

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



THE UNIVERSITY *of* EDINBURGH  
**informatics**

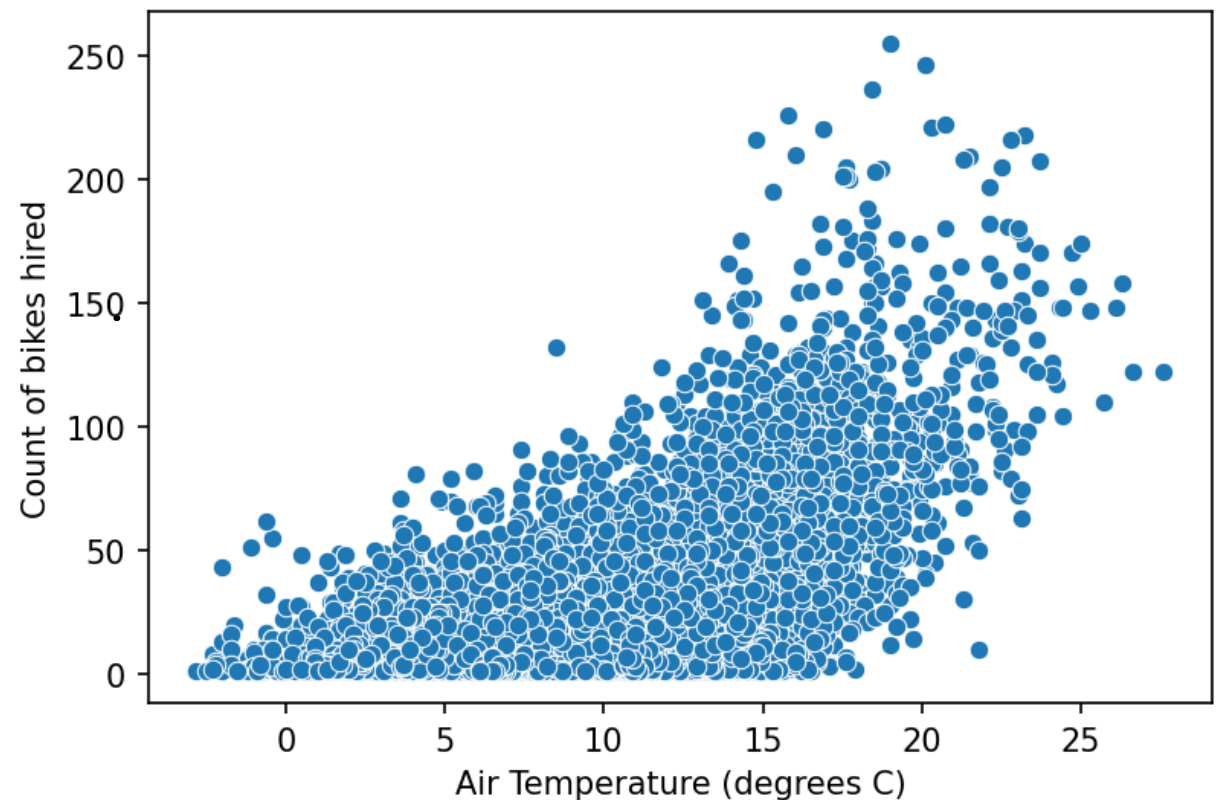
**FOUNDATIONS**  
**OF**  
**DATA**  
**SCIENCE**

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 + prediction uncertainty
1. Max likelihood with non-normal distributions
2. Poisson regression
3. Logistic regression and generalised linear models



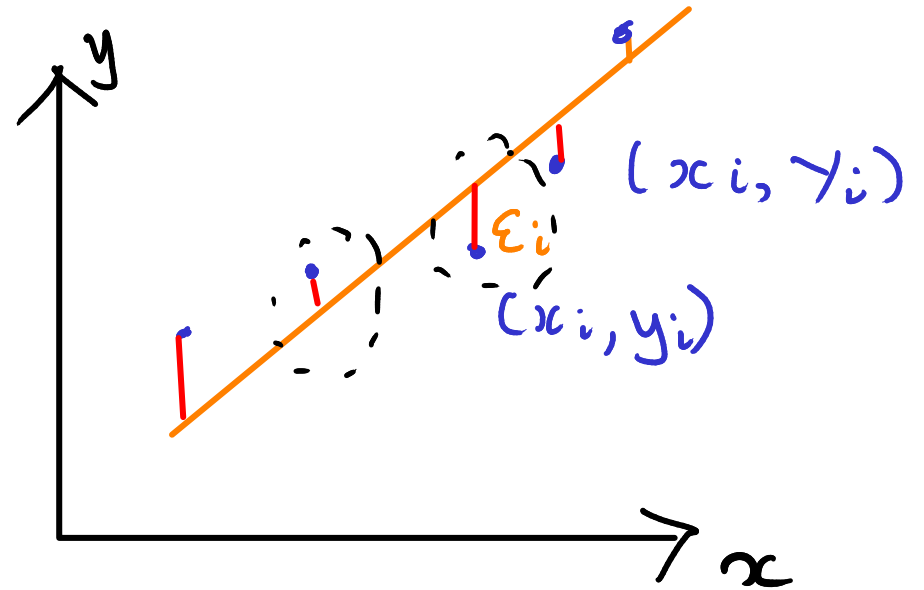
**Foundations of Data Science:  
Regression and inference -  
Recap of max likelihood applied to linear  
regression**

# 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

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

$$\ln p(\underline{y} = y_1, \dots, y_n; x_1, \dots, x_n \mid \underbrace{\beta_0, \beta_1, \sigma^2}_{\mu})$$

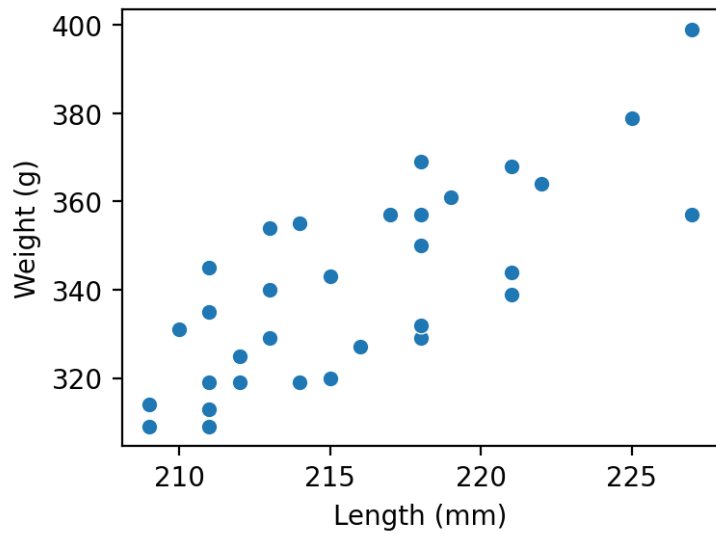
$$= \sum_{i=1}^n \left( -\frac{1}{2} \ln \pi \sigma^2 - \frac{1}{2} \left( \frac{y_i - \beta_0 - \beta_1 x_i}{\sigma} \right)^2 \right)$$



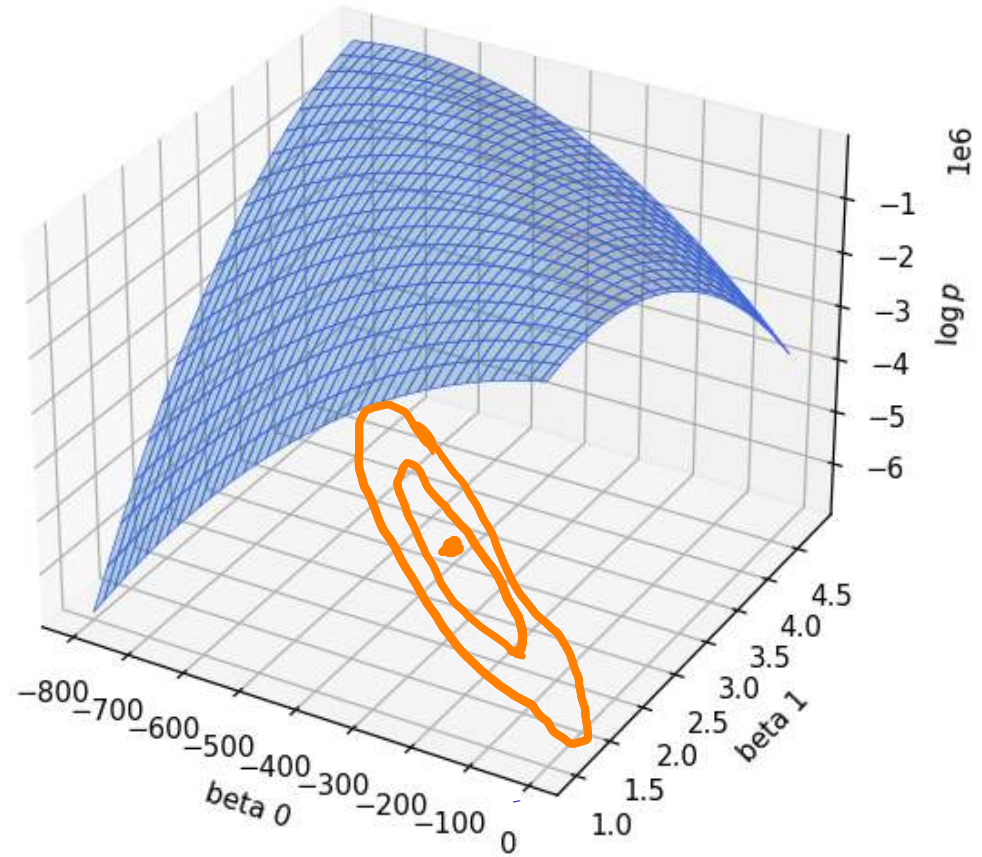
# Log likelihood of coefficients



Peter Trimming, Wikimedia Commons, CC BY 2.0



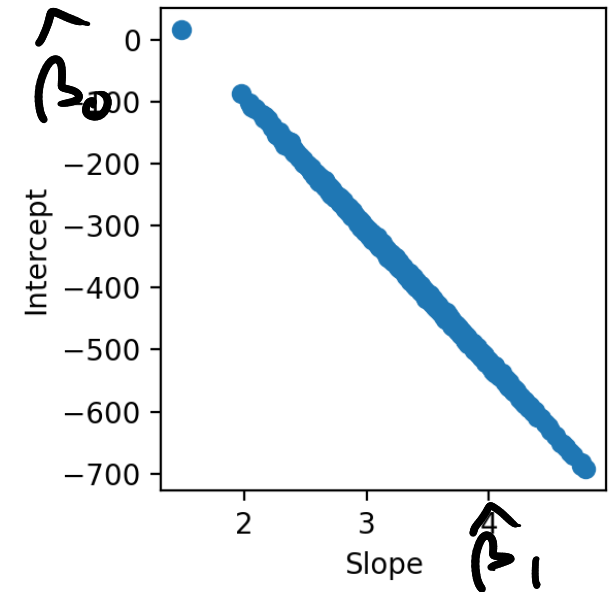
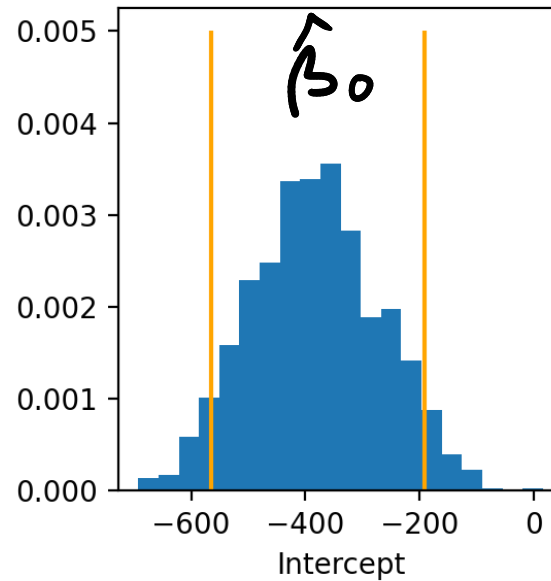
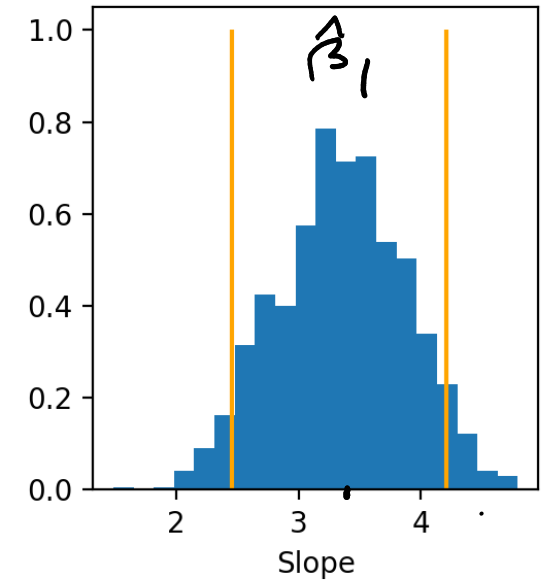
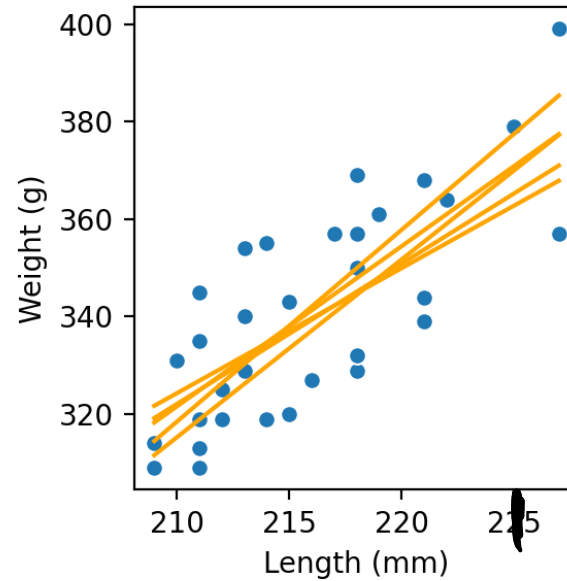
Data from Wauters and Dhondt 1989



# Bootstrap inference of coefficients

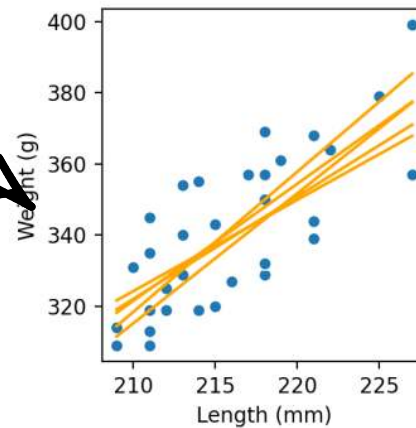
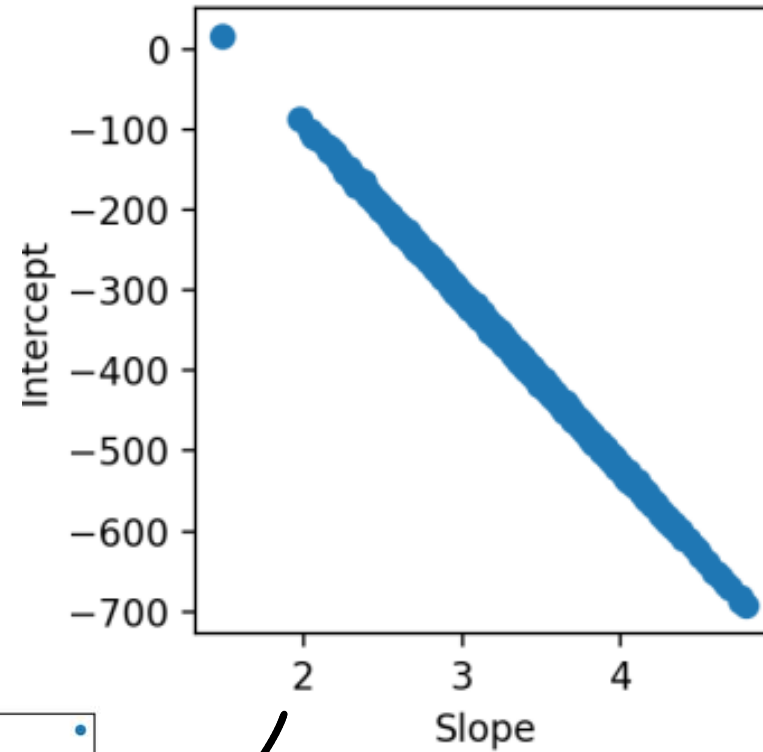
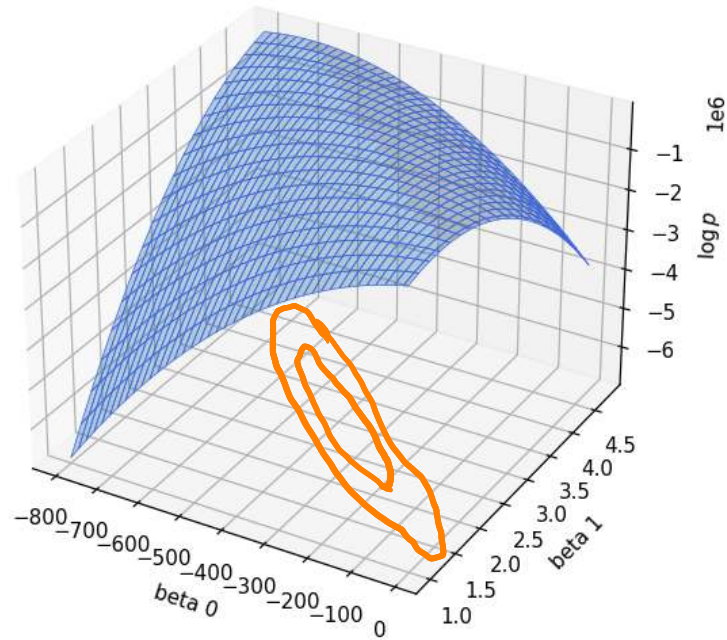


Peter Trimming, Wikimedia Commons, CC BY 2.0



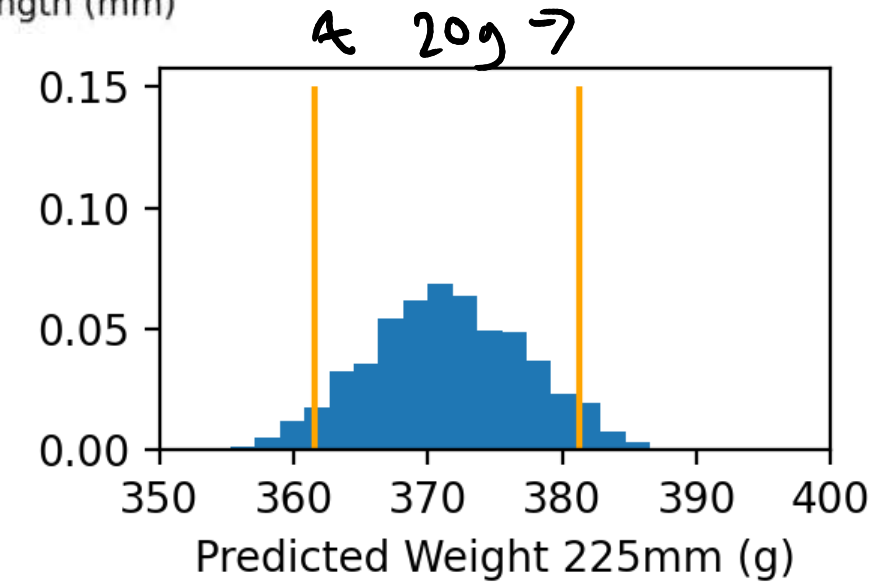
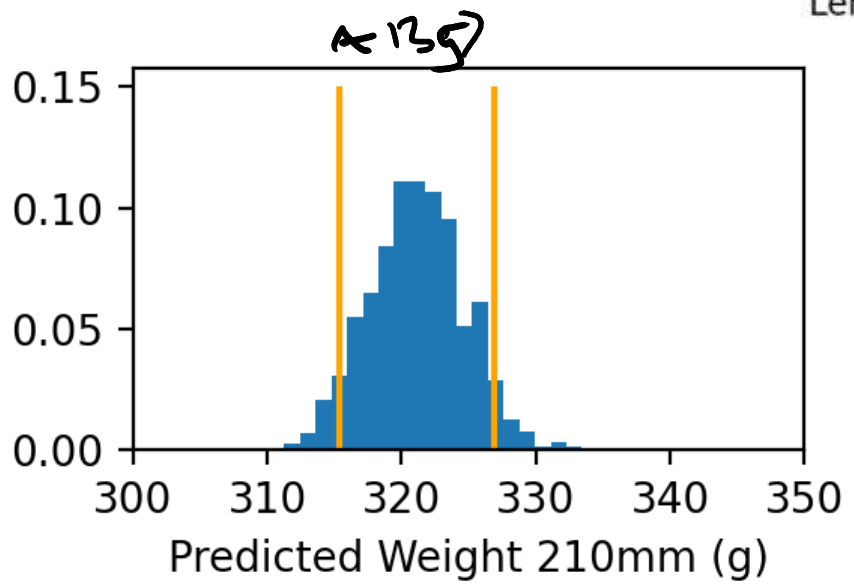
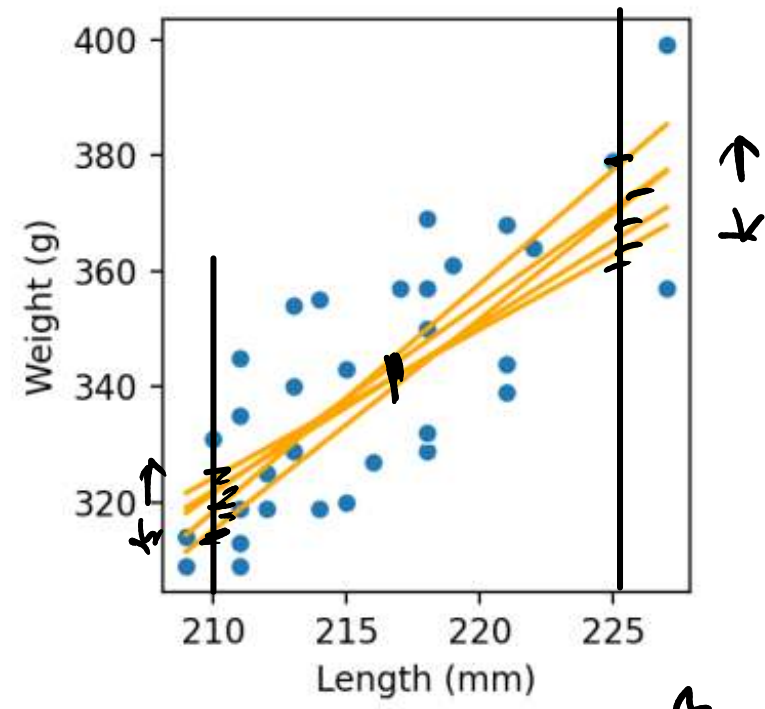
# (Log) Likelihood function

# Bootstrap samples





# Uncertainty in predictions (with Bootstrap)



$$y = \beta_0^{(1)} + \beta_1^{(1)} \cdot 210$$

$$y = \beta_0^{(2)} + \beta_1^{(2)} \cdot 210$$

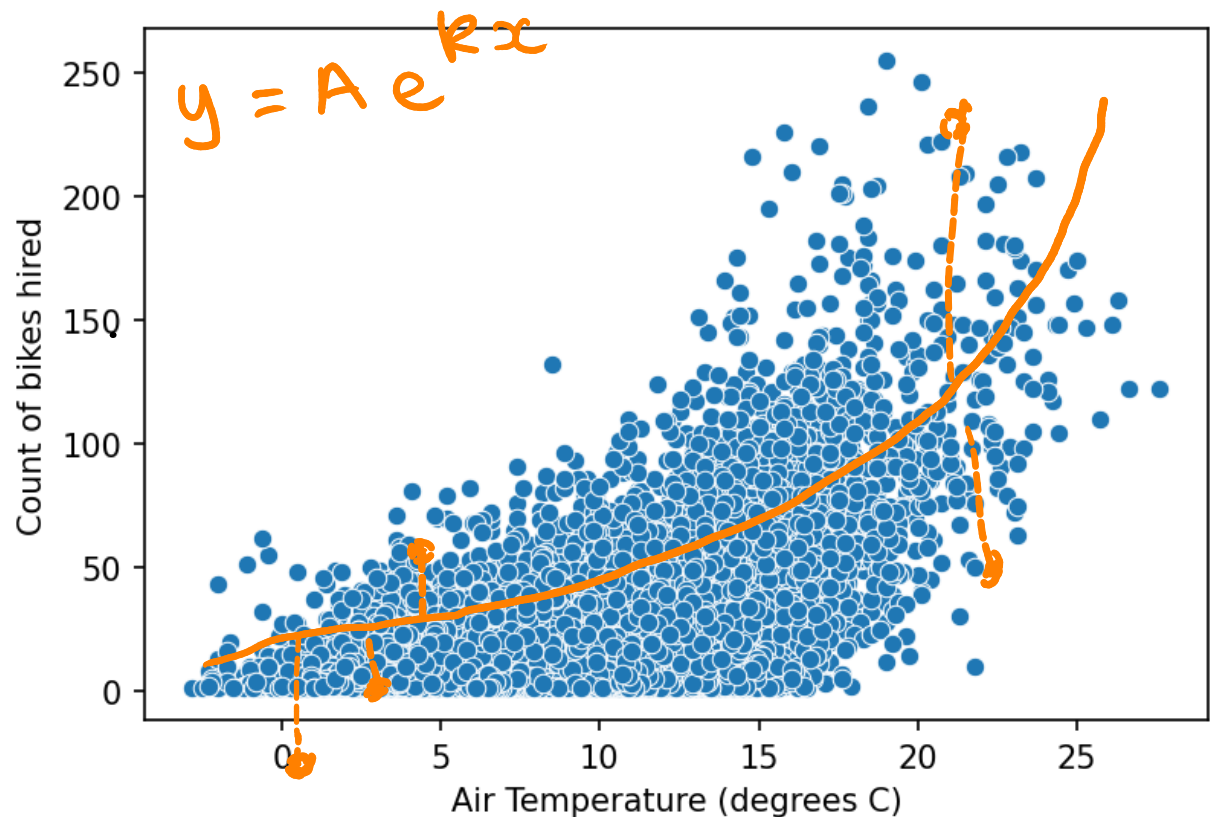
$$\vdots$$

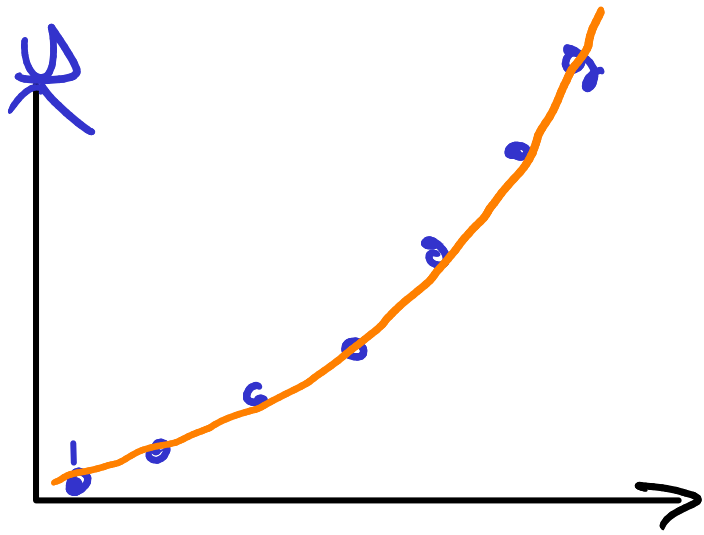
Bootstrap sample #1  
 " " " #2

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?

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

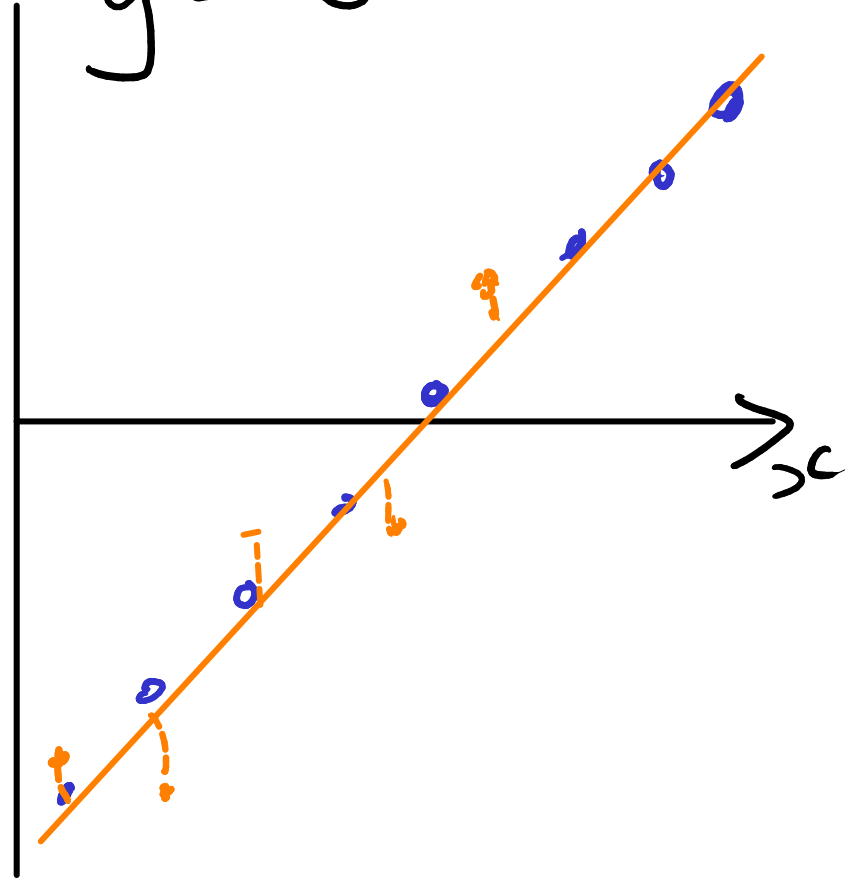




$x$   $\rightarrow$

$\ln y$

$$\ln y = \beta_0 + \beta_1 x$$
$$y = e^{\beta_0 + \beta_1 x}$$



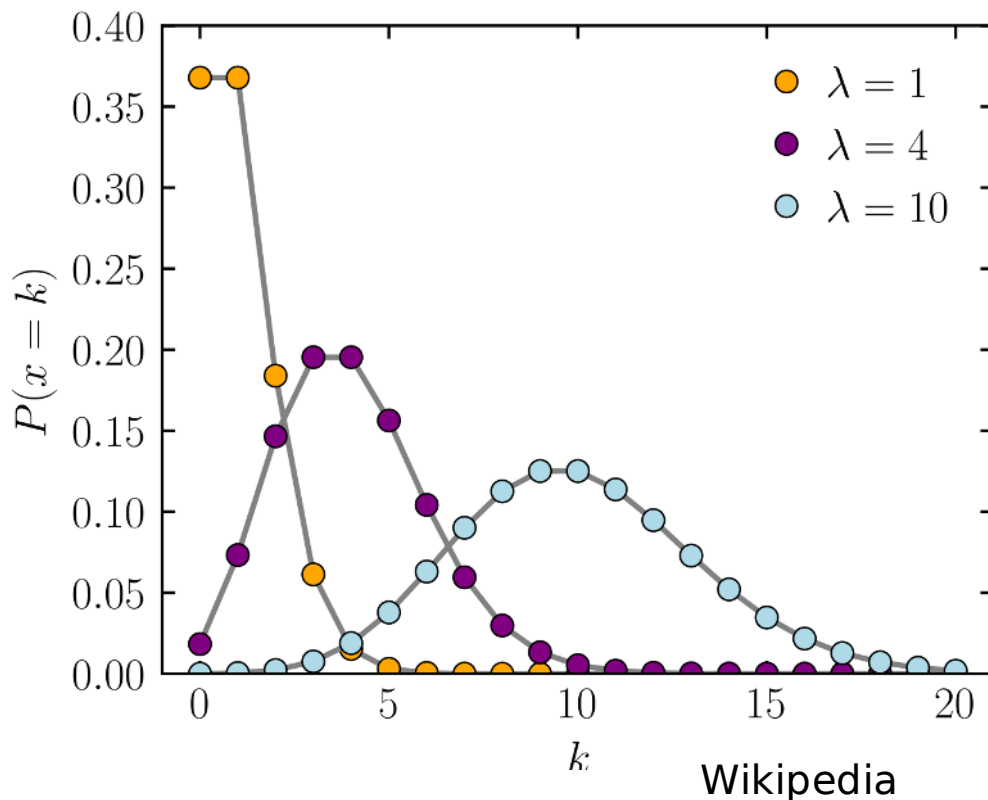


**Foundations of Data Science:  
Regression and inference -  
Max likelihood of univariate non-normal  
distributions**

# Max likelihood for models other than the normal

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

E.g. Poisson distribution



$$P(Y = k) = \frac{\lambda^k e^{-\lambda}}{k!}$$

$$k = 0, 1, 2, \dots$$

$$E(Y) = \lambda = \mu$$

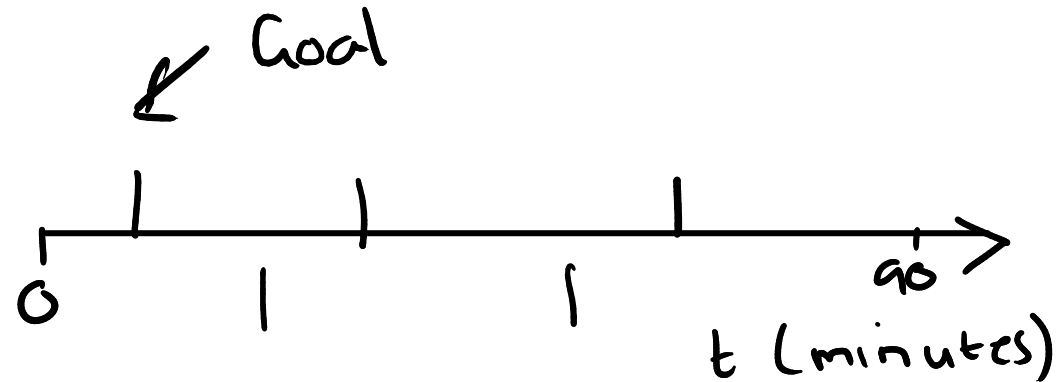
$$V(Y) = \lambda = \sigma^2$$



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



Wikipedia, CC-BY-SA 3.0



|||||

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

$$P(Y = k) = \frac{2.5^k e^{-2.5}}{k!}$$

$$P(Y = 0) = \frac{2.5^0 e^{-2.5}}{0!} = e^{-2.5} = 0.082$$

$$P(Y = 1) = 0.205$$

$$P(Y = 2) = 0.257$$

# Number of deaths by horse kicks in the Prussian army



Wikipedia, CC-BY 2.0

	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94
G	—	2	2	1	—	—	1	1	—	3	—	2	1	—	—	1	—	1	—	1
I	—	—	—	2	—	3	—	2	—	—	—	1	1	1	—	2	—	3	1	—
II	—	—	—	2	—	2	—	—	1	1	—	—	2	1	1	—	—	2	—	—
III	—	—	—	1	1	1	2	—	2	—	—	—	1	—	1	2	1	—	—	—
IV	—	1	—	1	1	1	1	—	—	—	—	1	—	—	—	—	1	1	—	—
V	—	—	—	—	2	1	—	—	1	—	—	1	—	1	1	1	1	1	1	—
VI	—	—	1	—	2	—	—	1	2	—	1	1	3	1	1	1	—	3	—	—
VII	1	—	1	—	—	—	1	—	1	1	—	—	2	—	—	2	1	—	2	—
VIII	1	—	—	—	1	—	—	1	—	—	—	—	1	—	—	—	1	1	—	1
IX	—	—	—	—	—	2	1	1	1	—	2	1	1	—	1	2	—	1	—	—
X	—	—	1	1	—	1	—	2	—	2	—	—	—	—	2	1	3	—	1	1
XI	—	—	—	—	2	4	—	1	3	—	1	1	1	1	2	1	3	1	3	1
XIV	1	1	2	1	1	3	—	4	—	1	—	3	2	1	—	2	1	1	—	—
XV	—	1	—	—	—	—	—	1	—	1	1	—	—	—	2	2	—	—	—	—

Bortkewitsch 1898

$$\underline{y} = (y_1, y_2, \dots, y_{280})$$

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

$$n_k = \sum_{i=1}^{280} I(y_i = k)$$

# Log likelihood calculation of Poisson distribution

$$\text{Log likelihood } l = \ln P(Y = y_1, \dots, y_n | \lambda)$$

$$= \sum_{i=1}^n \ln P(Y = y_i)$$

$$= \sum_{i=1}^n \ln \frac{\lambda^{y_i} e^{-\lambda}}{y_i!}$$

$$= \sum_{i=1}^n (y_i \ln \lambda + (-\lambda) - \ln y_i!)$$

$$l(\lambda) = \ln \lambda \sum_{i=1}^n y_i - n\lambda - \sum_{i=1}^n \ln y_i!$$

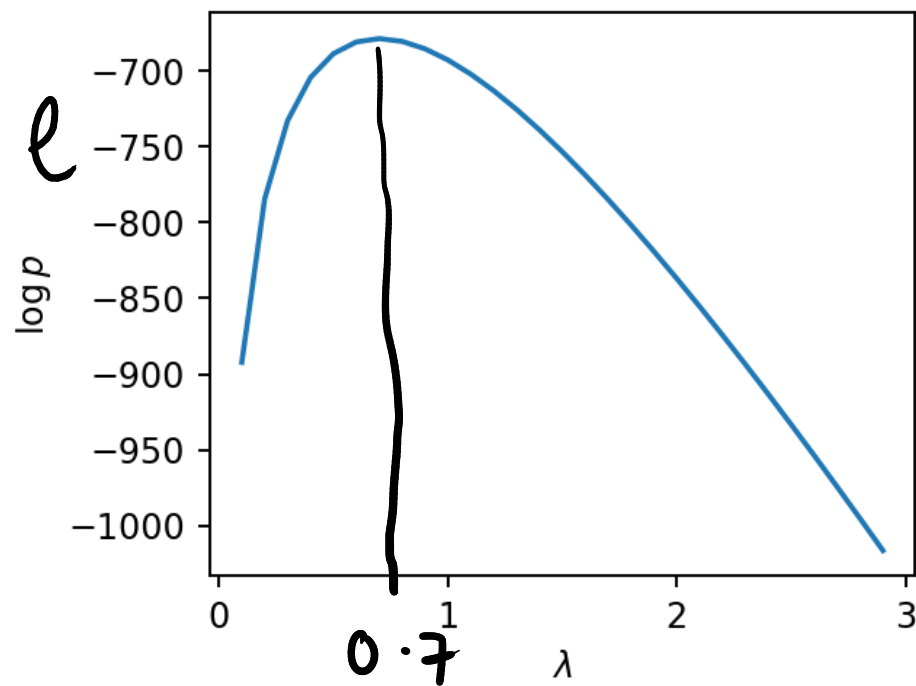
---

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

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

⋮  
⋮  
⋮  
⋮  
⋮

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

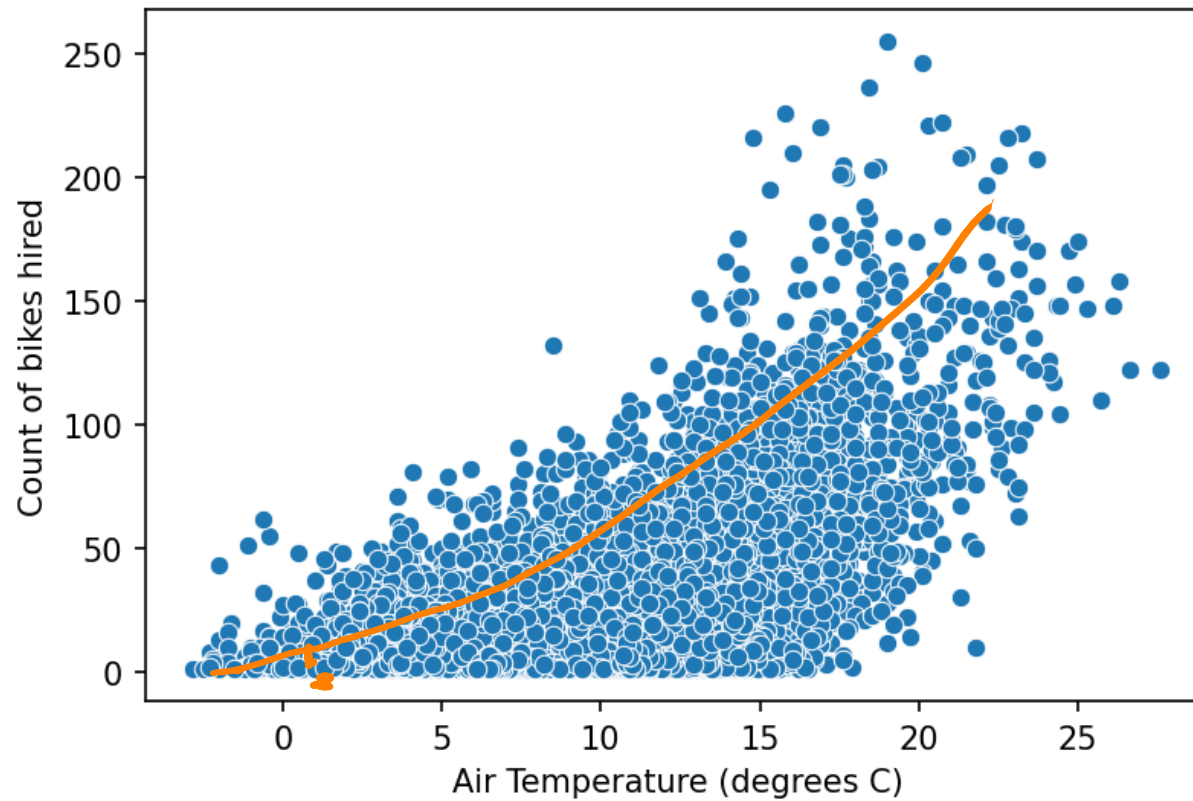




**Foundations of Data Science:  
Regression and inference -  
Poisson regression**



# Poisson regression



$$Y_i \sim \text{Poisson} \left( e^{\beta_0 + \beta_1 x_i} \right)$$

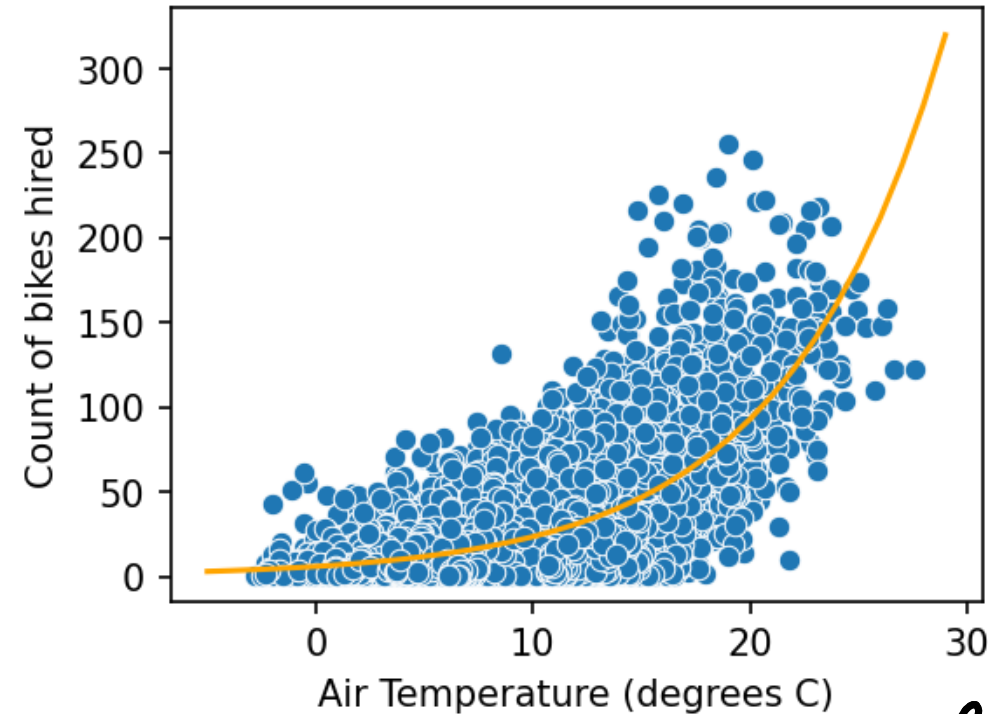
$$\ln \lambda = \beta_0 + \beta_1 x$$

# Results with statsmodels GLM

## Generalized Linear Model Regression Results

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



$\beta_0$   
 $\beta_1$

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

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

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

# Poisson regression

$$l = \ln P(\underline{Y} = y_1, \dots, y_n)$$

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

→ optimise  $\beta_0$  &  $\beta_1$

To my Valentine, Poisson Regression


Roses are red



Violets are blue

Some things aren't normal

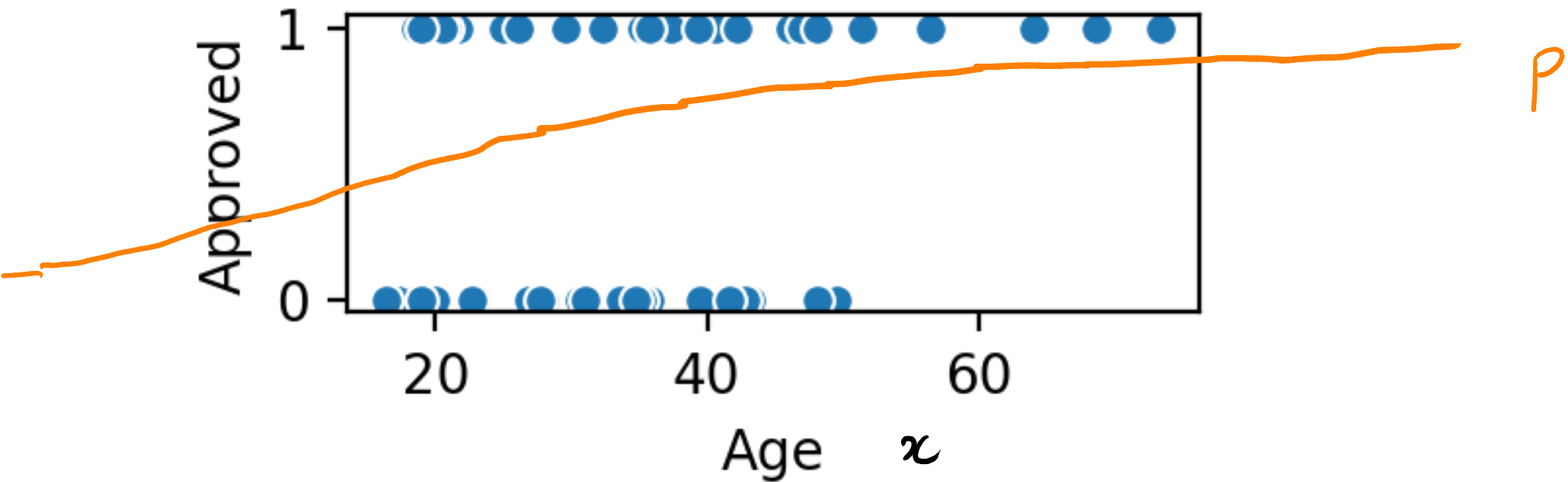
and nor are you



**Foundations of Data Science:  
Regression and inference -  
Logistic regression and  
generalised linear models**



# Excercise



What distribution would we use to model the data here?

Bernoulli - Parameter  $p$

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

Logistic function

$$p = \text{Logistic}(\beta_0 + \beta_1 x)$$

# Generalised linear models (GLMs)

	<u>Distribution</u>	<u>Link function</u>
linear regression	Normal	$\mu = \beta_0 + \beta_1 x, \sigma^2$
Poisson regression	Poisson	$\ln \lambda = \beta_0 + \beta_1 x$
logistic regression	Bernoulli	$\ln \frac{p}{1-p} = \beta_0 + \beta_1 x$

# 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:

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



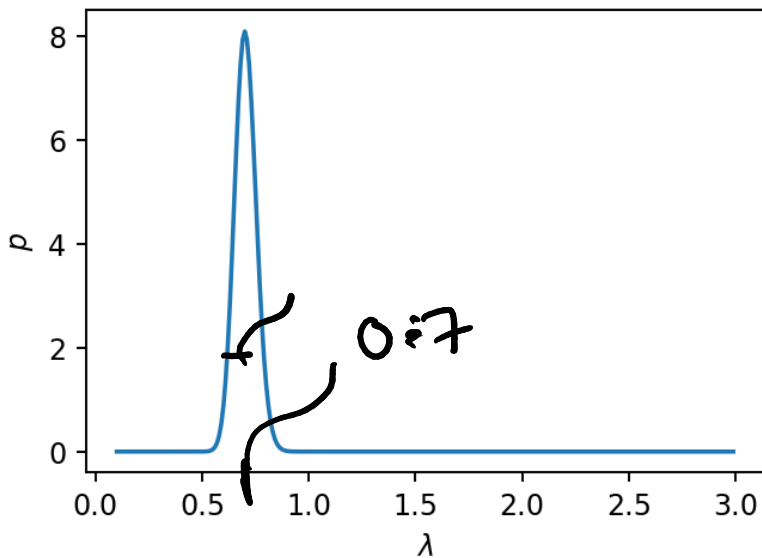
**Foundations of Data Science:  
Regression and inference -  
And finally...**

# Max likelihood -> Bayesian Inference

Bayes Theorem:

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

Horsekick posterior



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

# Summary

Motivated the probabilistic 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)
- MLPR (Bayesian approach; application to new problems)
- MCI (Causal inference)