

RICE CRUMB #2

$$\lim_{N \rightarrow \infty} P(\|\hat{\underline{\theta}}_{ML} - \underline{\theta}_0\| < \varepsilon) =$$

$$\lim_{N \rightarrow \infty} E[\|\hat{\underline{\theta}}_{ML} - \underline{\theta}_0\|^2] =$$

RICE CRUMB #2

$$\lim_{N \rightarrow \infty} P(\|\hat{\underline{\theta}}_{ML} - \underline{\theta}_0\| < \epsilon) = 1 \quad \forall \epsilon$$

CONVERGENCE IN PROBABILITY

$$\lim_{N \rightarrow \infty} E[\|\hat{\underline{\theta}}_{ML} - \underline{\theta}_0\|^2] = 0$$

CONVERGENCE IN MEAN SQUARE

(ASYMPTOTICALLY CONSISTENT)

§2.1

Fitting a
linear model

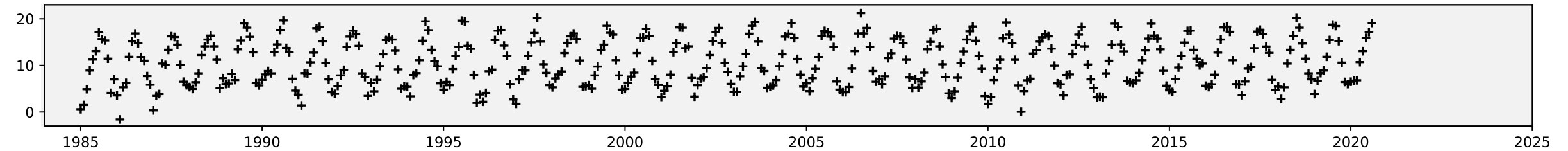
You've got to have models in your head. And you've got to array your experience – both vicarious and direct – on this latticework of models.

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.

Monthly average temperatures in Cambridge, UK

What's a good model for this dataset?

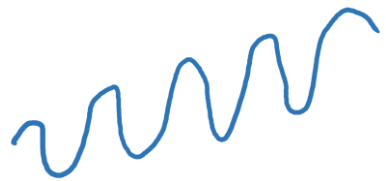


Climate is stable?

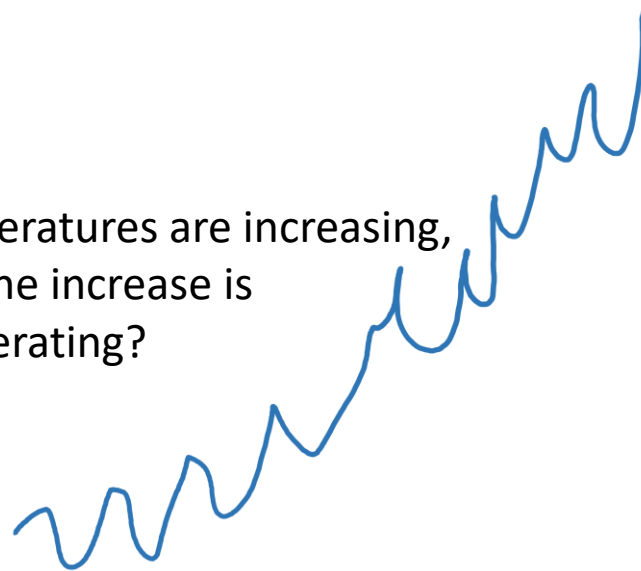
$$\text{Temp}(t) \sim a + b \sin(2\pi(t + \phi)) + N(0, \sigma^2)$$



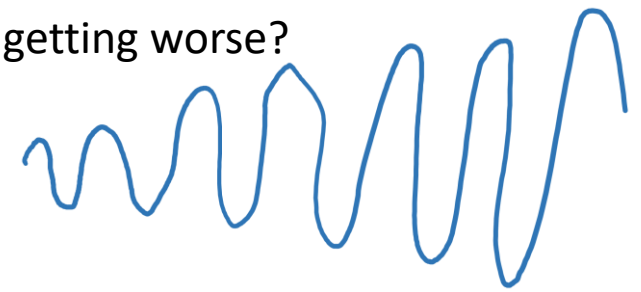
Temperatures are increasing?



Temperatures are increasing,
and the increase is
accelerating?



The extremes are
getting worse?



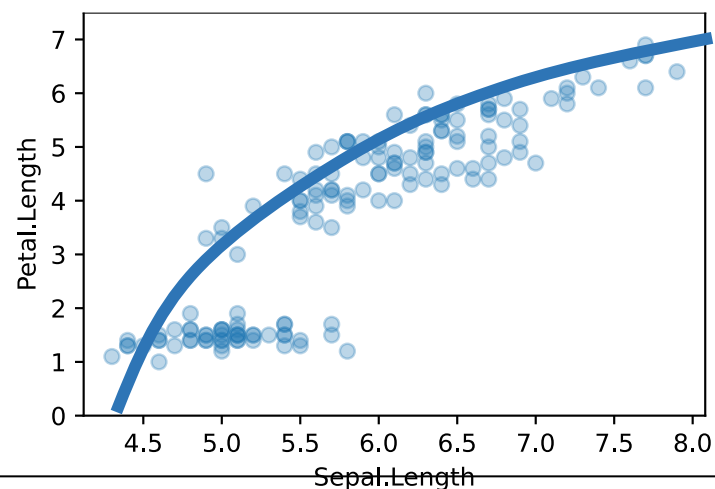
There are so many possible models. We want to make it easy to invent and fit new models, so we have time to explore all the possibilities.

Example 2.1.1

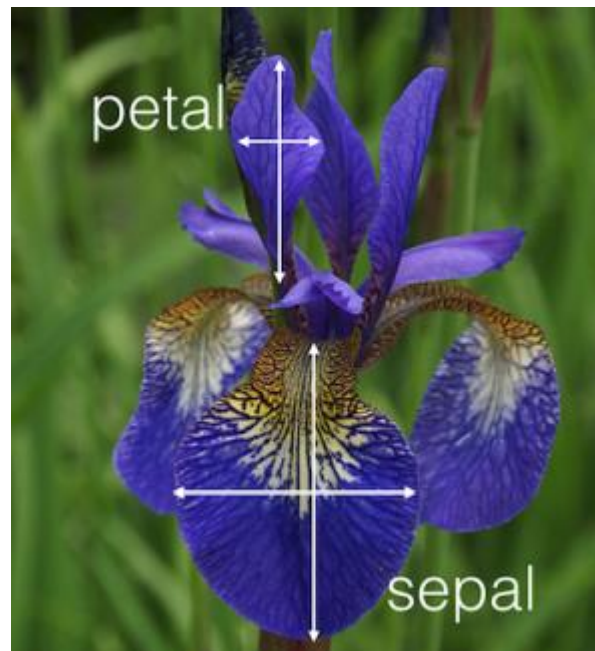
The Iris dataset has 50 records of iris measurements, from three species.

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 $\alpha, \beta, \gamma, \sigma$ (to be estimated),
Petal.Length $\approx \alpha + \beta$ **Sepal.Length** + $\gamma(\text{Sepal.Length})^2$

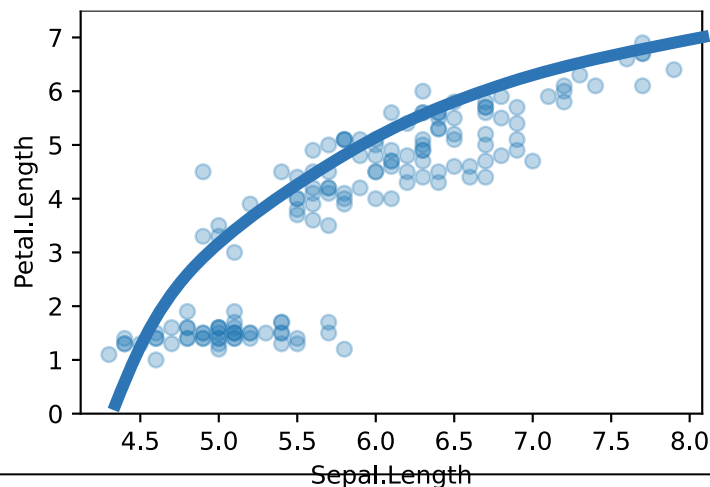


Example 2.1.1

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 $\alpha, \beta, \gamma, \sigma$ (to be estimated),
Petal.Length $\approx \alpha + \beta$ **Sepal.Length** $+ \gamma(\text{Sepal.Length})^2$

Linear Model

unknown parameters to be estimated

$$\begin{bmatrix} PL_1 \\ PL_2 \\ \vdots \\ PL_n \end{bmatrix} \approx \alpha \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} + \beta \begin{bmatrix} SL_1 \\ SL_2 \\ \vdots \\ SL_n \end{bmatrix} + \gamma \begin{bmatrix} (SL_1)^2 \\ (SL_2)^2 \\ \vdots \\ (SL_n)^2 \end{bmatrix}$$

response vector

feature vectors

Not a linear model

$$\text{Temp} \approx \alpha + \beta \sin(2\pi(t + \phi)) + \gamma t$$

Let's guess that for parameters $\alpha, \beta, \gamma, \sigma$ (to be estimated),
 $\text{Petal.Length} \approx \alpha + \beta \text{Sepal.Length} + \gamma(\text{Sepal.Length})^2$

Supervised learning

Response & features are numeric

Response is predicted by a
linear combination of
(known) feature vectors,
weighted by unknown parameters.

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

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

They are flexible, and very fast to optimize.

Least squares estimation

Consider a linear model

$$\mathbf{y} \approx \beta_1 \mathbf{e}_1 + \cdots + \beta_K \mathbf{e}_K$$

“All models are wrong.”

The vector of prediction errors is called the *residual vector*,

$$\boldsymbol{\varepsilon} = \mathbf{y} - (\beta_1 \mathbf{e}_1 + \cdots + \beta_K \mathbf{e}_K)$$

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

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

```
1 iris = pandas.read_csv(...)
```

$$\text{Petal.Length} \approx \alpha + \beta \text{ Sepal.Length} + \gamma (\text{Sepal.Length})^2$$

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

Fitting the model

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

Making predictions / getting fitted values from the model

```
6 newSL = np.linspace(4.2, 8.2, 20)
7 predPL = α + β*newSL + γ*(newSL**2)
```

```
1 iris = pandas.read_csv(...)
```

$$\text{Petal.Length} \approx \alpha + \beta \text{ Sepal.Length} + \gamma (\text{Sepal.Length})^2$$

$$\begin{bmatrix} PL_1 \\ PL_2 \\ \vdots \\ PL_n \end{bmatrix} \approx \alpha \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} + \beta \begin{bmatrix} SL_1 \\ SL_2 \\ \vdots \\ SL_n \end{bmatrix} + \gamma \begin{bmatrix} (SL_1)^2 \\ (SL_2)^2 \\ \vdots \\ (SL_n)^2 \end{bmatrix}$$

Fitting the model (cleaner code)

```
2 SL, PL = iris['Sepal.Length'], iris['Petal.Length']
3 model = sklearn.linear_model.LinearRegression()
4 model.fit(np.column_stack([SL, SL**2]), PL)
5 α, (β, γ) = model2.intercept_, model2.coef_
```

sklearn.linear_model puts in the $\begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$ 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)

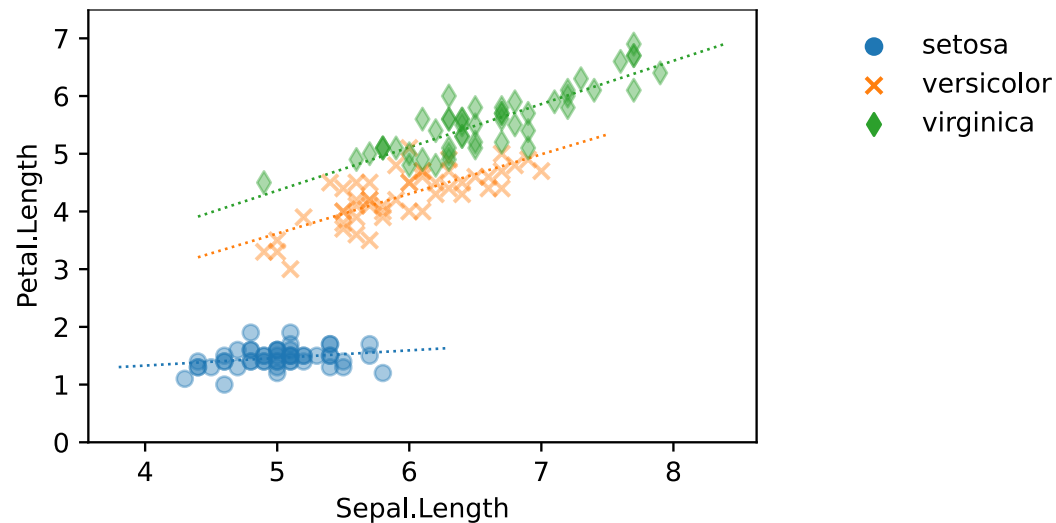
```
6 newSL = np.linspace(4.2, 8.2, 20)
7 predPL = model.predict(np.column_stack([newSL, newSL**2]))
```

This saves us from having to explicitly code up the prediction formula—so there's less chance we introduce bugs.

§2.2 Feature design

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

ONE-HOT CODING



$$PL \approx \alpha_{\text{species}} + \beta_{\text{species}} SL$$

This model has 6 unknown parameters:

$$\begin{matrix} \alpha_{\text{set}} & \alpha_{\text{vers}} & \alpha_{\text{virg}} \\ \beta_{\text{set}} & \beta_{\text{vers}} & \beta_{\text{virg}} \end{matrix}$$

So my model has 6 feature vectors.

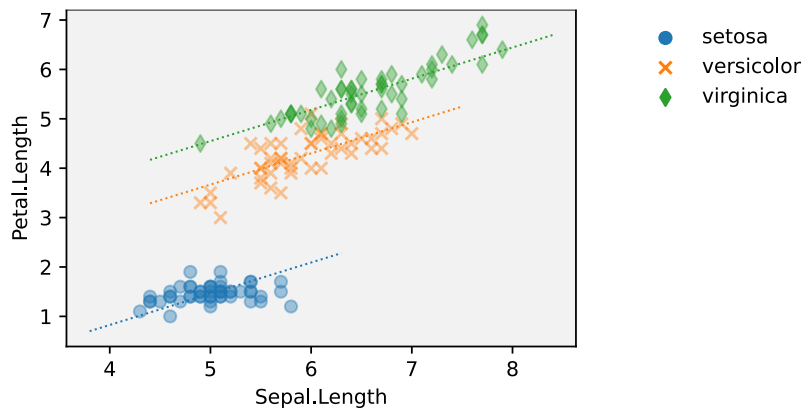
Row 1: setosa: $PL \approx \alpha_{\text{setosa}} + \beta_{\text{setosa}} SL$

$$\begin{matrix} * \text{ seto} \\ \text{virg} \\ \text{virg} \\ \text{seto} \\ \text{vers} \end{matrix} \begin{bmatrix} PL_1 \\ PL_2 \\ PL_3 \\ PL_4 \\ PL_5 \\ \vdots \end{bmatrix} \approx \alpha_{\text{seto}} \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \end{bmatrix} + \alpha_{\text{virg}} \begin{bmatrix} 0 \\ 1 \\ 1 \\ \vdots \end{bmatrix} + \alpha_{\text{vers}} \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \end{bmatrix} + \beta_{\text{seto}} \begin{bmatrix} SL_1 \\ 0 \\ 0 \\ \vdots \end{bmatrix} + \beta_{\text{virg}} \begin{bmatrix} 0 \\ SL_2 \\ SL_3 \\ \vdots \end{bmatrix} + \beta_{\text{vers}} \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \end{bmatrix}$$

$$\vec{PL} \approx \alpha_{\text{seto}} \vec{1}_{\text{spec}=\text{seto}} + \alpha_{\text{virg}} \vec{1}_{\text{spec}=\text{virg}} + \alpha_{\text{vers}} \vec{1}_{\text{spec}=\text{vers}} + \beta_{\text{seto}} (SL \times \vec{1}_{\text{spec}=\text{seto}}) + \dots$$

EXERCISE

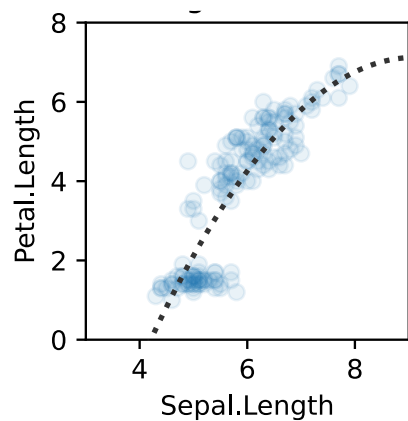
Fit the model with three parallel straight lines.



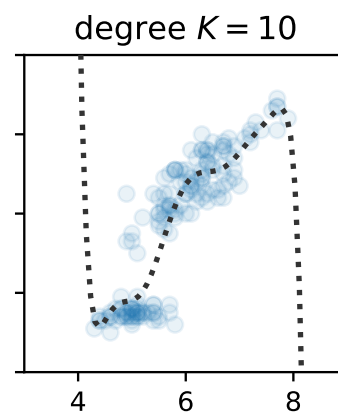
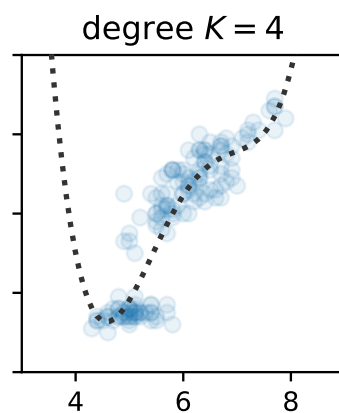
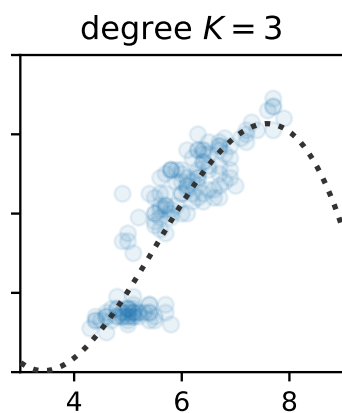
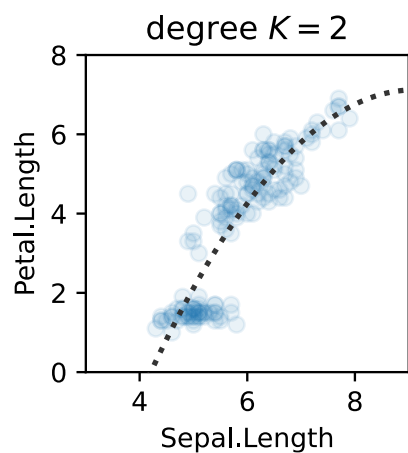
$$PL \approx \alpha_{\text{species}} + \beta SL$$

$$\approx \alpha_{\text{set}} 1_{\vec{\text{spec}} = \text{set}} + \dots + \beta \vec{SL}$$

NON-LINEAR RESPONSE



$$\text{Petal.Length} \approx \alpha + \beta \text{Sepal.Length} + \gamma (\text{Sepal.Length})^2$$



$$\text{Petal.Length} \approx \beta_0 + \sum_{k=1}^K \beta_k (\text{Sepal.Length})^k$$

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

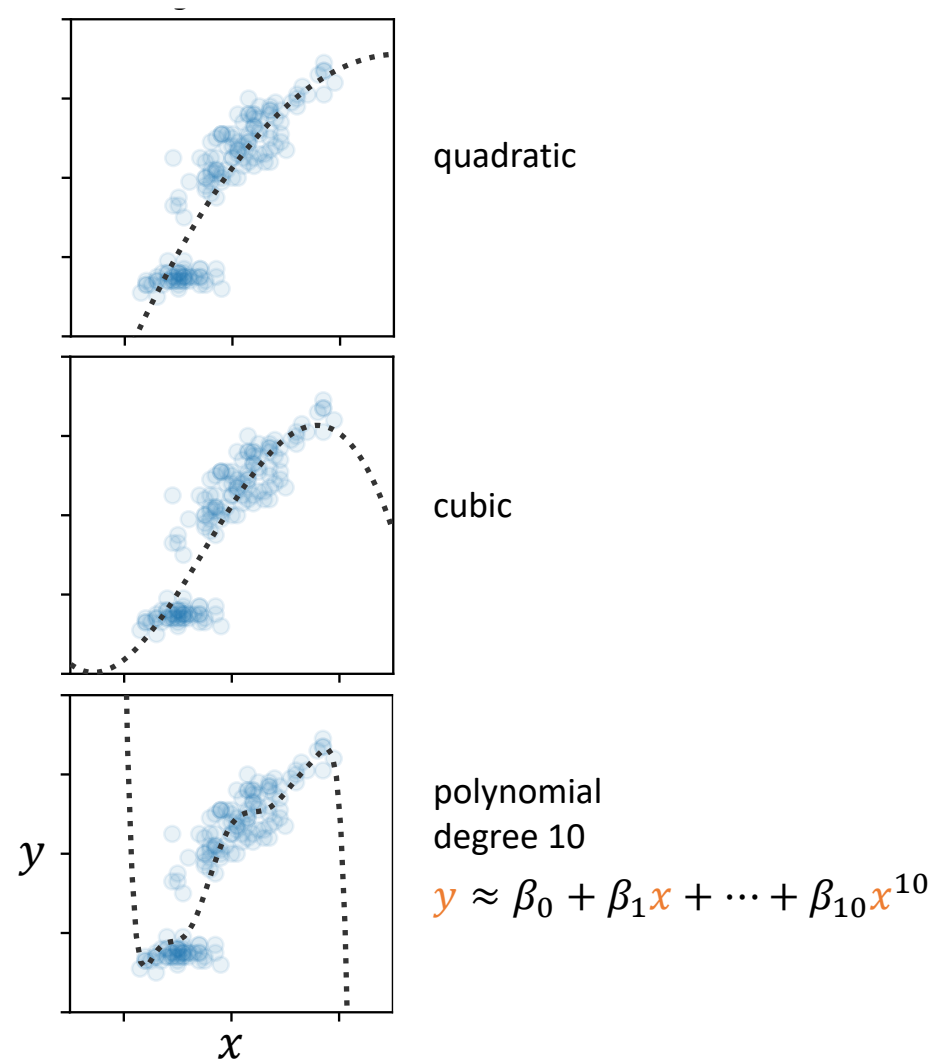
(seeing as the more features we add, the better we can fit the dataset)

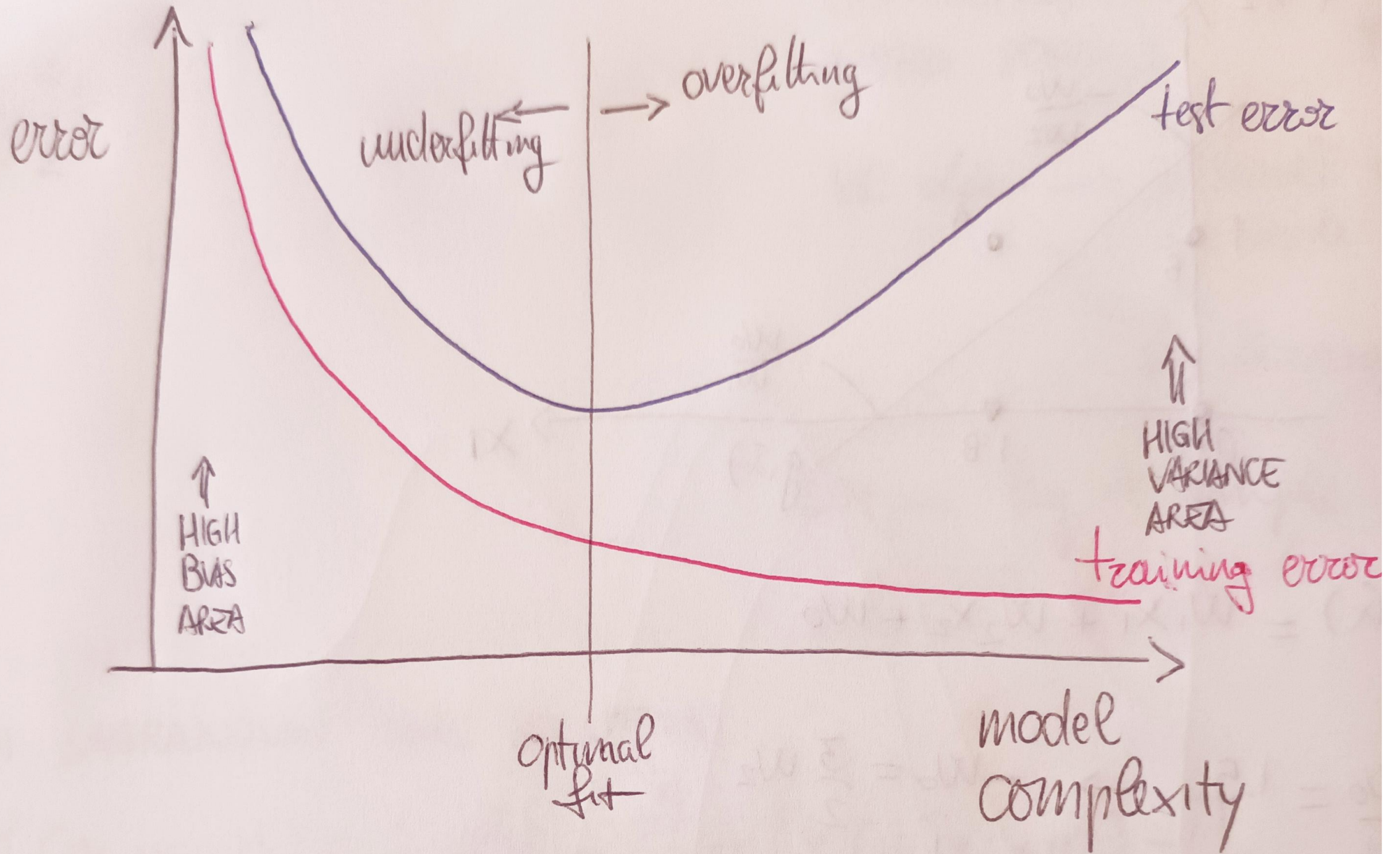
A. No. If we did, we'd *overfit*.

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

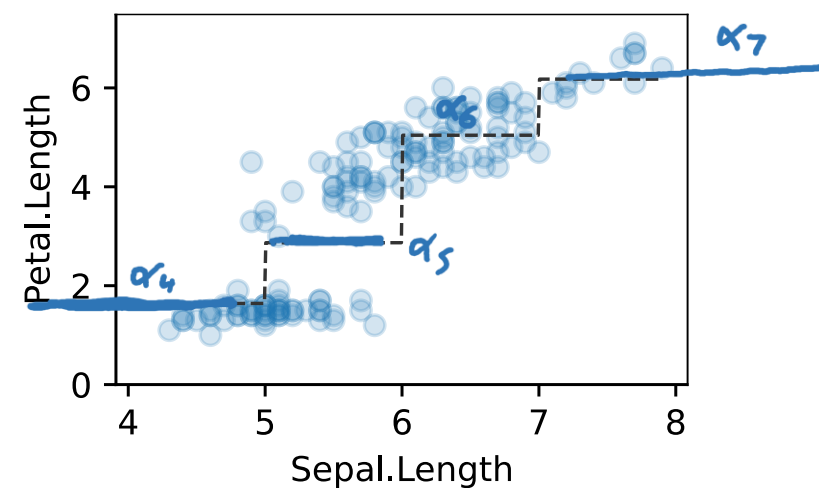
Or do model testing, e.g. evaluation on a holdout set.

[§5–7]





NON-LINEAR RESPONSE via one-hot coding



$$\vec{p}_L \approx \alpha_4 1_{\vec{S}L < 5} + \alpha_5 1_{[\vec{S}L] = 5} + \alpha_6 1_{[\vec{S}L] = 6} + \alpha_7 1_{\vec{S}L \geq 7}$$

e.g. for an observation with $SL = 5.3$ we predict α_5
 $SL = 2.7$ α_4

Memory lane





Listening modes:

- Rock
- Jazz
- Classical

How were listening modes calculated?

- Take a bunch of songs N_s of three different music genres
- In a linear model
 - the amplitudes for each frequency would be the feature vectors;
 - the amplitude of the bass at the speaker would be the response vector
 - (dimensionality = N_s)
- Calculate for each frequency the corresponding parameter

Scope of listening modes: improving audio quality

- By applying a frequency filter in the shape of the estimated distribution
 - Fostering some frequencies, reducing other frequencies

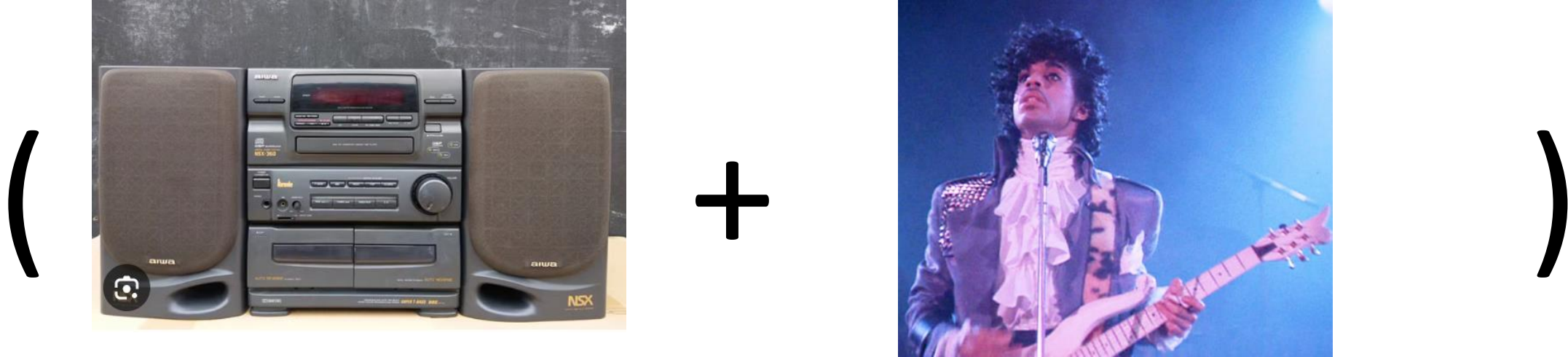


X

Accurate listening mode

=





X

Wrong listening mode

=

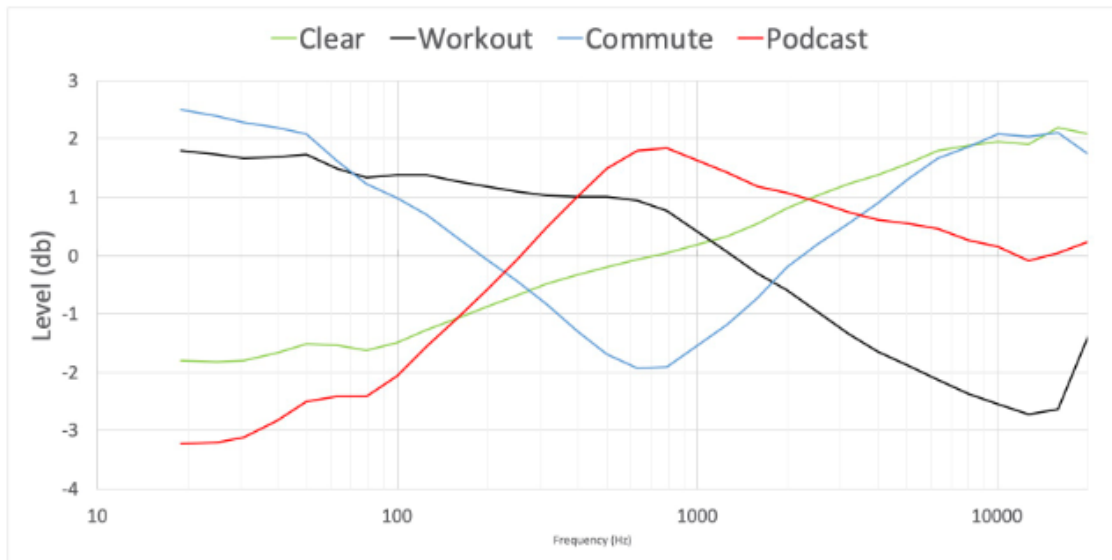


Listening modes, today

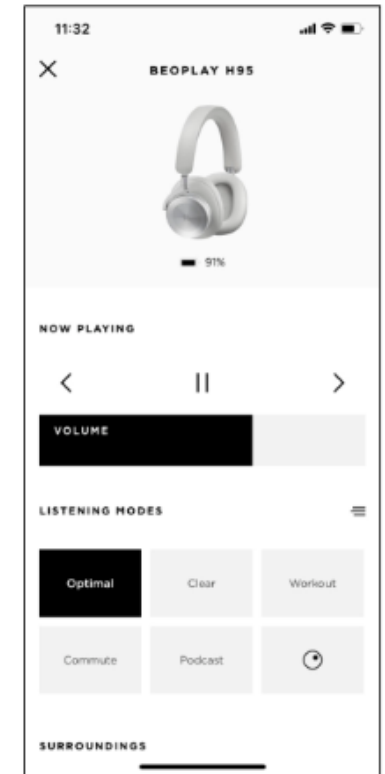


Listening Modes

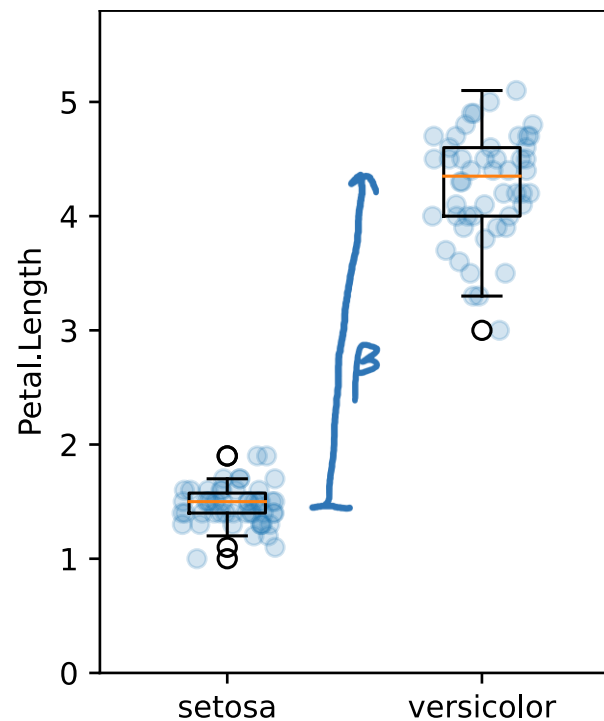
- These are the differences in frequency response relative to **Optimal Mode** (shown previous slide)
- **Clear Mode** is a upward sloped filter centered at 1kHz boosting treble above and cutting below
- **Workout Mode** does a similar EQ but sloped downwards
- **Commute Mode** boosts bass and treble
- **Podcast Mode** cuts bass below 250 Hz and boosts above that



Note: these curves are based on acoustical measurements of the left channel



COMPARING GROUPS



Measurements for condition A : $\mathbf{a} = [a_1, a_2, \dots, a_m]$

Measurements for condition B : $\mathbf{b} = [b_1, b_2, \dots, b_n]$

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

$$\vec{x} \approx \alpha_A \mathbf{1}_{\text{cond} = A} + \alpha_B \mathbf{1}_{\text{cond} = B}$$

Or

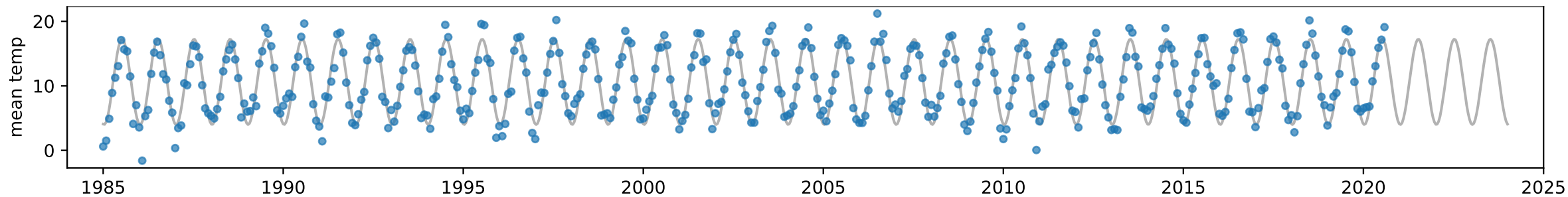
$$\vec{x} = \alpha + \beta \mathbf{1}_{\text{cond} = B}$$

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

β measures the difference between the two groups.

cond	x
A	a_1
A	\vdots
A	\vdots
A	a_m
B	b_1
B	b_2
\vdots	\vdots
B	b_n

PERIODIC PATTERN



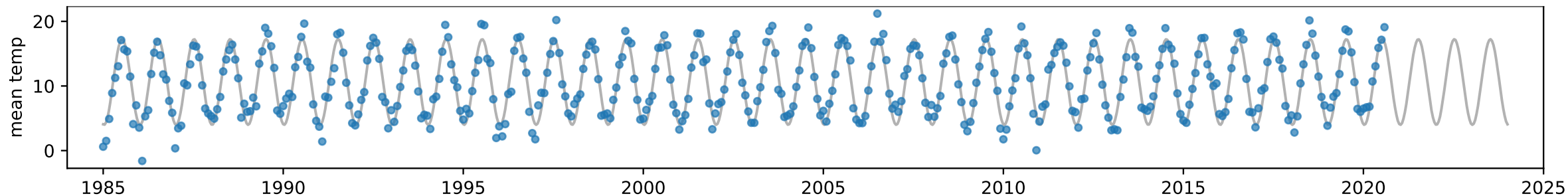
We'd like to fit the model:

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

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

$$\text{temp} \approx \alpha + \beta e + \gamma f$$

PERIODIC PATTERN



We'd like to fit the model:

$$\text{temp} \approx \alpha + \beta \sin(2\pi t + \phi)$$

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

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

$$= \alpha + \underbrace{\beta_1}_{\beta_1} \sin(2\pi t) + \underbrace{\beta_2}_{\beta_2} \cos(2\pi t)$$

$$= \alpha + \underbrace{\beta_1 e}_{\beta_1 e} + \underbrace{\beta_2 f}_{\beta_2 f}$$

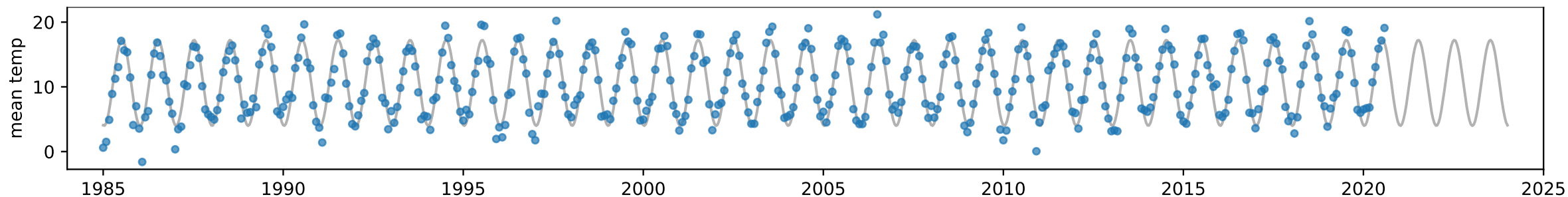
From school trigonometry,

$$\begin{aligned} \sin(A + B) \\ = \sin(A) \cos(B) + \cos(A) \sin(B) \end{aligned}$$

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

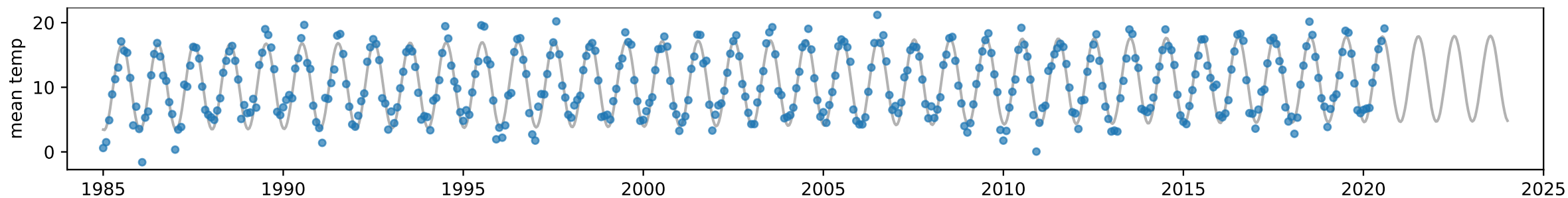
PERIODIC PATTERN

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



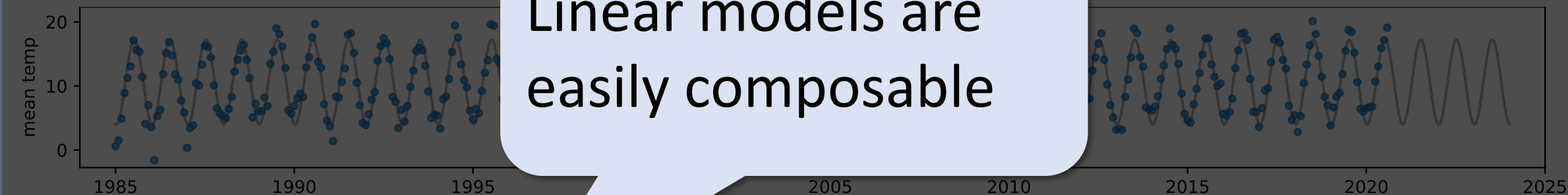
PERIODIC PATTERN + SECULAR TREND

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



PERIODIC PATTERN

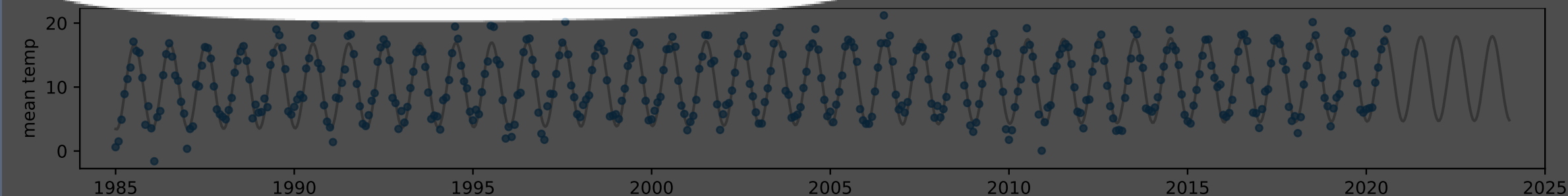
$$\text{temp} \approx \alpha + \beta_1 \sin(2\pi t) + \beta_2 \cos(2\pi t)$$



Linear models are easily composable

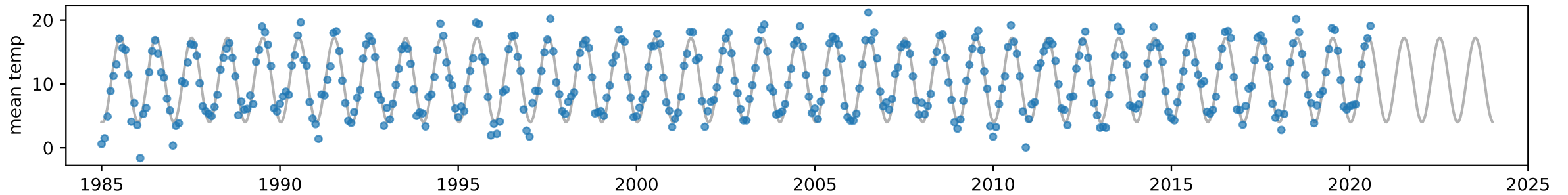
PERIODIC PATTERN + SECULAR TREND

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



With our periodic model ...

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

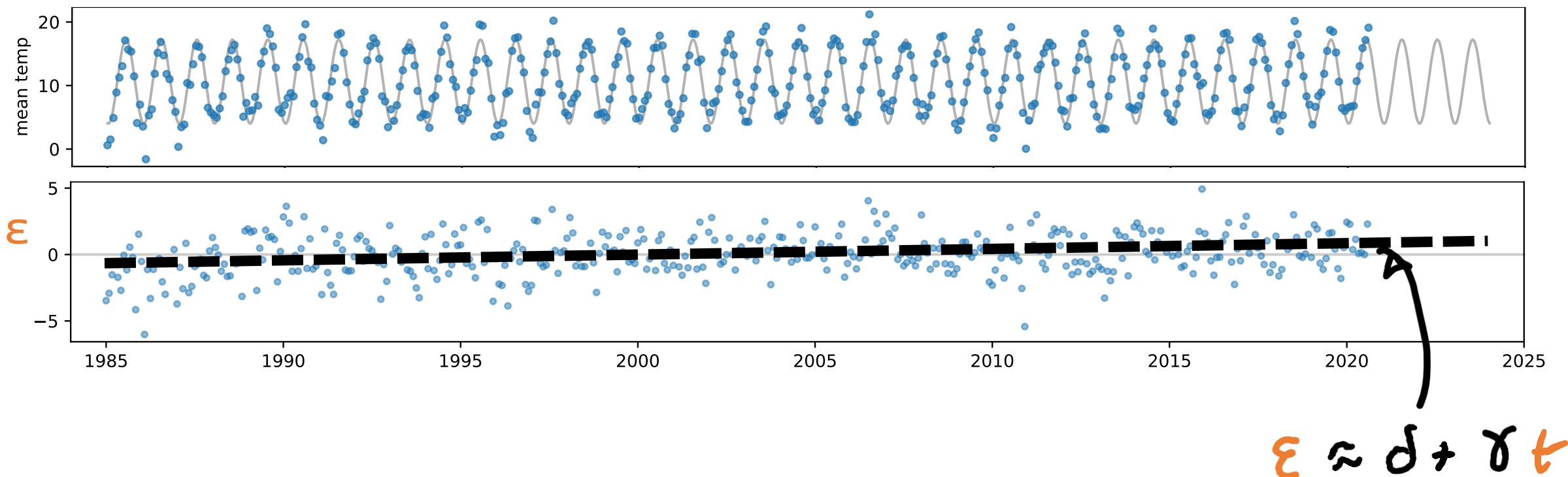


... how do we *discover* we should add a secular term $+\gamma t$?

If we hadn't thought to include climate change in our temperature model ...

$$\text{temp} \approx \alpha + \beta \sin(2\pi(\mathbf{t} + \phi))$$

$$\text{temp} = \alpha + \beta \sin(2\pi(\mathbf{t} + \phi)) + \varepsilon$$



This suggests a revised model ...

$$\text{temp} = \alpha' + \beta' \sin(2\pi(\mathbf{t} + \phi)) + \gamma \mathbf{t} + \varepsilon$$

Diagnosing a linear model

After fitting a model

```
model.fit(..., y)
```

1. Compute the residuals

```
 $\epsilon$  = y - model.predict()
```

2. Plot ϵ every way we can think of.
If there's a systematic pattern,
add a feature that describes it.