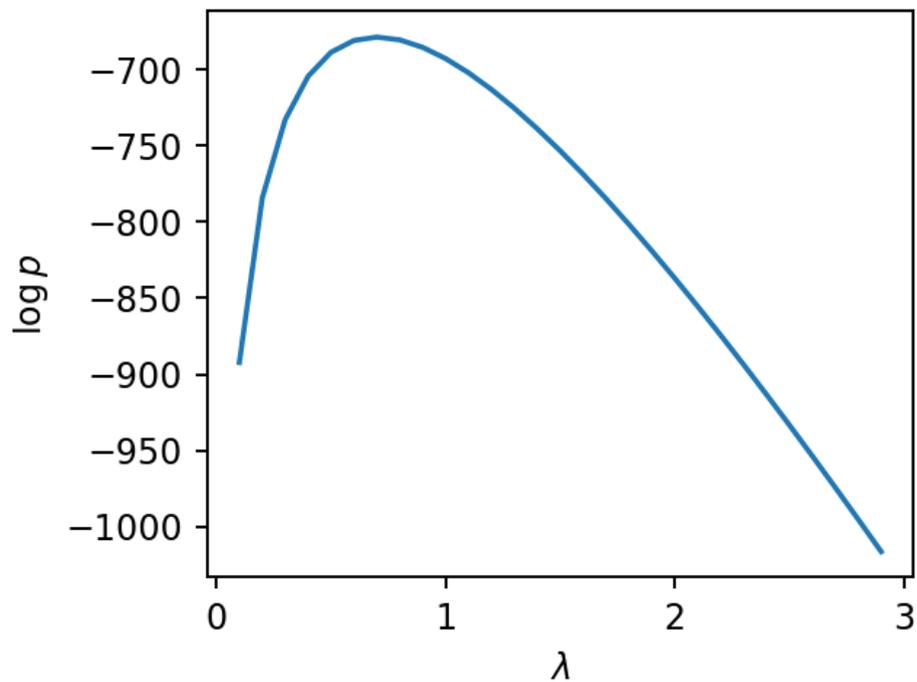# Log likelihood calculation of Poisson distribution

$$\ell = \ln \, P(Y = y_1, \ldots, y_n) = \ln \lambda \sum_{i=1}^{n} y_i - n\lambda - \sum_{i=1}^{n} \ln y_i!$$

$$\frac{d\ell}{d\lambda} = 0$$

$$\vdots$$



$$\Rightarrow \lambda_{MLE} = \frac{1}{n} \sum_{i=1}^{n} y_i$$

# Inf2 – Foundations of Data Science:
# Regression and inference –
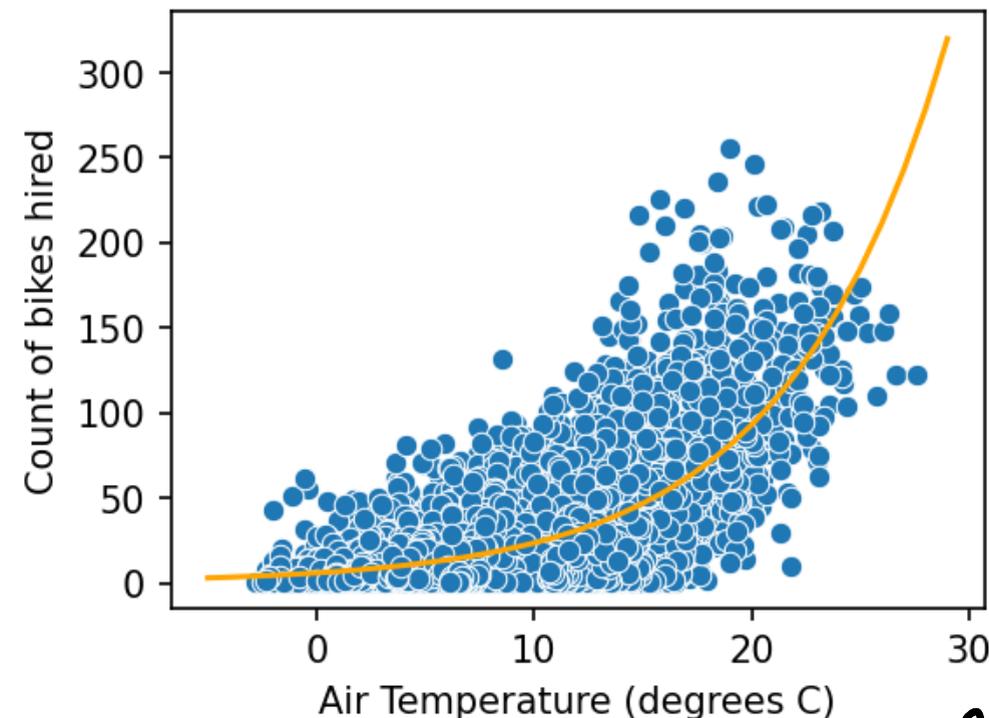# Poisson regression

# Poisson regression – generative model

# Results with statsmodels GLM



Generalized Linear Model Regression Results

| Dep. Variable: | count | No. Observations: | 8301 |
|---|---|---|---|
| Model: | GLM | Df Residuals: | 8299 |
| Model Family: | Poisson | Df Model: | 1 |
| Link Function: | Log | Scale: | 1.0000 |
| Method: | IRLS | Log-Likelihood: | -84533. |
| Date: | Wed, 01 Mar 2023 | Deviance: | 1.3111e+05 |
| Time: | 06:46:41 | Pearson chi2: | 1.40e+05 |
| No. Iterations: | 5 | Pseudo R-squ. (CS): | 1.000 |
| Covariance Type: | nonrobust | | |

| | coef | std err | z | P>|z| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| const | 1.7861 | 0.006 | 304.092 | 0.000 | 1.775 | 1.798 |
| air_temperature | 0.1373 | 0.000 | 323.057 | 0.000 | 0.136 | 0.138 |

$\beta_0$

$\beta_1$

$$\ln \lambda = \beta_0 + \beta_1 x$$

$$\lambda = e^{\beta_0 + \beta_1 x} = e^{\beta_0} e^{\beta_1 x} \qquad e^{0.1373} = 1.14$$

# Poisson regression

$$l = \ln P(Y = y_1, \ldots, y_n)$$

$$l(\beta_0, \beta_1) = \sum_{i=1}^{n}(\beta_0 + \beta_1 x_i) y_i \;-\; \sum_{i=1}^{n} e^{\beta_0 + \beta_1 x_i} \;-\; \sum_{i=1}^{n} \ln y_i!$$

To my Valentine, Poisson Regression

Roses are red

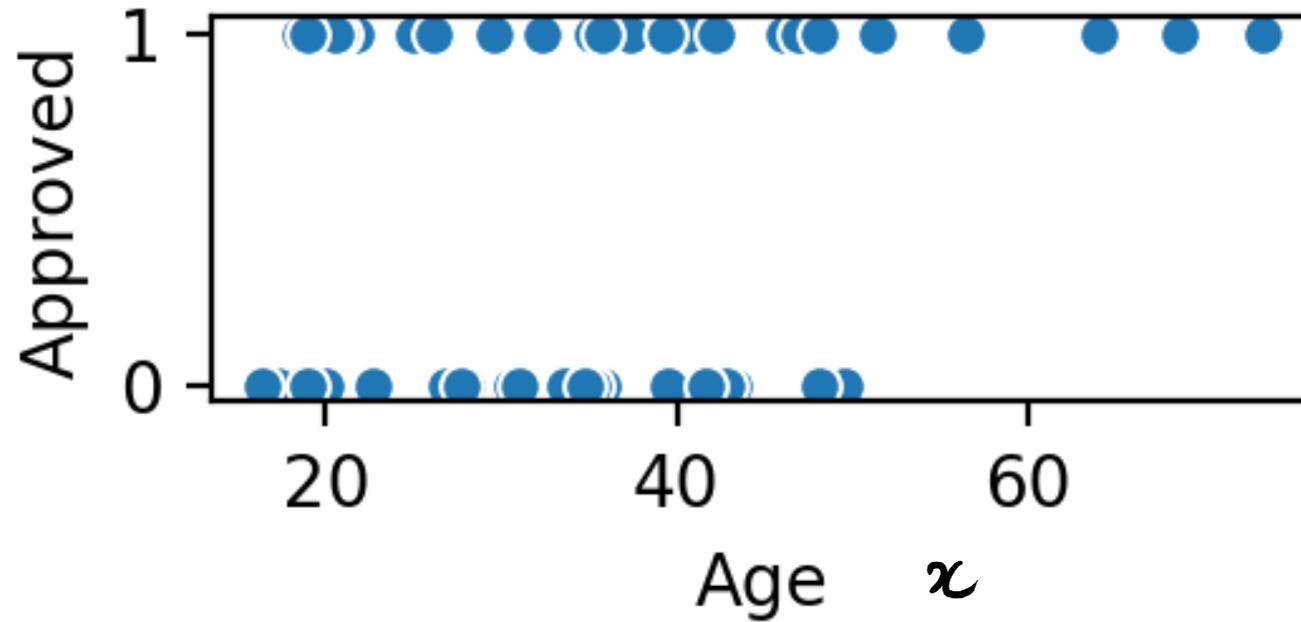Violets are blue

Some things aren't normal

and nor are you

# Inf2 – Foundations of Data Science: Regression and inference – Logistic regression and generalised linear models

THE UNIVERSITY *of* EDINBURGH

**informatics**

**FOUNDATIONS OF DATA SCIENCE**

# Excercise



What distribution would we use to model the data here?

How would the parameter of that distribution depend on x (Age)?

# Generalised linear models (GLMs)

|  | Distribution | Link function |
|---|---|---|
| linear regression | Normal | $\mu = \beta_0 + \beta_1 x \quad , \sigma^2$ |
| Poisson regression | Poisson | $\ln \lambda = \beta_0 + \beta_1 x$ |
| logistic regression | Bernoulli | $\ln \dfrac{p}{1-p} = \beta_0 + \beta_1 x$ |

# Link functions

Expected value $\mu = E(Y|x)$ of a Bernoulli dist is $p$

"        "        $\mu = E(Y|x)$ " " Poisson dist is $\lambda$

In general the link function is denoted $g(\mu)$
where $\mu = E(Y|x)$ for that distribution.

$$g(\mu) = \beta_0 + \beta_1 x$$

To make predictions, we invert the link function:

$$\mu = g^{-1}(\beta_0 + \beta_1 x)$$

# Inf2 – Foundations of Data Science: Regression and inference – And finally…

# Max likelihood -> Bayesian Inference

Bayes Theorem:

$$\underbrace{P(\vartheta \mid \underline{Y} = y)}_{\text{Posterior}} = \frac{\overbrace{P(\underline{Y} = y \mid \vartheta)}^{\text{Likelihood}} \; \overbrace{p(\vartheta)}^{\text{Prior}}}{\underbrace{P(Y = y)}_{\text{Evidence}}}$$

Horsekick posterior



$$P(Y = y) = \int_{-\infty}^{\infty} P(Y = y \mid \vartheta) \, p(\vartheta) \, d\vartheta$$

# Summary

Motivated the probabalistic basis of inference using
max likelihood                         .

Important: think of what distribution should
describe the data

Links to future courses:
– MLG (derivation of standard ML methods)
– ATML (new in 25–26: cutting–edge machine learning)
– MCI (Causal inference)