Example sheet 1

Question 8. For the climate data from section 2.2.5 of lecture notes, we proposed the model

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

in which the  $+\gamma t$  term asserts that temperatures are increasing at a constant rate. We might suspect though that temperatures are increasing non-linearly. To test this, we can create a nonnumerical feature out of t by

u = 'decade\_' + str(math.floor(t/10)) + '0s'

(which gives us values like 'decade\_1980s', 'decade\_1990s', etc.) and fit the model

 $\mathsf{temp} \approx \alpha + \beta_1 \sin(2\pi \mathsf{t}) + \beta_2 \cos(2\pi \mathsf{t}) + \gamma_{\mathsf{u}}.$ 

Write this as a linear model, and give code to fit it. [Note. You should explain what your feature vectors are, then give a one-line command to estimate the parameters.]







-10

0

### §2.3. Diagnosing a model

### After fitting a model,

- 1. Compute the prediction errors a.k.a. the residuals
- 2. Plot them every way we can think of. They're telling us where our model is poor.

Machine learning models don't fail with nice simple exceptions or incorrect answers. They fail by giving us fishy answers.

The only way to debug them is through data science investigation.

#### Lecture schedule

This is the planned lecture schedule. It will be updated as and when actual lectures deviate from schedule. Material marked \* is nonexaminable. **Slides** are uploaded the night before a lecture, and re-uploaded after the lecture with annotations made during the lecture.

Pre-recorded versions of each lecture are available on last year's version of this page.

#### Prerequisites

#### Example sheet 0 and solutions

§1-§4. Lear	rning with probability models		
Lecture 1	1. Learning with probability models		
[slides]	1.1 Specifying probability models		
Lecture 2	1.2 Standard random variables		
[slides]	1.3 Maximum likelihood estimation		
	1.4 Numerical optimization	ev1	
Lecture 3	1.5 Likelihood notation	Try the practical exercises, test your answers on Moodle, discuss with your supervisor.	
[slides]	1.6 Generative models		
	1.7 Supervised learning		
	3.1, 3.2 Prediction accuracy versus probability modelling (* non-examinable)		
Lecture 4	Mock exam duestion and walkthrough		
[slides]	3.3 Neural networks (* non-examinable)		
Lecture 5	2.1 Linear mode <del>iling</del>		
[slides]	2.2 Feature design	For your own fun, good if you want	
	2.3 Diagnosing a linear model	to do moro MI. Submit your answer	
Lecture 6	2.5 The geometry of linear models	to do more ML. Submit your answer	
[slides]	2.6 Interpreting parameters	on Moodle.	
Lecture 7	2.4 Probabilistic linear modelling		
	Example sheet 1 Example sheet	climate	
	OPTIONAL ex1 practical exercises [ex1.jpynb] (for supervisions)		
	OPTIONAL PyTorch introduction and challenge PyTorch introduction	Useful practice if you want to do real data science. Submit your answer on Moodle	
	OPTIONAL climate dataset challenge climate.ipynb		
	Code snippets: fitting.ipynb 🖉, Im.ipynb 🖉		
	Datasets investigated: climate.ipynb 🛃, stop-and-search.ipynb 🛃		

### TODAY'S AGENDA

- §2.3 Model diagnostics  $\checkmark$
- §2.6 Interpreting parameters
- §2.4 Least squares estimation & probability
- §4 Measuring model fit (\* non-examinable)

### §2.6 Interpreting parameters

- Write out the predicted response for a few typical / representative datapoints.
   This helps see what the parameters mean.
- Write out the features.
   If two models have different features but the same feature space, then (once fitted) they make the same predictions on the dataset.
- Check if the features are linearly dependent.
   If so, the parameters have no intrinsic meaning.
   We say the features are *confounded*, and the parameters are *non-identifiable*.

### COMPARING GROUPS

Or



Measurements for condition A:  $a = [a_1, a_2, ..., a_m]$ Measurements for condition B:  $b = [b_1, b_2, ..., b_n]$ Can we use a linear model to compare A and B?

$$x = \alpha_A 1_{and = A} + \alpha_B 1_{and = B}$$

$$\vec{x} = \alpha + \beta 1 \vec{cont} = \theta.$$

For a person of type 
$$A$$
,  $x \approx \alpha$   
 $B$ ,  $x \approx \alpha + \beta$ 

& measures the difference between the two groups.



Exercise 2.6.2 (Contrasts) In the dataset below, of measurements from two groups A and B, interpret the parameters from these models: $y \approx \alpha 1_{g=A} + \beta 1_{g=B} \qquad (M1)$ $y \approx \alpha' + \beta' 1_{g=B} \qquad (M2)$ $y \approx \alpha'' + \beta'' 1_{g=A} + \gamma'' 1_{g=B} \qquad (M3)$ $\frac{g  y}{A  0.5}$	What predictions de these models marke? MI MZ person from group A: X X' person from group B: B X'+B' A MI picks out the predicted responses in each group two groups.
A 1.9 B 3.5	M3: features are $\vec{1}, 1 \vec{j} = A, 1 \vec{j} = B.$
B 1.1 B 2.3	These are linearly dependent: $1_{g=A} + 1_{g=B} = 1$ So the corrangeters are not identifie by
Remark about notation. <b>1</b> means the constant vector [1,1,1,1,1] <b>g</b> is a vector from the dataset, [A,A,B,B,B] <b>f</b> ( $\vec{g}$ ) means "apply the function to each element o <b>1</b> $_{\vec{g}=A}$ means "apply the indicator to each element o	$e.g.  \vec{y} \propto \qquad i.2  l_{\vec{g}=A} + 2.3  l_{\vec{g}=B} \\ \approx \vec{1} + 0.2  l_{\vec{g}=A} + 1.3  l_{\vec{g}} = R \\ \approx 2.3  \vec{1} - 1.1  l_{\vec{g}=A} \end{cases}$

§2.6

## 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

# Guandana News website of the year

Can I set up a model with a parameter that *measures* the quantity I'm interested in?

#### Example 2.6.4

The UK Home Office makes available a dataset of police stop-and-search incidents. We wish to investigate whether there is racial bias in police decisions to stop-and-search. Consider the linear model

$$y_i \approx \alpha + \beta_{eth_i}$$

where  $eth_i$  is the officer-defined ethnicity for record *i*, and  $y_i$  records the outcome:  $y_i = 1$  if the police found something, 0 otherwise.

- a) Write this as a linear equation using one-hot coding.
- b) Are the parameters identifiable? If not, rewrite the model so that they are.
- c) Does the model suggest there is racial bias in policing actions?

(a)  

$$y \approx x \ 1 + \beta_{AS} e_{AS} + \beta_{BI} e_{BI} + \beta_{MI} e_{MI} + \beta_{DM} e_{OH} + \beta_{WH} e_{WH}$$
 where  $e_{R} = 1_{eH} = k$   
(b) They are linearly dependent:  $1 = e_{AS} + e_{BI} + e_{MI} + e_{OH} + e_{WH}$   
So the powarmetos are not identifiable, is we're lidely to  
 $q_{VL}$  silly answers are of linear-model fitting.  
(b)  $q_{VL} = e_{AS} + e_{BI} + e_{MI} + e_{AL} + e_{WH}$ 

The non-identifiable model that was proposed by the question:

 $y \approx \alpha \mathbf{1} + \beta_{\text{As}} \mathbf{e}_{\text{As}} + \beta_{\text{Bl}} \mathbf{e}_{\text{Bl}} + \beta_{\text{Mi}} \mathbf{e}_{\text{Mi}} + \beta_{\text{Oth}} \mathbf{e}_{\text{Oth}} + \beta_{\text{Wh}} \mathbf{e}_{\text{Wh}}$ 

(b) Rewrite it to have identifiable parameters.

(c) Interpret the parameters.

For a person with eth = As predicted 
$$y = \alpha'$$
  
eth = Bl  $= \alpha' + \beta' e l$   
eth = Mi  $= \alpha' + \beta' n i$   
eth = Oth  $= \alpha' + \beta' n n$   
eth = Wh  $= \alpha' + \beta' n n$ .

These Bern marine differences with respect to the barchine of people with eth = Asian,

e.g. if  $\beta_{BE} > 0$ , then the ang response for people with eth = P1 is higher than that for people with eth = As.

### Output from the identifiable model

 $y \approx \alpha' \mathbf{1} + \beta'_{Bl} \mathbf{e}_{Bl} + \beta'_{Mi} \mathbf{e}_{Mi} + \beta'_{Oth} \mathbf{e}_{Oth} + \beta'_{Wh} \mathbf{e}_{Wh}$ 



### Output from the non-identifiable model

8

9

10

11

 $y \approx \alpha + \beta_{\text{As}} 1_{\text{eth}=\text{As}} + \beta_{\text{Bl}} 1_{\text{eth}=\text{Bl}} + \beta_{\text{Mi}} 1_{\text{eth}=\text{Mi}} + \beta_{\text{Oth}} 1_{\text{eth}=\text{Oth}} + \beta_{\text{Wh}} 1_{\text{eth}=\text{Wh}}$ 

Asian Black Mixed Other White

```
ethnicity_levels = np.unique(eth)

eth_onehot = [np.where(eth==k,1,0) for k in ethnicity_levels]

model = sklearn.linear_model.LinearRegression()

model.fit(np.column_stack(eth_onehot), y)

\alpha,\beta s = model.intercept_, model.coef_

print(f'\alpha = {\alpha}')

for k,\beta in zip(ethnicity_levels, \beta s):

print(f'\beta[{k}] = {\beta}')

\alpha = -34037792910.00365

\beta[Asian] = 34037792910.26522

\beta[Black] = 34037792910.265717

\beta[Mixed] = 34037792910.2939

\beta[Other] = 34037792910.2604

\beta[White] = 34037792910.2604
```

# §2.4 Least squares estimation & probability



Carl Friedrich Gauss 1777–1855

### Least squares estimation

Fit the linear model

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

i.e.

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

by choosing the parameters  $\beta_1, ..., \beta_K$  so as to minimize the mean square error

mse = 
$$\frac{1}{n} \sum_{i=1}^{n} \varepsilon_i^2$$

### Maximum likelihood estimation

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 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

 $\mathsf{PL} \approx \alpha + \beta \, \mathsf{SL} + \gamma \, \mathsf{SL}^2$ 

#### 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<sub>i</sub> ~ 
$$\alpha + \beta$$
 SL<sub>i</sub> +  $\gamma$  SL<sup>2</sup><sub>i</sub> +  $N(0, \sigma^2)$   
pewrike it as  
 $\gamma_i \sim \alpha + \beta e_i + \sigma f_i + N(\sigma, \sigma^2)$   
~  $N(\alpha + \beta e_i + \sigma f_i, \sigma^2)$ 

Log likelihood of the dataset:

$$\log \Pr(y_{1}, ..., y_{n}; x_{1}\beta, \gamma, \sigma) = -\frac{n}{2} \log (2\pi\sigma^{2}) - \frac{1}{2\sigma^{2}} \sum_{i=1}^{n} \left[ y_{i} - (\alpha + \beta e_{i} + \delta f_{i}) \right]^{2}$$

We want to maximize this over 0, \$, 8, 5

Maximize over the unknown parameters,  $\alpha$ ,  $\beta$ ,  $\gamma$ , and  $\sigma$ :

 $\max_{(x_{1},y_{2},y_{1},y_{1},y_{2},$ = max max  $\begin{cases} -\frac{N}{2}\log(2\pi\sigma^2) - \frac{1}{2\sigma^2}\sum_{i}(y_i - (x+\beta e_i + \sigma f_i))^2 \end{cases}$  $= \max \left[ -\frac{1}{2} \left( \log \left( 2\pi\sigma^2 \right) + \max \left( -\frac{1}{2\sigma^2} \sum_{i} \left( y_i - \left( \alpha + \beta e_i + \delta f_i \right) \right)^2 \right) \right] \right]$  $= \max_{\sigma} \left[ -\frac{n}{2} \log \left( 2\pi \sigma^2 \right) - \frac{1}{2\sigma^2} \left\{ \min_{\sigma, \beta; \sigma} \sum_{i} \left( y_i - \left( x + \beta e_i + \delta f_i \right) \right)^2 \right\} \right]$  $= \max\left[-\frac{n}{2}\log\left(2\pi\sigma^{2}\right) - \frac{1}{2\sigma^{2}}\sum_{i}\left(y_{i}-\hat{y}_{i}\right)^{2}\right] \quad \text{where} \quad \hat{y}_{i} = \hat{\alpha} + \hat{\beta}e_{i} + \hat{\delta}f_{i}$ obtained by least squares estimation  $\Rightarrow \hat{\sigma} = \int \hat{n} \sum (y_i - \hat{y}_i)$ 

§2.4



In a linear model where the errors have a mean of zero, are uncorrelated, normally distributed, and have equal variances, the best linear unbiased estimator of the parameters is the least-squares estimator.

Maximize over the unknown parameters,  $\alpha$ ,  $\beta$ ,  $\gamma$ , and  $\sigma$ :

max X, B1

Max X, B, 7,0

max

max

= max

マ ら=

Least squares estimation *derives* from a Gaussian probability model.

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

A sensible model diagnostic is to plot a histogram of the residuals, and check they look Gaussian.





### e 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



Alardian

News website of the year

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

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

Fit  $\alpha$  and  $\beta_{Bl}$ ,  $\beta_{Mi}$ , ... 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  $\theta_i$  depends on multiple features. It uses softmax. See the code in [stop-and-search.ipynb], or Part II Advanced Data Science.

### Why is Gaussian distribution important?

- Computational tractability
- Models several operational scenarios
  - Central limit theorem

### Central limit theorem

Let's consider n'independent random variables X., .... Xn with mean and variances M: and C?

Let's consider a neu random variable  $z = \sum_{i=1}^{n} x_i$ z has  $\longrightarrow$  mean  $w = \Sigma \mu$ : Variance  $G^2 = \sum_{i=1}^{n} G^2_i$ 1=1

### Central limit theorem



### Question

Does the central limit theorem have limitations for applied scenarios?

• Well, yes.

### Why is Gaussian distribution important?



Horizontal steam boiler, Augsburg machines, early 19<sup>th</sup> century

Joseph Fourier (1768-1830)

Geille Seult

### "Don't believe me, just watch"

Bruno Mars

### THÉORIE

ANALYTIQUE

### DE LA CHALEUR,

PAR M. FOURIER.



#### A PARIS,

CHEZ FIRMIN DIDOT, PÈRE ET FILS,

LIBRAIRES POUR LES MATHÉMATIQUES, L'ARCHITECTURE HYDRAULIQUE ET LA MARINE, RUE JACOB, Nº 24.

1822.

### Two main results:

- Any function of a variable, whether
   continuous or discontinuous, can be
   expanded in a series of sines of multiples
   of the variable → Fourier transform
- partial differential equation for conductive diffusion of heat

The heat diffusion mechanism is summarized by a stochastic differential equation of this form:

 $\dot{\mathbf{x}} = -\nabla U(\mathbf{x}) + \sqrt{2} \dot{\mathbf{w}},$ 

where  $\dot{\mathbf{x}}$  and  $\dot{\mathbf{w}}$  identify the derivatives with respect to t of  $\mathbf{x}$  and  $\mathbf{w}$ , respectively. Moreover, U is the free energy at  $\mathbf{x}$  (which can be also called the potential at  $\mathbf{x}$ ), and  $\mathbf{w}(t)$  is an *n*-dimensional Brownian motion process.

This says to what rate  $(\dot{x})$ the material at a point will heat up (or cool down) is proportional to how much hotter (or cooler) the surrounding material is

This equation can be expressed in terms of the probability density function of the heat diffusion process (Fokker-Planck derivation):

$$\frac{\partial p}{\partial t} = \nabla \cdot (\nabla p + p \nabla U(\mathbf{x})),$$

A. Marinoni, C. Jutten, M. Girolami, «A graph representation based on fluid diffusion model for data analysis: theoretical aspects and enhanced community detection,» subm. to IEEE Trans. On Pattern Analaysis and Machine Intelligence (TPAMI), <a href="https://arxiv.org/abs/2112.04388">https://arxiv.org/abs/2112.04388</a>, 2023

The solution of this equation is:

$$p(\mathbf{x}(t+\epsilon) = \mathbf{x}_j | \mathbf{x}(t) = \mathbf{x}_i) = p_{ij} \propto \exp\left[-\frac{||\mathbf{x}_i - \mathbf{x}_j||_2}{2\sigma}\right]$$

Which is the probability of the given thermodynamic system to transition from state  $\mathbf{x}_i$  to state  $\mathbf{x}_j$ 





# §4. How should we measure how well a model fits the data? (\* non-examinable)



#### Climate is stable:

Temp(
$$t$$
) ~  $a + b \sin(2\pi(t + \phi)) + N(0, \sigma^2)$ 

Temperatures are increasing linearly: Temp(t) ~  $\cdots + \gamma t$ 



Temperatures are increasing, and the rate is nonlinear:



Temperatures are increasing, and the rate is increasing piecewise-linearly:



And if so, when is the tipping point?





### RICE CRUMB #3 – v2.0

• How can we write the distribution of the parameters estimated by MLE for  $N \rightarrow +\infty$ ?

... keep the central limit theorem in mind ...