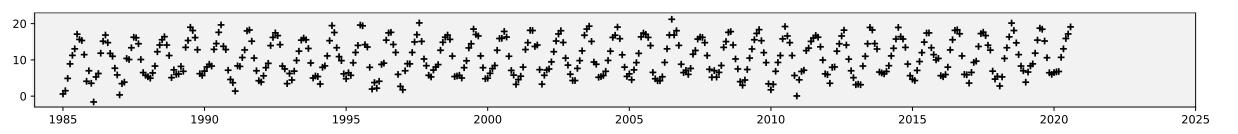
§2.1 Fitting a linear model

Monthly average temperatures in Cambridge, UK

What's a good model for this dataset?



Climate is stable?

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

Temperatures are increasing?

Temperatures are increasing, and the increase is accelerating?

The extremes are getting worse?

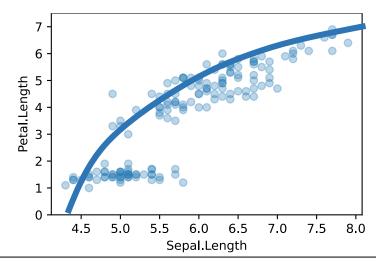
There are so many possible models. We want to make it easy to invent and fit new models, so we have time to explore all the possibilities.

Example 2.1.1

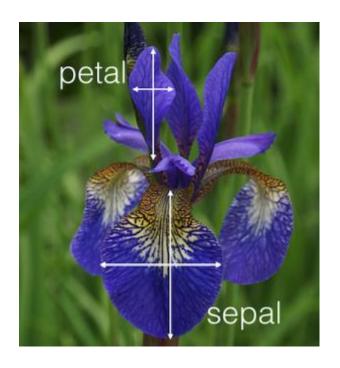
The Iris dataset has 50 records of iris measurements, from three species.

Petal. Length	Petal. Width	Sepal. Length	Sepal. Width	Species
1.0	0.2	4.6	3.6	setosa
5.0	1.9	6.3	2.5	virginica
5.8	1.6	7.2	3.0	virginica
4.2	1.2	5.7	3.0	versicolor

How does Petal.Length depend on Sepal.Length?



Let's guess that for parameters α , β , γ , σ (to be estimated), Petal.Length $\approx \alpha + \beta$ Sepal.Length + γ (Sepal.Length)²

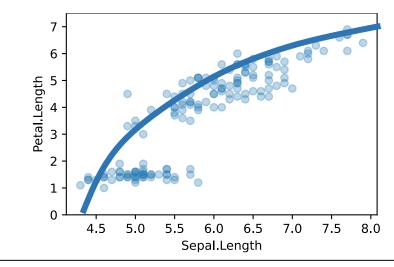


Example 2.1.1

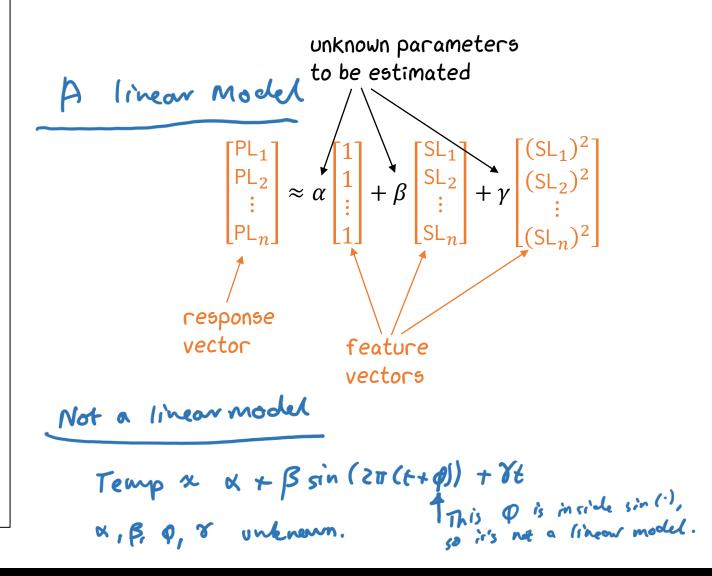
The Iris dataset has 50 records of iris measurements, from three species.

ies	Specie	Sepal. Width	Sepal. Length	Petal. Width	Petal. Length
osa	setos	3.6	4.6	0.2	1.0
nica	virginio	2.5	6.3	1.9	5.0
nica	virginio	3.0	7.2	1.6	5.8
olor	versicol	3.0	5.7	1.2	4.2

How does Petal.Length depend on Sepal.Length?



Let's guess that for parameters α , β , γ , σ (to be estimated), Petal.Length $\approx \alpha + \beta$ Sepal.Length + γ (Sepal.Length)²



Let's guess that for parameters α , β , γ , σ (to be estimated), Petal.Length $\approx \alpha + \beta$ Sepal.Length + γ (Sepal.Length)²

A linear model is

- a supervised-learning model [it has a response variable, and feature variables]
- in which the response and the features are numeric
- and the response vector is predicted by a linear combination of (known) feature vectors weighted by (unknown) parameters

$$\begin{bmatrix} \mathsf{PL}_1 \\ \mathsf{FL}_2 \\ \mathsf{PL}_1 \end{bmatrix} \approx \alpha \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} + \beta \begin{bmatrix} \mathsf{SL}_1 \\ \mathsf{SL}_2 \\ \vdots \\ \mathsf{SL}_n \end{bmatrix} + \gamma \begin{bmatrix} (\mathsf{SL}_1)^2 \\ (\mathsf{SL}_2)^2 \\ \vdots \\ (\mathsf{SL}_n)^2 \end{bmatrix}$$

Models of this form are called *linear models* (because they're based on linear algebra).

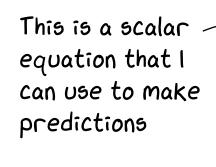
They are flexible, and very fast to optimize.



Petal.Length $\approx \alpha + \beta$ Sepal.Length $+ \gamma$ (Sepal.Length)²



This is a vector equation describing my dataset, one row per datapoint



Least squares estimation

Consider a linear model

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

"All models are wrong."

The vector of prediction errors is called the residual vector,

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

We can fit the model using *least squares estimation*. This means finding parameters β_1, \dots, β_K to minimize the mean square error

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

iris = pandas.read_csv(...)

Petal.Length $\approx \alpha + \beta$ Sepal.Length + γ (Sepal.Length)²

$$\begin{bmatrix} \mathsf{PL}_1 \\ \mathsf{PL}_2 \\ \vdots \\ \mathsf{PL}_n \end{bmatrix} \approx \alpha \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} + \beta \begin{bmatrix} \mathsf{SL}_1 \\ \mathsf{SL}_2 \\ \vdots \\ \mathsf{SL}_n \end{bmatrix} + \gamma \begin{bmatrix} (\mathsf{SL}_1)^2 \\ (\mathsf{SL}_2)^2 \\ \vdots \\ (\mathsf{SL}_n)^2 \end{bmatrix}$$

Fitting the model

```
one, SL, PL = np.ones(len(iris)), iris['Sepal.Length'], iris['Petal.Length']
model = sklearn.linear_model.LinearRegression(fit_intercept=False)
model.fit(np.column_stack([one, SL, SL**2]), PL)
(α,β,γ) = model.coef_
```

Making predictions / getting fitted values from the model

```
newSL = np.linspace(4.2, 8.2, 20)

predPL = \alpha + \beta*newSL + \gamma*(newSL**2)
```

Petal.Length $\approx \alpha + \beta$ Sepal.Length $+ \gamma$ (Sepal.Length)²

$$\begin{bmatrix} \mathsf{PL}_1 \\ \mathsf{PL}_2 \\ \vdots \\ \mathsf{PL}_n \end{bmatrix} \approx \alpha \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} + \beta \begin{bmatrix} \mathsf{SL}_1 \\ \mathsf{SL}_2 \\ \vdots \\ \mathsf{SL}_n \end{bmatrix} + \gamma \begin{bmatrix} (\mathsf{SL}_1)^2 \\ (\mathsf{SL}_2)^2 \\ \vdots \\ (\mathsf{SL}_n)^2 \end{bmatrix}$$

```
Fitting the model (cleaner code)
```

```
SL, PL = iris['Sepal.Length'], iris['Petal.Length']
```

model = sklearn.linear_model.LinearRegression()

model.fit(np.column_stack([SL, SL**2]), PL)

 α ,(β , γ) = model2.intercept_, model2.coef_

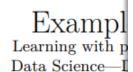
sklearn.linear_model puts in the [I,...,I] feature for us by default. If we don't want it, we have to specify fit intercept=False.

Making predictions / getting fitted values from the model (cleaner code)

```
6 newSL = np.linspace(4.2, 8.2, 20)
```

7 predPL = model.predict(np.column_stack([newSL, newSL**2]))

This saves us from having to explicitly code up the prediction formula — less chance we introduce bugs



Some of the questions ask for pseudocode. I suggonline tester. There is a notebook with template the course materials webpage.

Question 1. Given a dataset $[x_1, \ldots, x_n]$, we wrandom variable with a single parameter $\lambda > 0$,

$$\Pr(x; \lambda) = \frac{\lambda^x e^{-\lambda}}{x!}$$

Show that the maximum likelihood estimator fo

Question 2. Give pseudocode to fit the mod [Optional.] If you want to test your code using PoissonModel.

In practice it'd be daft to use numerical opt answer. But it's good to get used to numerical problem where you what the answer should be.

Question 3. Given a dataset $[x_1, \ldots, x_n]$, we is unknown. Show that the maximum likelihood

ex1: practical exercises for Example Sheet 1.

The example sheet asks you to implement the three classes given below: PoissonModel, PiecewiseLinearModel, and StepPeriodModel. The class skeletons are given, and you should fill in the missing pieces. To test your answers on Moodle, please upload either a Jupyter notebook called ex1.ipynb or a plain Python file called ex1.py.

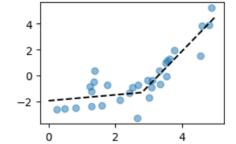
```
import numpy as np
import scipy.optimize
import sklearn.linear_model
```

Poisson model. Suppose we're given a dataset $[x_1, \ldots, x_n]$. We wish to fit the model that says each x_i is an independent sample from the Poisson(λ) distribution. Estimate λ using scipy optimize. fmin.

Note. If the tester reports that your answer is a little bit off, try increasing the precision that <code>scipy.optimize.fmin</code> is using by e.g. passing in the argument <code>xtol=0.00001</code>.

```
class PoissonModel():
    def __init__(self):
        self.\lambda_ = np.nan
    def fit(self, x):
        # Input: x is a numpy vector of integers
        # TODO: set self.\lambda_
```

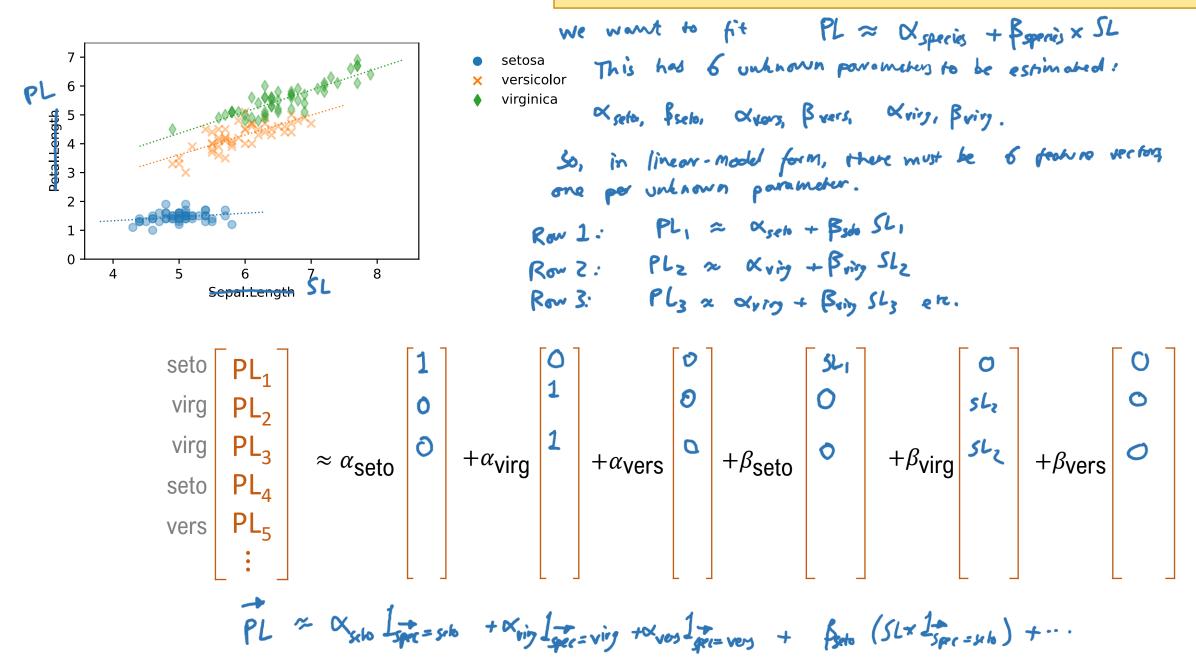
Piecewise linear response. Suppose we're given a dataset of (x_i, y_i) pairs. We wish to fit a model for y as a function of x, made up of two straight lines. The function must be continuous, i.e. the two straight lines must meet at an inflection point. The x-coordinate of the inflection point is given.



§2.2 Feature design

How do we design features, so that linear models answer the questions we want answered?

Build a model to answer science questions machine learning data building Build a model to models science make predictions from data



I'm using numpy's notation for handling vectors:

- $\vec{1}$ is the constant vector [1,1,1,1,1]
- spec is a vector from the dataset, ["seto", "virg", "virg", "seto", ...]
- $f(\vec{x})$ means "apply the function to each element of \vec{x} "
- $1_{\overrightarrow{\text{spec}}=\text{"seto"}}$ means "apply the indicator to each element of $\overrightarrow{\text{spec}}$ "
- \overrightarrow{SL} * $1_{\overrightarrow{spec}="seto"}$ uses element-wise multiplication

```
import sklearn.linear_model

species, SL, PL = iris['Species'], iris['Petal.Length'], iris['Sepal.Length']

species_levels = ['setosa', 'virginica', 'versicolor']

i1,i2,i3 = [np.where(species==k, 1, 0) for k in species_levels]

X = np.column_stack([i1, i2, i3, i1*SL, i2*SL, i3*SL])

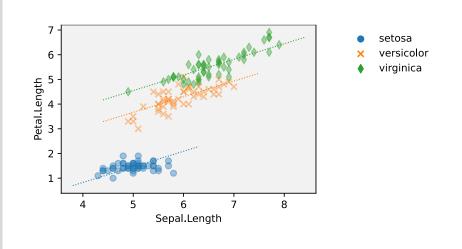
m = sklearn.linear_model.LinearRegression(fit_intercept=False)

m.fit(X, PL)

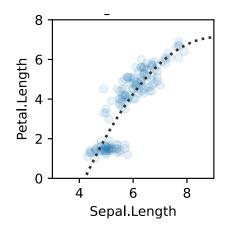
m.coef_
```

EXERCISE

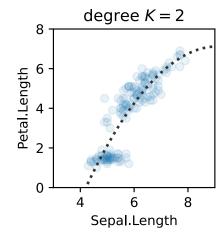
Fit the model with three parallel straight lines.

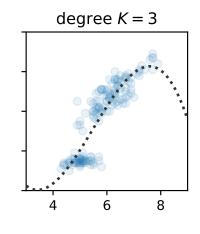


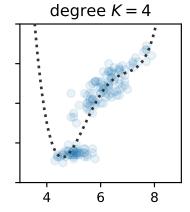
NON-LINEAR RESPONSE

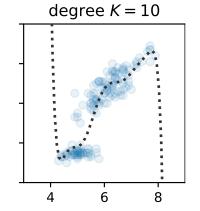


Petal.Length \approx $\alpha + \beta$ Sepal.Length $+ \gamma$ (Sepal.Length)²









Petal.Length
$$\approx$$

$$\beta_0 + \sum\nolimits_{k=1}^K \beta_k (\text{Sepal.Length})^k$$

Q. Should we just keep adding more and more features to our model?

(seeing as the more features we add, the better we can fit the dataset)

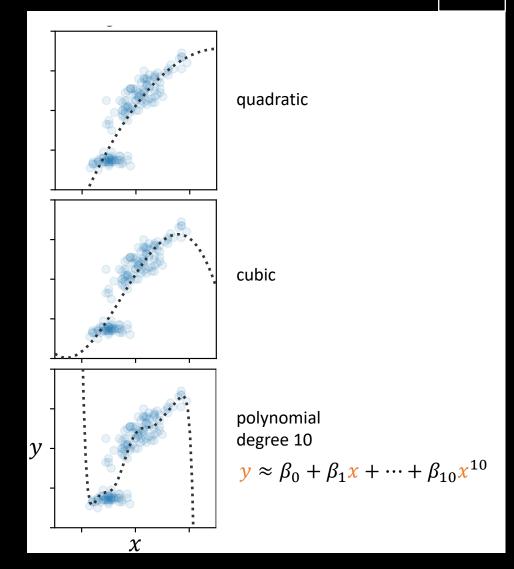
A. No. If we did, we'd *overfit*.

Only add in features that you (as a scientist) believe are relevant.

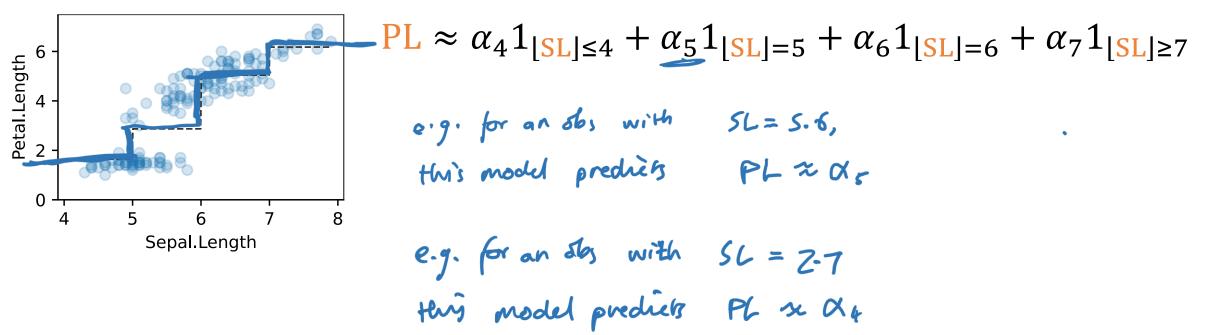
Or, evaluate your model on a holdout set.

If your model is overfitted to the *training* data, it'll perform poorly on *holdout* data.

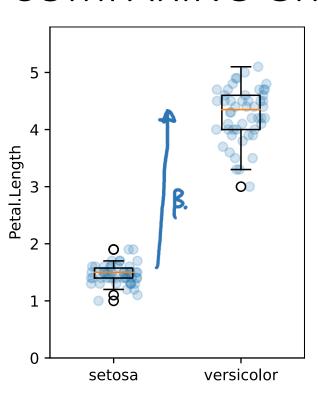
[\$4, 8-10]



NON-LINEAR RESPONSE via one-hot coding



COMPARING GROUPS



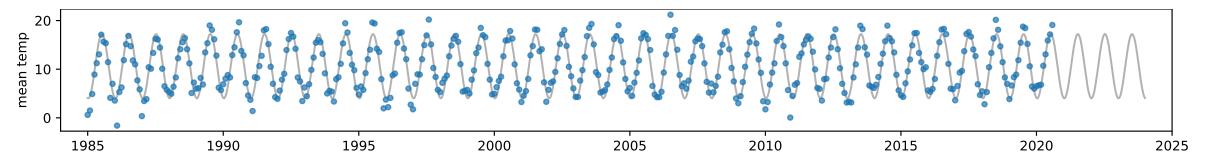
Measurements for condition A: $a = [a_1, a_2, ..., a_m]$ Measurements for condition $B: b = [b_1, b_2, ..., b_n]$

Can we use a linear model to compare A and B?

$$\frac{1}{x} = x + \beta \frac{1}{3} = B.$$
For a person of group A,

 $\dot{x} = \alpha + \beta 1_{\dot{g}=B}$.

For a person of group A, $x \approx \alpha$ the difference between the groups.

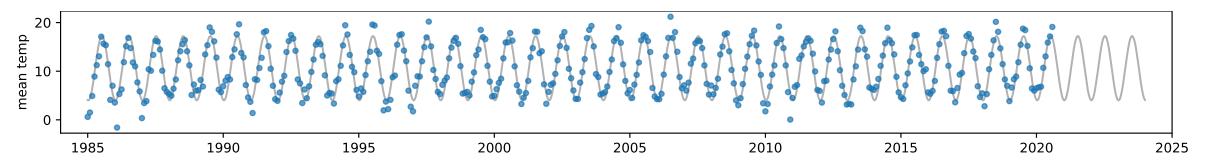


We'd like to fit the model:

$$temp \approx \alpha + \beta \sin(2\pi(t+\varphi))$$

It looks like we can't use sklearn. LinearRegression. That's only for linear models, e.g.

temp $\approx \alpha + \beta e + \gamma f$ Instead, we could just bruteforce optimize it with scipy.optimize.fmin.



We'd like to fit the model:

$$temp \approx \alpha + \beta \sin(2\pi t + \varphi)$$

$$\approx \alpha + \beta \left\{ \sin(2\pi t) \cos \phi + \cos(2\pi t) \sin \phi \right\}$$

$$= \alpha + (\beta \cos \phi) \sin(2\pi t) + (\beta \sin \phi) \cos(2\pi t)$$

$$= \alpha + \beta_1 \sin(2\pi t) + \beta_2 \cos(2\pi t)$$

$$= \alpha + \beta_1 e + \beta_2 f$$

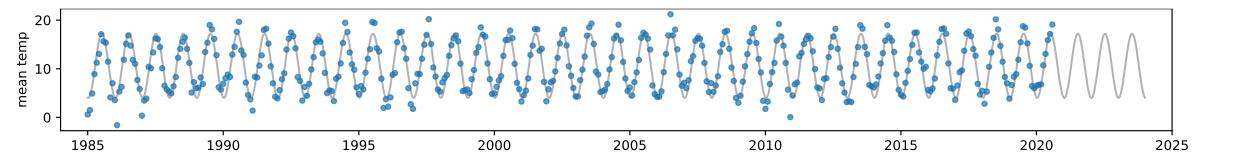
From secondary school trigonometry,

$$sin(A + B)$$

$$= sin(A) cos(B) + cos(A) sin(B)$$

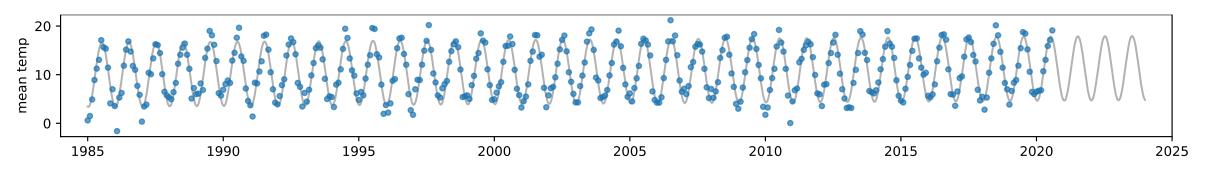
a linear model with feature vectors 1, $sin(2\pi t)$, $cos(2\pi t)$

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

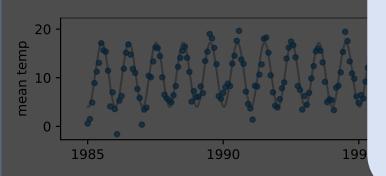


PERIODIC PATTERN + SECULAR TREND

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

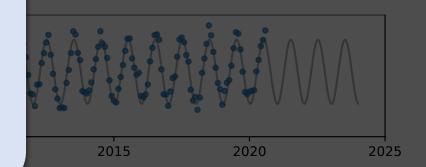


 $temp \approx \alpha + \beta_1 \sin(2\pi t) +$



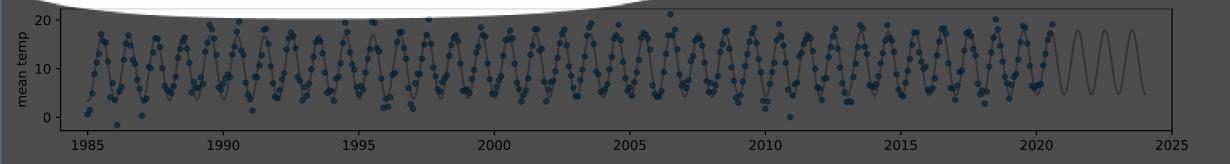
Linear models are easily composable.

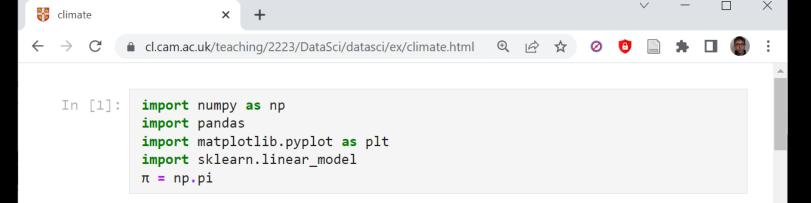
Here I added on the "linear trend" feature. We could easily have stuck in the step-function response instead.



PERIODIC PATTERN + SECULAR TREND

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





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
for presentation / discussion
next week

You've got to have models in your head. And you've got to array your experience – both vicarious and direct – on this latticework of models.

You may have noticed students who just try to remember and pound back what is remembered. Well, they fail in school and in life. You've got to hang experience on a latticework of models in your head.

Charlie Munger (business partner of Warren Buffet), A lesson on elementary, worldly wisdom as it relates to investment management & business.