

# Automatic construction and description of nonparametric models

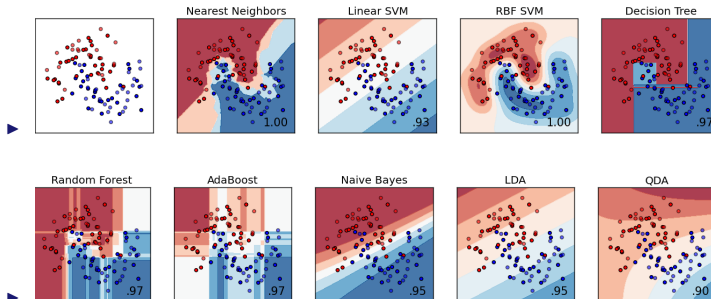


James Robert Lloyd, David Duvenaud, Roger Grosse,  
Josh Tenenbaum, Zoubin Ghahramani

December 2, 2013

# MOTIVATION

- ▶ Models today built by hand, or chosen from a fixed set.
  - ▶ Example: Scikit-learn



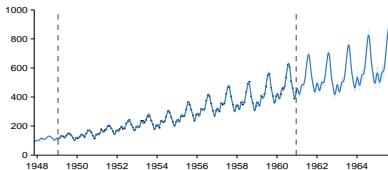
# PROBLEMS WITH THIS APPROACH

- ▶ Building by hand requires expertise, understanding of the dataset
- ▶ just being nonparametric isn't good enough
- ▶ can silently fail.

# MOTIVATION

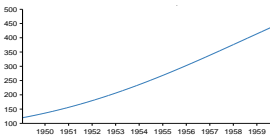
- ▶ Andrew Gelman asks: How would an AI do statistics?
- ▶ would need:
  - ▶ a language for describing arbitrarily complicated models
  - ▶ a way to search over those models
  - ▶ a way of checking model fit
- ▶ We built such a language over regression models, a procedure to search over them, and a method to describe in english language the properties of the resulting models.
  - ▶ Working on automatic model-checking.

# EXAMPLE

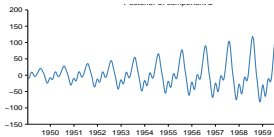


=

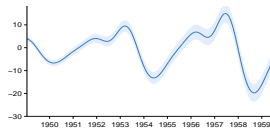
entire signal



+



+



A very smooth, monotonically increasing function

An approximately periodic function with a period of 1.0 years and with approximately linearly increasing amplitude

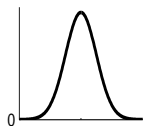
An exactly periodic function with a period of 4.3 years but with linearly increasing amplitude

# HOW TO BUILD A LANGUAGE OF MODELS?

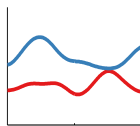
- ▶ We'll do this by defining a language on GP kernels
- ▶ Simple rules to combine them give diverse structures

# KERNEL DETERMINES STRUCTURE

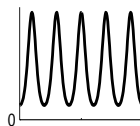
- ▶ Kernel determines almost all the properties of the prior.
- ▶ Many different kinds, with very different properties:



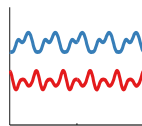
Squared-exp  
(SE)



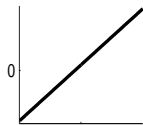
local variation



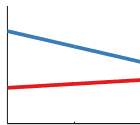
Periodic (PER)



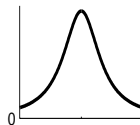
repeating  
structure



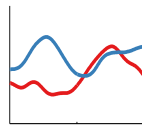
Linear (LIN)



linear func-  
tions



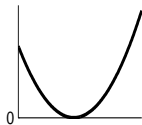
Rational-  
quadratic(RQ)



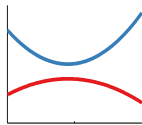
multi-scale  
variation

# KERNELS CAN BE COMPOSED

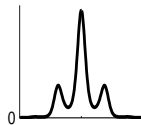
- Two main operations: adding, multiplying



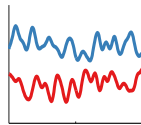
$\text{LIN} \times \text{LIN}$



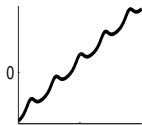
quadratic  
functions



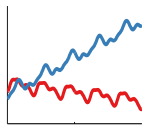
$\text{SE} \times \text{PER}$



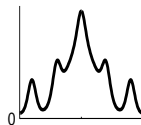
locally  
periodic



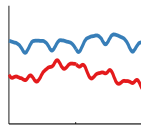
$\text{LIN} + \text{PER}$



periodic with  
trend



$\text{SE} + \text{PER}$

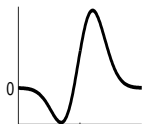


periodic with  
noise

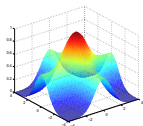


# KERNELS CAN BE COMPOSED

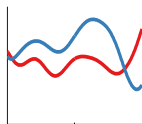
- Can be composed across multiple dimensions



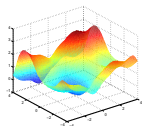
$\text{LIN} \times \text{SE}$



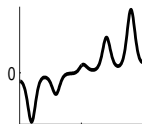
$\text{SE}_1 + \text{SE}_2$



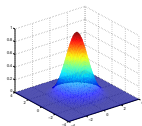
increasing  
variation



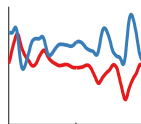
$f_1(x_1) + f_2(x_2)$



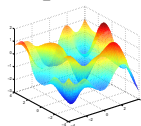
$\text{LIN} \times \text{PER}$



$\text{SE}_1 \times \text{SE}_2$



growing  
amplitude



$f(x_1, x_2)$

# SPECIAL CASES

Bayesian linear regression

LIN

Bayesian polynomial regression

LIN  $\times$  LIN  $\times$  ...

Generalized Fourier decomposition

PER + PER + ...

Generalized additive models

$\sum_{d=1}^D \text{SE}_d$

Automatic relevance determination

$\prod_{d=1}^D \text{SE}_d$

Linear trend with deviations

LIN + SE

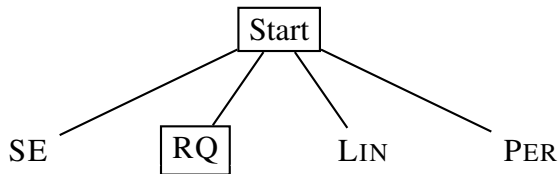
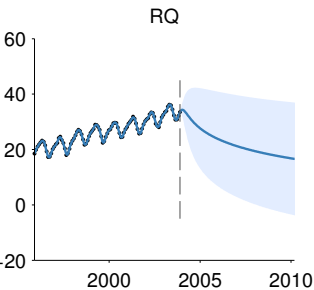
Linearly growing amplitude

LIN  $\times$  SE

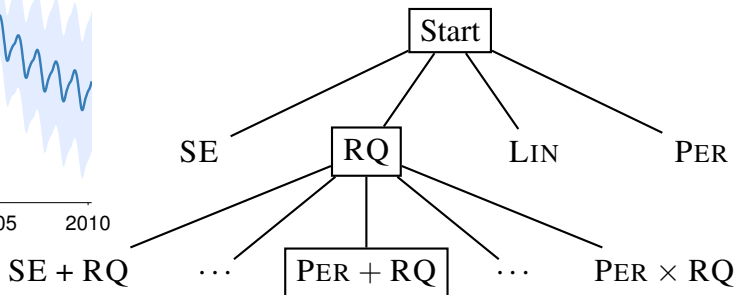
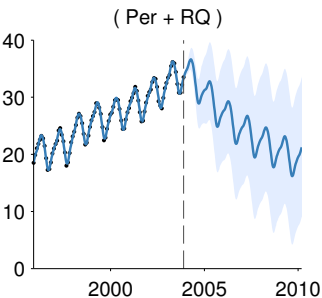
# COMPOSITIONAL STRUCTURE SEARCH

- ▶ Define grammar over kernels:
  - ▶  $K \rightarrow K + K$
  - ▶  $K \rightarrow K \times K$
  - ▶  $K \rightarrow \{\text{SE}, \text{LIN}, \text{PER}\}$
- ▶ Search the space of kernels greedily by applying production rules, checking model fit (approximate marginal likelihood).

# COMPOSITIONAL STRUCTURE SEARCH

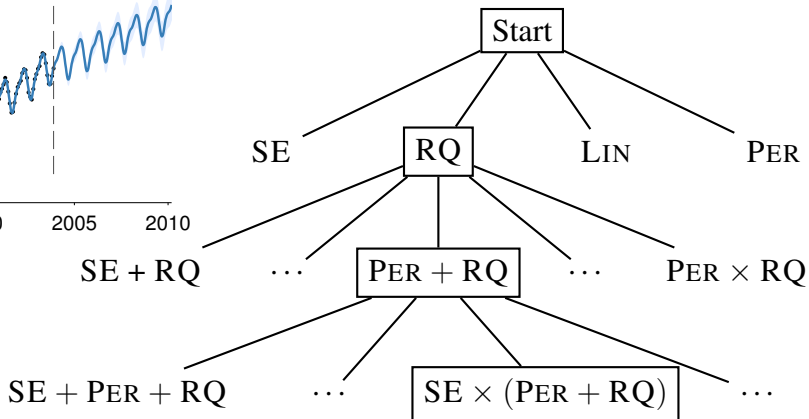
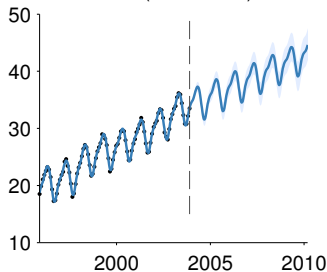


# COMPOSITIONAL STRUCTURE SEARCH

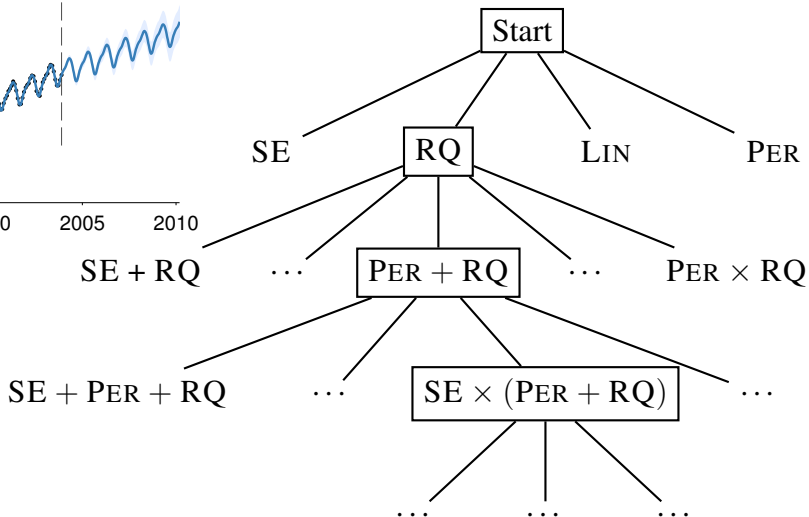
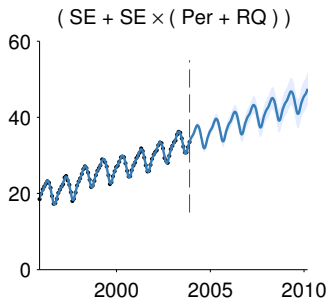


# COMPOSITIONAL STRUCTURE SEARCH

$SE \times (Per + RQ)$



# COMPOSITIONAL STRUCTURE SEARCH



# DISTRIBUTIVITY HELPS INTERPRETABILITY

We can write all kernels as sums of products of base kernels:

$$\text{SE} \times (\text{RQ} + \text{LIN}) = (\text{SE} \times \text{RQ}) + (\text{SE} \times \text{LIN}).$$

Sums of kernels are equivalent to sums of functions.

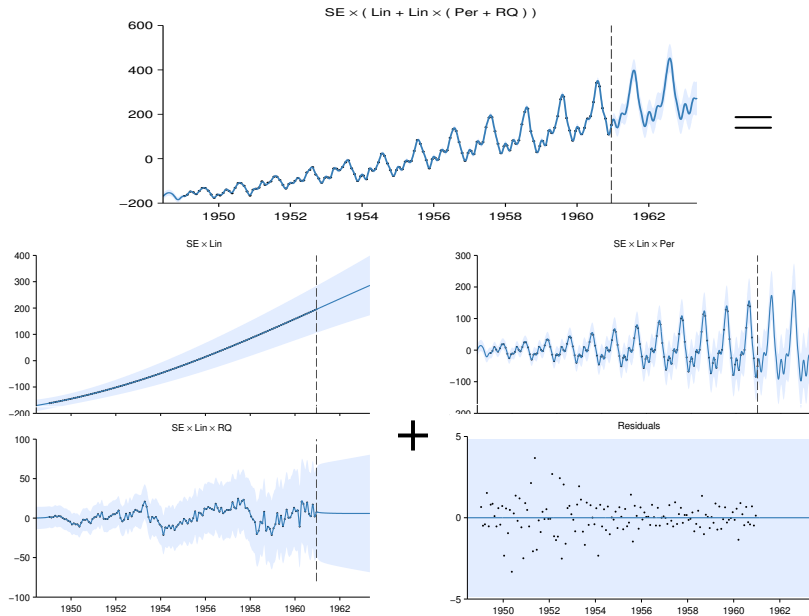
If  $f_1, f_2$  are independent, and  $f_1 \sim \mathcal{GP}(\mu_1, k_1), f_2 \sim \mathcal{GP}(\mu_2, k_2)$ .

Then it follows that

$$(f_1 + f_2) \sim \mathcal{GP}(\mu_1 + \mu_2, k_1 + k_2)$$



# EXAMPLE DECOMPOSITION: AIRLINE



# EXAMPLE KERNEL DESCRIPTIONS

Product of Kernels	Description
PER	An exactly periodic function
PER $\times$ SE	An approximately periodic function
PER $\times$ SE $\times$ LIN	An approximately periodic function with linearly varying amplitude
LIN	A linear function
LIN $\times$ LIN	A quadratic function
PER $\times$ LIN $\times$ LIN	An exactly periodic function with quadratically varying amplitude

# THIS ANALYSIS WAS AUTOMATICALLY GENERATED

The raw data and full model posterior with extrapolations are shown in figure 1.

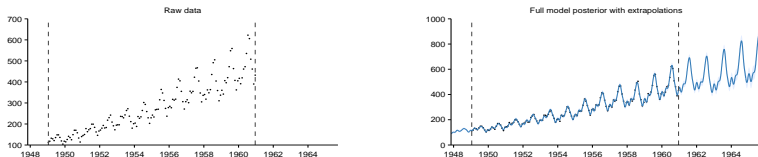


Figure 1: Raw data (left) and model posterior with extrapolation (right)

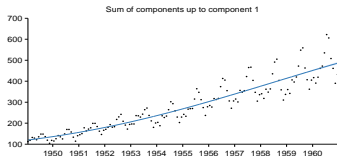
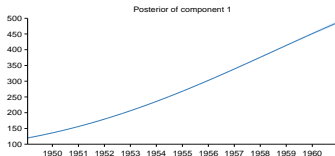
The structure search algorithm has identified four additive components in the data:

- A very smooth monotonically increasing function.
- An approximately periodic function with a period of 1.0 years and with approximately linearly increasing amplitude.
- An exactly periodic function with a period of 4.3 years but with linearly increasing amplitude.
- Uncorrelated noise with linearly increasing standard deviation.

# THIS ANALYSIS WAS AUTOMATICALLY GENERATED

## 2.1 Component 1

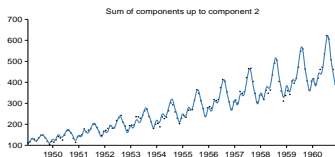
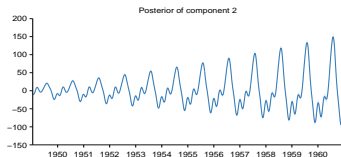
This component is a very smooth and monotonically increasing function.



# THIS ANALYSIS WAS AUTOMATICALLY GENERATED

## 2.2 Component 2

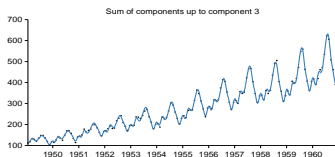
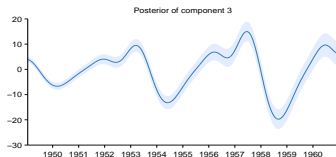
This component is approximately periodic with a period of 1.0 years and varying amplitude. Across periods the shape of this function varies very smoothly. The amplitude of the function increases approximately linearly. The shape of this function within each period has a typical lengthscale of 6.0 weeks.



# THIS ANALYSIS WAS AUTOMATICALLY GENERATED

## 2.3 Component 3

This component is exactly periodic with a period of 4.3 years but with varying amplitude. The amplitude of the function increases linearly. The shape of this function within each period has a typical lengthscale of 7.4 months.



# SUMMARY

- ▶ Compositions of kernels give a language of models.
- ▶ Can search over models automatically.
- ▶ Kernels modify prior in predictable ways, allowing automatic english description of models.

# SUMMARY

- ▶ Compositions of kernels give a language of models.
- ▶ Can search over models automatically.
- ▶ Kernels modify prior in predictable ways, allowing automatic english description of models.

Thanks!