§2.4 Least squares estimation & probability



Least squares estimation

Given observed data y_1, \dots, y_n , fit the linear model

$$y \approx \beta_1 e_1 + \cdots + \beta_K e_K$$

i.e.

$$y_i \approx \beta_1 e_{1,i} + \dots + \beta_K e_{K,i}$$

by choosing the parameters β_1, \dots, β_K so as to minimize the mean square error

$$mse = \frac{1}{n} \sum_{i=1}^{n} (y_i - pred_i)^2$$

where $\operatorname{pred}_i = \beta_1 e_{1,i} + \cdots + \beta_K e_{K,i}$

Example 2.1.1

The Iris dataset has 50 records of iris measurements, from three species.

How does Petal.Length (PL) depend on Sepal.Length (SL)?

We fitted the linear model

$$PL \approx \alpha + \beta SL + \gamma SL^2$$

Maximum likelihood estimation

Given observed data y_1, \dots, y_n , fit the probability model $Y_i \sim \cdots$

by choosing the model parameters so as to maximize the log likelihood of the observed data

$$\log \Pr(y_1, \dots, y_n) = \sum_{i=1}^n \log \Pr_Y(y_i; \dots)$$

Example

Let's fit the probability model

$$PL_i \sim \alpha + \beta SL_i + \gamma SL_i^2 + Normal(0, \sigma^2)$$

Model for a single observation:

$$PL_{i} \sim \alpha + \beta SL_{i} + \gamma SL_{i}^{2} + N(0, \sigma^{2})$$

$$y_{i} \sim \alpha + \beta e_{i} + \gamma f_{i} + N(0, \sigma^{2})$$

$$\sim N(\alpha + \beta e_{i} + \gamma f_{i}, \sigma^{2})$$

Likelihood of a single observation:

Pry:
$$(y_i : \alpha, \beta, \gamma, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y_i^2 - (\alpha + \beta e_i + \gamma + \gamma_i^2))^2/2\sigma^2}$$

Log likelihood of the dataset:

leg Pr
$$(y_1, \dots, y_n)$$
 $\propto (\beta, \gamma, \sigma) = -\frac{n}{2}\log(2\pi \sigma^2) - \frac{1}{2\sigma^2}\sum (y_1 - (\kappa + \beta r_1 + \delta f_1))^2$

We then seek $\kappa_1\beta_1\beta_2$, σ so maximize this.

Maximize over the unknown parameters,

 α , β , γ , and σ :

and
$$\sigma$$
:

were

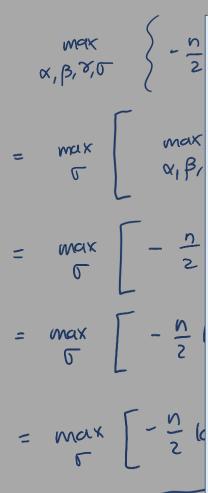
 $\alpha_{1}\beta_{1}, \overline{\gamma}_{1}\overline{\zeta}$
 $\begin{cases}
-\frac{n}{2} (\log(2\pi\sigma^{2})) - \frac{1}{2\sigma^{2}} \sum_{i=1}^{n} (y_{i} - (\alpha + \beta e_{i} + \overline{\gamma}_{1}e_{i}))^{2} \\
= \max_{\alpha_{1}\beta_{1}, \overline{\gamma}_{1}\overline{\zeta}}
\end{cases}$
 $= \max_{\alpha_{1}\beta_{1}, \overline{\gamma}_{2}}
\begin{cases}
-\frac{n}{2} \log(2\pi\sigma^{2}) - \frac{1}{2\sigma^{2}} \sum_{i} (y_{i} - (\alpha + \beta e_{i} + \overline{\gamma}_{1}e_{i}))^{2} \\
= \max_{\alpha_{1}\beta_{1}, \overline{\gamma}_{2}}
\end{cases}$
 $= \max_{\alpha_{1}\beta_{1}, \overline{\gamma}_{2}}
\begin{cases}
-\frac{n}{2} \log(2\pi\sigma^{2}) + \max_{\alpha_{1}\beta_{1}, \overline{\gamma}_{2}} \\
-\frac{1}{2\sigma^{2}} \sum_{i} (y_{i} - (\alpha + \beta e_{i} + \overline{\gamma}_{1}e_{i}))^{2} \\
= \max_{\alpha_{1}\beta_{2}, \overline{\gamma}_{2}}
\end{cases}$
 $= \max_{\alpha_{1}\beta_{2}}
\begin{cases}
-\frac{n}{2} \log(2\pi\sigma^{2}) - \frac{1}{2\sigma^{2}} \sum_{i} (y_{i} - \hat{y}_{i})^{2} \\
= \max_{\alpha_{1}\beta_{2}, \overline{\gamma}_{2}}
\end{cases}$

where

 $\hat{y}_{i} = \hat{\alpha} + \hat{\beta}e_{i} + \hat{\delta}f_{i}$

obtained by least squares estimation

 α , β , γ , and σ :



Least squares estimation derives from a Gaussian probability model.

If that model doesn't fit the data, then don't use least squares estimation!

Sign in



Stop and search

• This article is more than 3 years old

Met police 'disproportionately' use stop and search powers on black people

London's minority black population targeted more than white population in 2018 - official figures



Let $y_i \in \{0,1\}$ be the outcome for stop-and-search incident i.



$$y_i \approx \alpha + \beta_{\mathrm{eth}_i}$$
 i.e. $Y_i \sim \alpha + \beta_{\mathrm{eth}_i} + N(0, \sigma^2)$

Fit α and β_{Bl} , β_{Mi} , ... using least squares estimation or, equivalently, fit using maximum likelihood estimation

 $Y_i \sim \text{Bin}(1, \alpha + \beta_{\text{eth}_i})$ Fit the parameters using maximum likelihood estimation

There's a more advanced version called *Logistic Regression*, for $Bin(1, \theta_i)$ where θ_i depends on multiple features. It uses softmax. See the code in [stop-and-search.ipynb], or Part II Advanced Data Science.

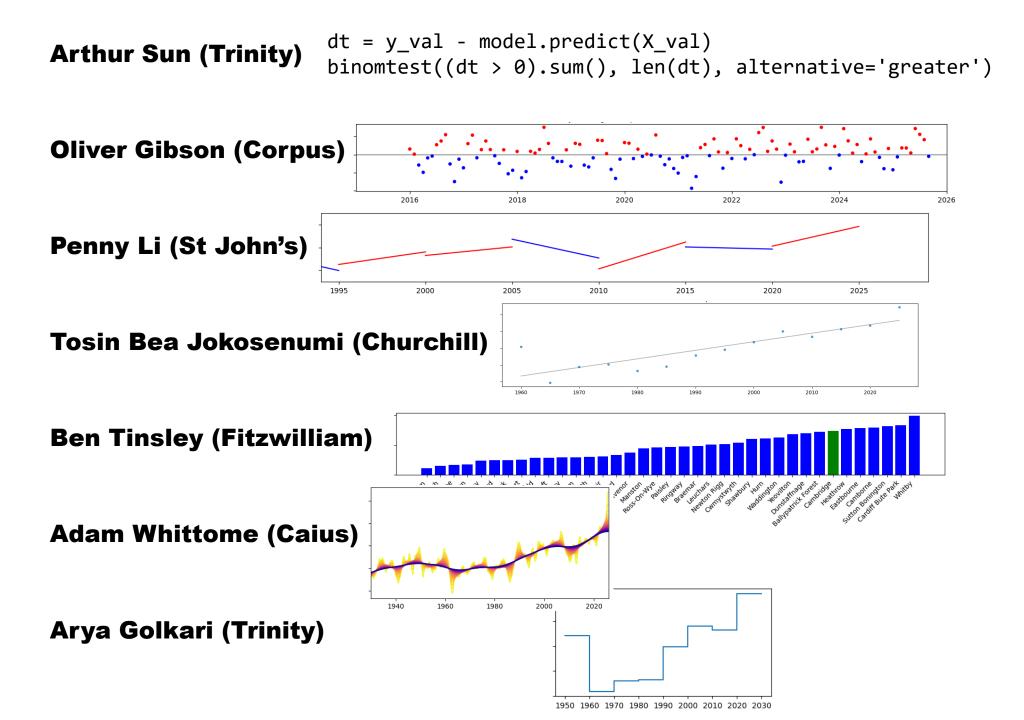


Climate dataset challenge

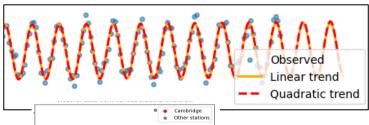
- What is the rate of temperature increase in Cambridge?
- Are temperatures increasing at a constant rate, or has the increase accelerated?
- How do results compare across the whole of the UK?

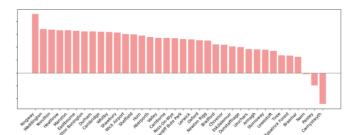
Your task is to answer these questions using appropriate linear models, and to produce elegant plots to communicate your findings. Please submit a Jupyter notebook, or a pdf. Include explanations of what your models are, and of what your plots show.

The dataset is from https://www.metoffice.gov.uk/pub/data/weather/uk/climate/. Code for retrieving the dataset is given at the bottom.



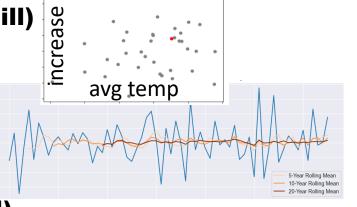
Ed Read (Churchill)





Justin Leung (Churchill)

Zerui Chen (Wolfson)



Isaac Chan (Churchill)

$$\gamma_{1990} = -1.019$$

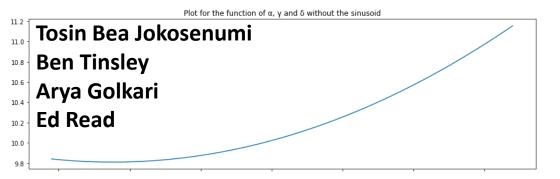
$$\gamma_{1990} = -1.019, \qquad \gamma_{20000} = -0.649, \qquad \gamma_{2010} = -0.737, \qquad \gamma_{2020} = 0$$

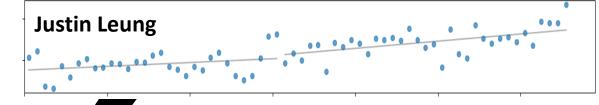
$$\gamma_{2010} = -0.737$$

$$\gamma_{2020} = 0$$

This model assumes a simple functional form, which makes it easy to summarize the nonlinearity – but which may be wrong.

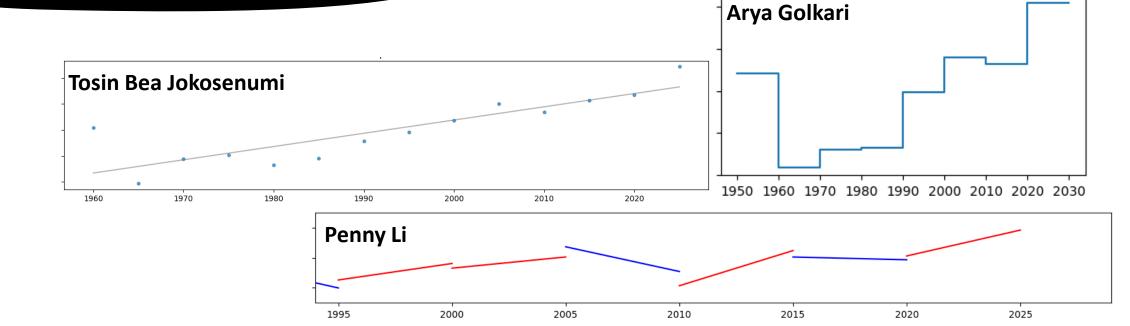


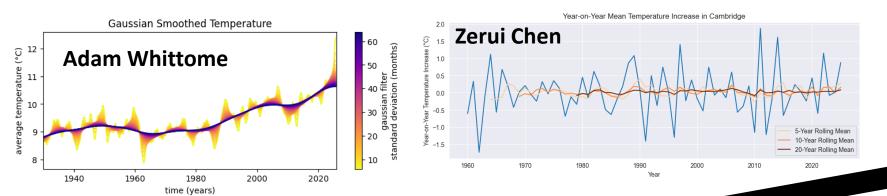




These models don't presume the shape of the response; they're open to anything. But we still need to find a way to summarize all the coefficients.

temp $\approx \beta_1 \sin(2\pi t) + \beta_2 \cos(2\pi t) + \gamma_{\text{decade}}$



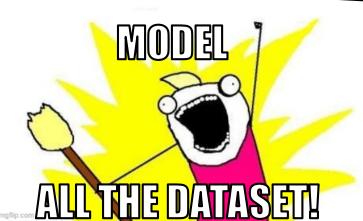


"Don't trust models, just invent algorithms to process the plain honest data"

"I lowkey kind of trust models ..."

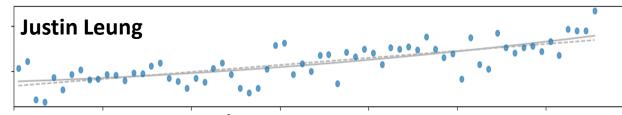
for *s* in stations:

model data from station s as Temp $\sim \alpha + \beta_1 \sin(2\pi t) + \beta_2 \cos(2\pi t) + \gamma t + N(0, \sigma^2)$



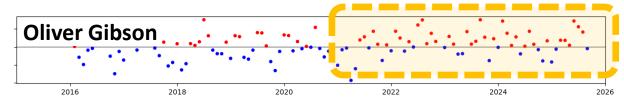
Temp $\approx \alpha_{\text{station}} + \beta_1 \sin(2\pi t) + \beta_2 \cos(2\pi t) + \gamma_{\text{station}} t$

temp
$$\approx \alpha + \beta_1 \sin(2\pi t) + \beta_2 \cos(2\pi t) + \gamma t + \delta t^2$$

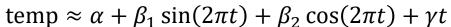


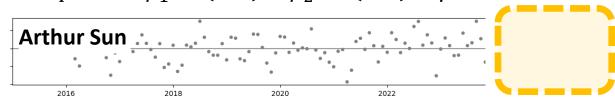
"I found the mle to be $\hat{\delta}$ =0.000757, which is positive, indicating that temperature increase is accelerating"

temp
$$\approx \alpha + \beta_1 \sin(2\pi t) + \beta_2 \cos(2\pi t) + \gamma t$$



"I fitted the model, then examined the signs of the later residuals, and did a hypothesis test to see if they're consistent with Bin(1,½)."





"I did the hypothesis test on a holdout set, to see if the signs of the residuals are consistent with $Bin(1,\frac{1}{2})$."



model choice \approx parameter estimation



generalization



