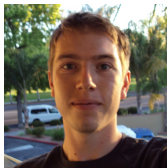# Automatic construction and description of nonparametric models
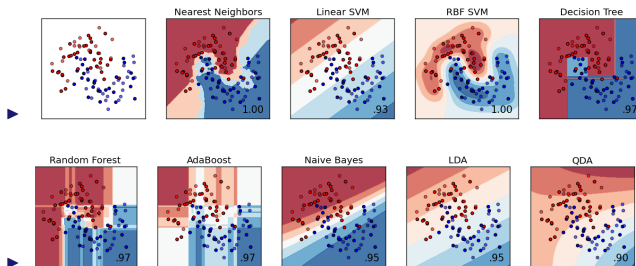
James Robert Lloyd, David Duvenaud, Roger Grosse,
Josh Tenenbaum, Zoubin Ghahramani

November 29, 2013

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



- Just being nonparametric sometimes isn't good enough
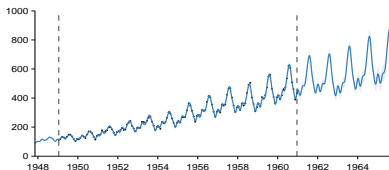- Building by hand requires expertise, understanding of the dataset.

# MOTIVATION

- Models today built by hand, or chosen from a fixed set.
  - Building by hand requires expertise, understanding of the dataset.
  - Follows cycle of: propose model, do inference, check model fit
    - Propose new model
    - Do inference
    - Check model fit
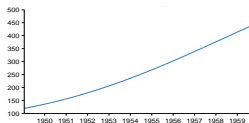  - for high-dimensional data, this can silently fail
- Andrew Gelman asks: How would an AI do statistics?
- It would need a language for describing arbitrarily complicated models, a way to search over those models, nad a way of checking model fit.

# MOTIVATION

- Andrew Gelman asks: How would an AI do statistics?
- It would need a language for describing arbitrarily complicated models, a way to search over those models, nad 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.
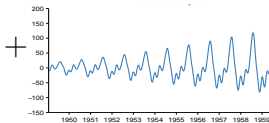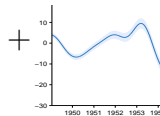
# EXAMPLE



entire signal

=


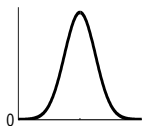
A very smooth monotonically increasing function

+

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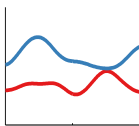
+

An exactly [...]tion with a [...] years but w[...] creasing an[...]
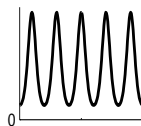
# KERNEL CHOICE IS IMPORTANT

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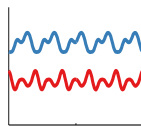- ▶ Many different kinds, with very different properties:
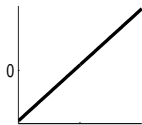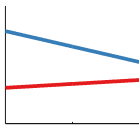


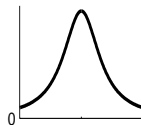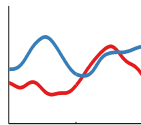Squared-exp (SE)

local variation
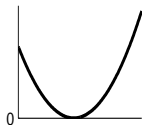
Periodic (PER)

repeating structure

Linear (LIN)

linear functions

Rational-quadratic(RQ)

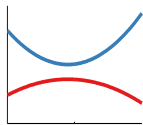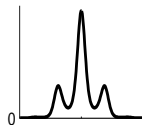multi-scale variation

▶ Two main operations: adding, multiplying



LIN × LIN

quadratic functions

SE × PER

locally periodic

LIN + PER

periodic with trend

SE + PER

periodic with noise

- Can be composed across multiple dimensions



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

increasing variation

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

growing amplitude

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

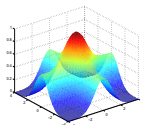$f_1(x_1) + f_2(x_2)$

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

$f(x_1, x_2)$

# SPECIAL CASES

| | |
|---|---|
| Bayesian linear regression | LIN |
| Bayesian polynomial regression | LIN $\times$ LIN $\times \ldots$ |
| Generalized Fourier decomposition | PER $+$ PER $+ \ldots$ |
| 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 |

- SE kernel $\rightarrow$ basic smoothing.
- Richer kernels means richer structure can be captured.

# KERNELS ARE HARD TO CHOOSE

- Given the diversity of priors available, how to choose one?
- Standard GP software packages include many base kernels and means to combine them, but *no default kernel*
- Software can't choose model for you, you're the expert (?)

# KERNELS ARE HARD TO CONSTRUCT

- Carl devotes 4 pages of his book to constructing a custom kernel for CO2 data
- requires specialized knowledge, trial and error, and a dataset small and low-dimensional enough that a human can interpret it.
- In practice, most users can't or won't make custom kernel, and SE kernel became *de facto* standard kernel through inertia.

# RECAP

- GP Regression is a powerful tool
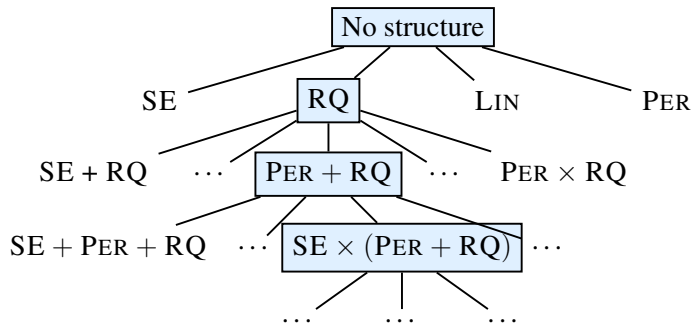- Kernel choice allows for rich structure to be captured - different kernels express very different model classes
- Composition generates a rich space of models
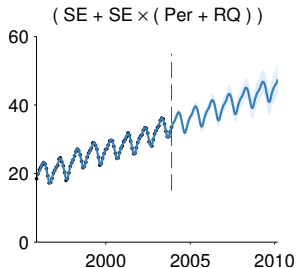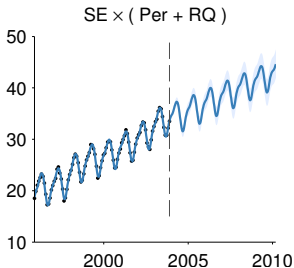- Hard & slow to search by hand
- Can kernel specification be automated?

# COMPOSITIONAL STRUCTURE SEARCH

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

Tree diagram of compositional kernel structure search:

- No structure
  - SE
  - RQ
    - SE + RQ
    - ⋯
    - PER + RQ
      - SE + PER + RQ
      - ⋯
      - SE × (PER + RQ)
        - ⋯   ⋯   ⋯
    - ⋯
    - PER × RQ
  - LIN
  - PER

# EXAMPLE SEARCH: MAUNA LUA $CO_2$

# EXAMPLE DECOMPOSITION: MAUNA LOA CO$_2$

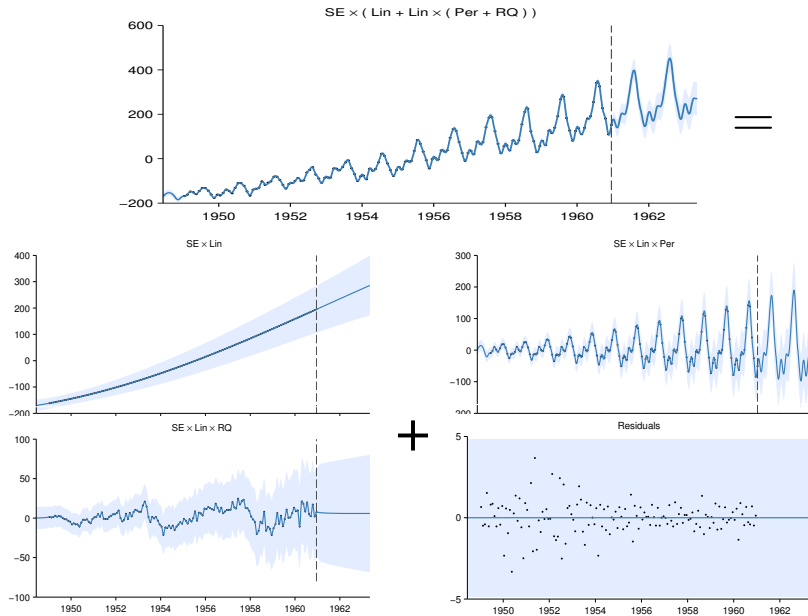Suppose functions $f_1, f_2$ are draw from independent GP priors, $f_1 \sim \mathcal{GP}(\mu_1, k_1), f_2 \sim \mathcal{GP}(\mu_2, k_2)$. Then it follows that
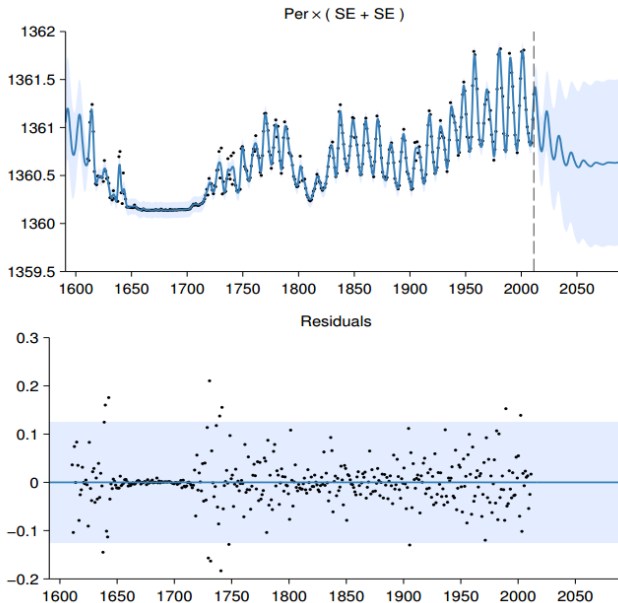
$$f := f_1 + f_2 \sim \mathcal{GP}(\mu_1 + \mu_2, k_1 + k_2)$$

Sum of kernels is equivalent to sum of functions. Distributivity means we can write compound kernels as sums of products of base kernels:

$$SE \times (RQ + LIN) = SE \times RQ + SE \times LIN.$$

Per × ( SE + SE )

Residuals

Can express change in covariance:



Periodic changing to SE

Can express change in covariance:



SE changing to linear

# SUMMARY

- Choosing form of kernel is currently done by hand.
- Compositions of kernels lead to more interesting priors on functions than typically considered.
- A simple grammar specifies all such compositions, and can be searched over automatically.
- Composite kernels lead to interpretable decompositions.

# CONCLUSIONS

- ▶ Model-building is currently done mostly by hand.
- ▶ Grammars over composite structures are a simple way to specify open-ended model classes.
- ▶ Composite structures often imply interpretable decompositions of the data.
- ▶ Searching over these model classes is a step towards automating statistical analysis.

# CONCLUSIONS

- ▶ Model-building is currently done mostly by hand.
- ▶ Grammars over composite structures are a simple way to specify open-ended model classes.
- ▶ Composite structures often imply interpretable decompositions of the data.
- ▶ Searching over these model classes is a step towards automating statistical analysis.

Thanks!