This week we're building up our skills at inventing useful probability models.

MONDAY

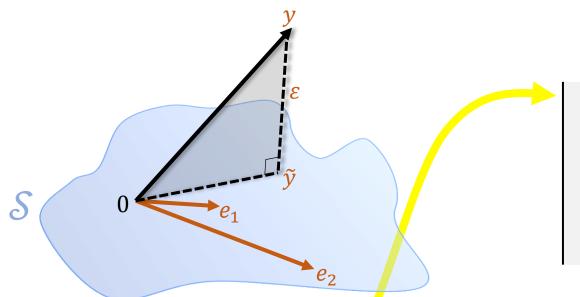
linear models & feature design (to quickly turn our ideas into easy-to-fit models)

WEDNESDAY

- "debugging" models
- linear algebra interpretation

FRIDAY

- identifiability of parameters
- the link between least squares and likelihood



What does "closest" even mean? It means: find $\tilde{y} \in \mathcal{S}$ to minimize $||y - \tilde{y}||$

i.e. to minimize $\sqrt{\sum_i \varepsilon_i^2}$ where $\varepsilon = y - \tilde{y}$.

The minimization is over all $\tilde{y} \in \mathcal{S}$, i.e. over all linear combinations of $\{e_1, \dots, e_K\}$.

GEOMETRY OF PROJECTIONS

Let S be the span of $\{e_1, \dots, e_K\}$.

What's the closest we can get to y, while staying in S?

mean square error, $\frac{1}{\text{\#datapoints}} \sum_{i} \varepsilon_i^2$

where $\varepsilon = y - (\beta_1 e_1 + \cdots + \beta_K e_K)$

What does "best approximation" mean?

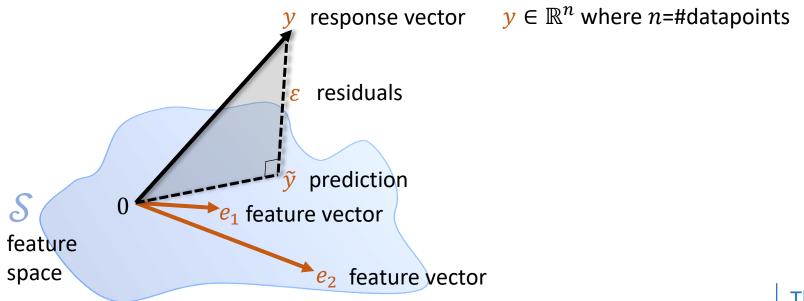
It means: find β_1, \dots, β_K to minimize the

LINEAR MODELLING / LEAST SQUARES ESTIMATION

Given features $\{e_1, \dots, e_K\}$ let's approximate $y \approx \beta_1 e_1 + \dots + \beta_K e_K$.

What parameters give us the best approximation?

These are exactly the same thing!



The span of $\{e_1, ..., e_K\}$ is called the *feature* space

- Fitting a linear model \equiv projecting the data y onto the feature space S
- lacktriangle The fitted parameters eta_k are the "coordinates" of the predicted response, $ar{y}=\sum_k eta_k m{e_k}$
- If two models have different features but the same feature space, they'll make exactly the same prediction \tilde{y} (but of course have different β_k)
- If the feature vectors are linearly independent, then the β_k are uniquely identified. Otherwise, the parameters have no intrinsic meaning, and we say they're *confounded*.

Exercise 2.6.2 (Contrasts)

For the dataset below, of measurements from two groups A and B, discuss the differences between these three models:

$$y \approx \alpha 1_{g=A} + \beta 1_{g=B} \tag{M1}$$

$$y \approx \alpha' + \beta' 1_{g=B} \tag{M2}$$

$$y \approx \alpha'' + \beta'' 1_{g=A} + \gamma'' 1_{g=B}$$
 (M3)

g	у
Α	0.5
Α	1.9
В	3.5
В	1.1
В	2.3

These span the same sporce, identical predictions on the dataset (though they usk different parameterizations)

Those are linearly dependent, so the parameter one not identifiable — there are multiple equivalent ways of writing out the same response vector.

$$\vec{y} \propto \frac{1.2 \, 1_{\vec{g}=A} + 2.3 \, 1_{\vec{g}=B}}{2 \, 1 + 0.2 \, 1_{\vec{g}=A} + 1.3 \, 1_{\vec{g}=B}}$$

$$\approx 2.3 \, \vec{1} - 1.1 \, 1_{\vec{g}=A}.$$

§2.6 Interpreting parameters

- Check if the features are linearly dependent.
 If they are, the parameters have no intrinsic meaning.
 We say the features are confounded, and the parameters are non-identifiable.
- Write out the predicted response for a few typical / representative datapoints.
 This helps see what the parameters mean.
- It may be useful to use different features
 (with the same feature space).
 This doesn't change the model's predictions, but it does change the parameterization. It lets us extract the contrasts i.e. comparisons we're interested in.

🖰 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

Giardian News website of the year

Do I trust this finding?

- Proportionate to what?
- Is it cherry picking?
- Why the scare quotes?

🖰 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

Giardian News website of the year

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

force	Date	LatLng	Object of search	Gender	Age range	Officer-defined ethnicity	Outcome		
cambridgeshire	2023-08-31 15:44:04+00:00	(52.43,-0.142)	Controlled drugs				A no further action disposal		
cambridgeshire	2023-08-31 15:35:41+00:00	(52.43,-0.142)	Firearms	Male	25-34	White	Khat or Cannabis warning		
cambridgeshire	2023-08-31 14:44:04+00:00	(52.43,-0.142)	Firearms	Male	25-34	White	Khat or Cannabis warning		
cambridgeshire	2023-08-31 03:44:14+00:00	(52.58,-0.244)	Offensive weapons	Male		Other	A no further action disposal		
cambridgeshire	2023-08-31 02:34:16+00:00	(52.59,-0.247)	Controlled drugs	Male	25-34	White	Arrest		
cambridgeshire	2023-08-31 02:27:10+00:00	(52.21,0.124)	Controlled drugs	Male	18-24	White	A no further action disposal		
cambridgeshire	2023-08-30 22:28:13+00:00	(52.45,-0.117)	Controlled drugs	Female	over 34	White	A no further action disposal		
cambridgeshire	2023-08-30 20:24:13+00:00	(52.32,-0.0708)	Controlled drugs	Male	10-17	White	Summons / charged by post		
cambridgeshire	2023-08-30 14:26:58+00:00	(52.57,-0.24)	Controlled drugs	Male	over 34	Asian	A no further action disposal		
cambridgeshire	2023-08-30 14:13:45+00:00	(52.57,-0.24)	Controlled drugs	Male	25-34	Black	Arrest		
Log of England+Wales stop-and-search incidents, from the UK home office https://data.police.uk/									

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 model 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?

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.

For records with $eth_i = k$, the average response is

$$\frac{\sum_{i: \operatorname{eth}_{i} = k} y_{i}}{|\{i: \operatorname{eth}_{i} = k\}|}$$

This is just #finds/#stops, i.e. $\mathbb{P}(\text{find something})$

For records with $eth_i = k$, the predicted response is $\alpha + \beta_k$.

Actual responses y_i may be above or below this prediction, but their average is $\alpha + \beta_k$.

This model is thus proposing: $\mathbb{P}(\text{find something}) \approx \alpha + \beta_k$ for a person in ethnic group k

If $\beta_k < 0$, that means $\mathbb{P}(\text{find something})$ is low compared to other ethnic groups, i.e. the police are stopping relatively more innocent people.

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

$$y_i \approx \alpha + \beta_{eth_i}$$
 \hat{a} $\hat{y} \approx \alpha \hat{1} + \sum_{k} \beta_k \hat{1}_{eth_i} = k$

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.

e.g. for eth=mixed,
$$y \approx \propto + \beta_{Mi}$$

a) Write this as a linear model using one-hot coding.

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

α,βs = model.intercept_, model.coef_

print(f'α = {α}')
for k,β in zip(ethnicity_levels, βs):
    print(f'β[{k}] = {β}')
```

```
\alpha = -34037792910.00365

\beta[Asian] = 34037792910.26522

\beta[Black] = 34037792910.265717

\beta[Mixed] = 34037792910.2939

\beta[Other] = 34037792910.260

\beta[White] = 34037792910.26
```

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

$$y_i \approx \alpha + \beta_{\text{eth}_i}$$

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 model using one-hot coding. $\vec{j} \approx \alpha \vec{l} + \vec{l}_{eff_i} = k$
- b) Are the parameters identifiable? If not, rewrite the model so that they are.

(a)
$$y \approx \alpha \mathbf{1} + \beta_{As} \mathbf{1}_{eth=As} + \beta_{Bl} \mathbf{1}_{eth=Bl} + \beta_{Mi} \mathbf{1}_{eth=Mi} + \beta_{Oth} \mathbf{1}_{eth=Oth} + \beta_{Wh} \mathbf{1}_{eth=Wh}$$

Not identificable, since the features one linearly dependent:
$$\vec{I} = 1_{eH} = A_S + 1_{eH} = B_L + \cdots + 1_{eH} = Wh$$
. To get linear independence:

— we could direct the \vec{I} feature, $\vec{y} \approx \beta_{AS} 1_{eH} = A_S + \cdots + \beta_{Wh} 1_{eH} = Wh$

— of we could direct one of the other, $\vec{y} \approx \beta_{BL} 1_{eH} = B_L + \cdots + \beta_{Wh} 1_{eH} = Wh$

c) Does this model suggest there is racial bias in policing actions?

First, let's interpret the parameters of the rewritten model. $\vec{J} \approx \alpha + \beta_{BL} 1_{eH} = BL + \cdots + \beta_{Wh} 1_{eH} = Wh$

For a person with eth = As predicted y = x

$$eth = Bl$$

$$eth = Mi$$

$$eth = Oth$$

$$eth = Wh$$

These β_k' measure <u>differences</u> with respect to the baseline of people with eth=Asian.

e.g. if we find $\beta_{\rm Bl}'>0$ it's telling us that the average response for people with eth=Bl is higher than that for people with eth=As

c) Does this model suggest there is racial bias in policing actions?

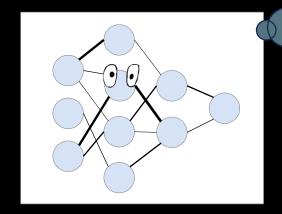
```
This or parameter is P(police findsthy) for eth=Asian
                                                                                      i.e. P(poliefiedsthg) is 0.00 082 longer

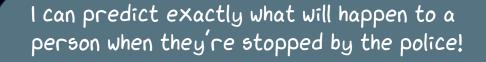
for eth=Black than for eth = Asian.
     some_levels = [k for k in ethnicity_levels if k != 'Asian']
     eth_onehot = [np.where(eth==k,1,0) for k in some_levels]
     model = sklearn.linear_model.LinearRegression()
     model.fit(np.column_stack(eth_onehot), y)
     \alpha, \beta s = model.intercept_, model.coef_
                                                                                                          Are these meaningful?
                                                           \alpha' = 0.261423626892284
     print(f'\alpha = \{\alpha\}')
                                                                                                          Is this evidence of bias,
                                                           \beta'[Black] = 0.0008154705731564
10
     for k, \beta in zip(some_levels, \betas):
                                                                                                          or is it just noise?
                                                           \beta'[Mixed] = 0.02871561721115
11
          print(f'\beta[\{k\}] = \{\beta\}')
                                                                                                          [See next two weeks]
                                                           \beta'[Other] = -0.0044710573665
                                                           \beta'[White] = -0.00372037824708
```

In a dataset of police stop-and-search records, we've looked for evidence of ethnic bias. What about gender bias? Is the net bias simply additive, or do these biases *intersect*?

Additive model: Yix & + Beth: + Vgender;

Intersectional model: Yi & Sethi, gendering

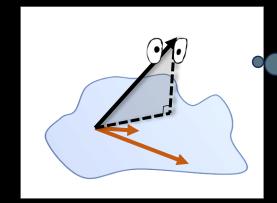




Just tell me their gender. And ethnicity.

And location. And whether they're left or right handed. And whether they have a pet cat or a dog. And what their pet is called. ...

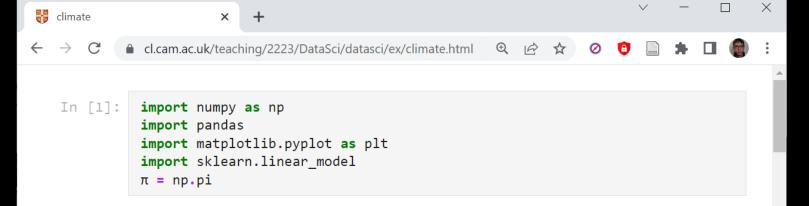
What was the question again?



Great question! I can answer it directly; all I need is a model with cunningly chosen features.

(just don't ask me whether my model fits well)

(There's a big research area called AI Explainability / Alignment, trying to bridge this gap.)



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.

Upload your answers to
Moodle by Sunday evening
for presentation / discussion
next week