

How to build an automatic statistician



James Robert Lloyd¹, David Duvenaud¹, Roger Grosse²,



Joshua Tenenbaum², Zoubin Ghahramani¹

1: Department of Engineering, University of Cambridge, UK

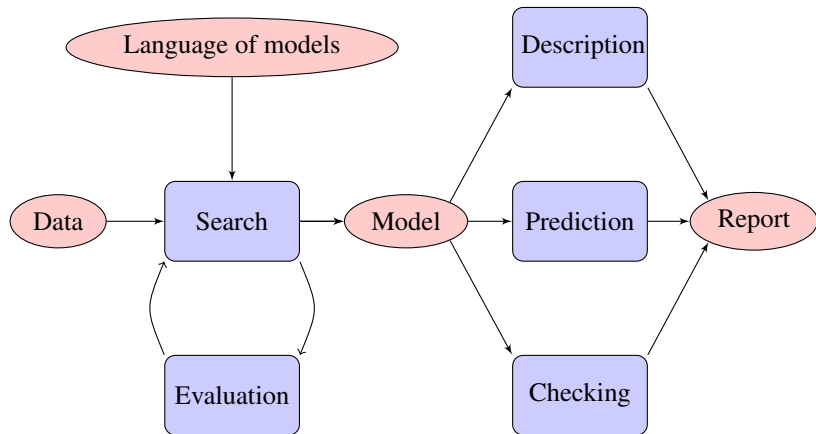
2: Massachusetts Institute of Technology, USA

November 19, 2014

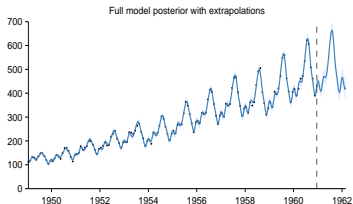
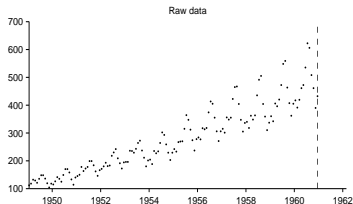
THERE IS A GROWING NEED FOR DATA ANALYSIS

- ▶ We live in an era of abundant data
- ▶ The McKinsey Global Institute claim
 - ▶ *“The United States alone faces a shortage of 140,000 to 190,000 people with analytical expertise and 1.5 million managers and analysts with the skills to understand and make decisions based on the analysis of big data.”*
- ▶ Diverse fields increasingly relying on expert statisticians, machine learning researchers and data scientists e.g.
 - ▶ Computational sciences (e.g. biology, astronomy, ...)
 - ▶ Online advertising
 - ▶ Quantitative finance
 - ▶ ...

WHAT WOULD AN AUTOMATIC STATISTICIAN DO?



PREVIEW: AN ENTIRELY AUTOMATIC ANALYSIS



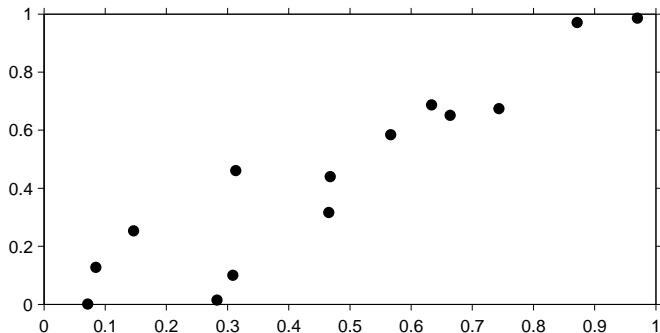
Four additive components have been identified in the data

- ▶ A linearly increasing function.
- ▶ An approximately periodic function with a period of 1.0 years and with linearly increasing amplitude.
- ▶ A smooth function.
- ▶ Uncorrelated noise with linearly increasing standard deviation.

AGENDA

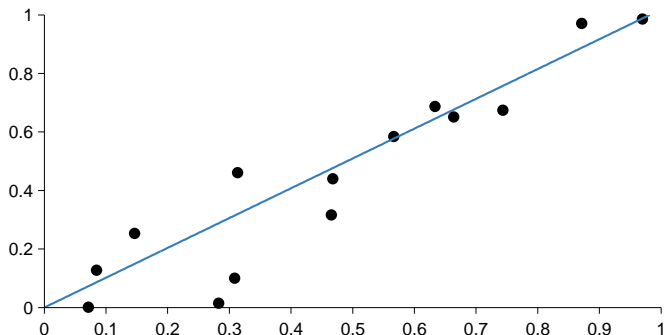
- ▶ Introduction to Bayesian statistics
- ▶ Introduction to Gaussian processes via linear regression
- ▶ Description of automatic statistician
- ▶ Examples of automatically generated output

HOW TO DO LINEAR REGRESSION



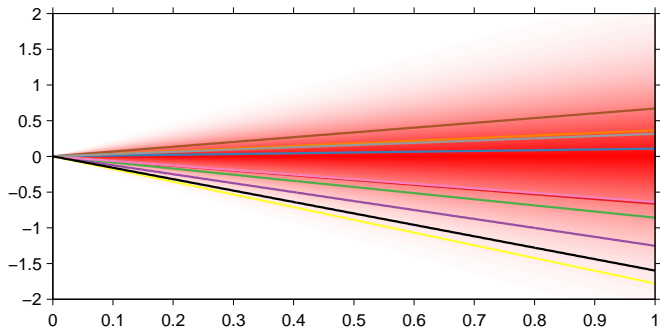
- ▶ Linear regression starts by assuming a model of the form $y_i = mx_i + \varepsilon_i$ where x and y are inputs and outputs respectively and ε are errors or noise

HOW TO DO LINEAR REGRESSION



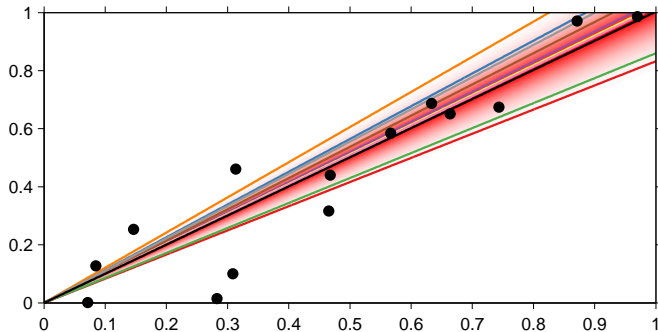
- ▶ Common approach: Estimate m by least squares

HOW TO DO BAYESIAN LINEAR REGRESSION



- ▶ Bayesian approach:
 - ▶ Specify beliefs about m using probability distributions

HOW TO DO BAYESIAN LINEAR REGRESSION



- ▶ Bayesian approach:
 - ▶ Specify prior beliefs about m using probability distributions
 - ▶ Follow rules of probability to update beliefs rationally after observing data

- ▶ Suppose we wish to fit a model M with parameters θ to data D
 - ▶ M defined by $p(D | \theta, M)$ - the likelihood

BAYESIAN STATISTICS

- ▶ Suppose we wish to fit a model M with parameters θ to data D
 - ▶ M defined by $p(D | \theta, M)$ - the likelihood
- ▶ We represent our uncertainty about θ with a probability distribution
 - ▶ $p(\theta | M)$ - the prior

BAYESIAN STATISTICS

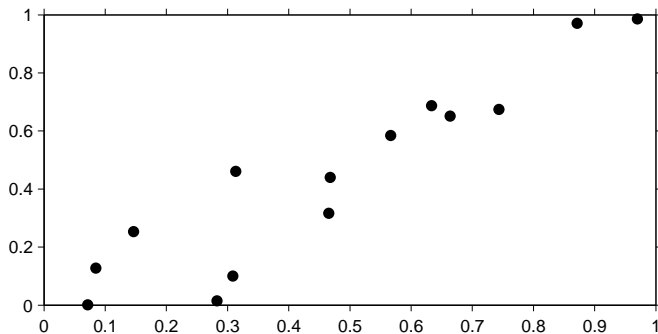
- ▶ Suppose we wish to fit a model M with parameters θ to data D
 - ▶ M defined by $p(D | \theta, M)$ - the likelihood
- ▶ We represent our uncertainty about θ with a probability distribution
 - ▶ $p(\theta | M)$ - the prior
- ▶ We now derive $p(\theta | D, M)$ - the posterior

BAYES RULE

- ▶ Suppose we wish to fit a model M with parameters θ to data D
 - ▶ M defined by $p(D | \theta, M)$ - the likelihood
- ▶ We represent our uncertainty about θ with a probability distribution
 - ▶ $p(\theta | M)$ - the prior
- ▶ We then derive $p(\theta | D, M)$ - the posterior

$$\underbrace{p(\theta | D, M)}_{\text{posterior}} = \frac{\overbrace{p(D | \theta, M)}^{\text{likelihood}} \overbrace{p(\theta | M)}^{\text{prior}}}{\underbrace{\int p(D | \theta, M) p(\theta | M) d\theta}_{\text{marginal likelihood}}}$$

BAYESIAN LINEAR REGRESSION



- ▶ Linear regression starts by assuming a model of the form $y_i = mx_i + \varepsilon_i$ where x and y are inputs and outputs respectively and ε are errors or noise

BAYESIAN LINEAR REGRESSION

- ▶ Linear regression starts by assuming a model of the form $y_i = mx_i + \varepsilon_i$ where x and y are inputs and outputs respectively and ε are errors or noise
- ▶ Represent prior beliefs about m using a probability distribution

BAYESIAN LINEAR REGRESSION

- ▶ Linear regression starts by assuming a model of the form $y_i = mx_i + \varepsilon_i$ where x and y are inputs and outputs respectively and ε are errors or noise
- ▶ Represent prior beliefs about m using a probability distribution
 - ▶ e.g. $m \sim \text{Normal}(\text{mean} = 0, \text{variance} = 1)$

BAYESIAN LINEAR REGRESSION

- ▶ Linear regression starts by assuming a model of the form $y_i = mx_i + \varepsilon_i$ where x and y are inputs and outputs respectively and ε are errors or noise
- ▶ Represent prior beliefs about m using a probability distribution
 - ▶ e.g. $m \sim \mathcal{N}(0, 1)$

BAYESIAN LINEAR REGRESSION

- ▶ Linear regression starts by assuming a model of the form $y_i = mx_i + \varepsilon_i$ where x and y are inputs and outputs respectively and ε are errors or noise
- ▶ Represent prior beliefs about m using a probability distribution
 - ▶ e.g. $m \sim \mathcal{N}(0, 1)$
- ▶ Define likelihood by representing prior beliefs about ε using a probability distribution
 - ▶ e.g. $\varepsilon_i \sim \mathcal{N}(0, \sigma_\varepsilon^2)$ independently for all i

BAYESIAN LINEAR REGRESSION

- ▶ Linear regression starts by assuming a model of the form $y_i = mx_i + \varepsilon_i$ where x and y are inputs and outputs respectively and ε are errors or noise
- ▶ Represent prior beliefs about m using a probability distribution
 - ▶ e.g. $m \sim \mathcal{N}(0, 1)$
- ▶ Define likelihood by representing prior beliefs about ε using a probability distribution
 - ▶ e.g. $\varepsilon_i \sim \mathcal{N}(0, \sigma_\varepsilon^2)$ independently for all i
- ▶ Apply Bayes rule

BAYES RULE APPLIED TO LINEAR REGRESSION

$$p(m | y, x) = \frac{p(y | m, x) p(m)}{\int p(y | m, x) p(m) \, dm}$$

BAYES RULE APPLIED TO LINEAR REGRESSION

$$\begin{aligned} p(m | y, x) &= \frac{p(y | m, x) p(m)}{\int p(y | m, x) p(m) \, dm} \\ &\propto p(y | m, x) p(m) \end{aligned}$$

BAYES RULE APPLIED TO LINEAR REGRESSION

$$\begin{aligned} p(m | y, x) &= \frac{p(y | m, x) p(m)}{\int p(y | m, x) p(m) \, dm} \\ &\propto p(y | m, x) p(m) \\ &\propto \left(\prod_i \frac{1}{2\pi\sigma_\varepsilon} e^{-(y_i - mx_i)^2 / (2\sigma_\varepsilon^2)} \right) \frac{1}{2\pi} e^{-m^2/2} \end{aligned}$$

BAYES RULE APPLIED TO LINEAR REGRESSION

$$\begin{aligned} p(m | y, x) &= \frac{p(y | m, x) p(m)}{\int p(y | m, x) p(m) \, dm} \\ &\propto p(y | m, x) p(m) \\ &\propto \left(\prod_i \frac{1}{2\pi\sigma_\varepsilon} e^{-(y_i - mx_i)^2 / (2\sigma_\varepsilon^2)} \right) \frac{1}{2\pi} e^{-m^2/2} \\ &\propto e^{-\sum_i (y_i - mx_i)^2 / (2\sigma_\varepsilon^2) - m^2/2} \end{aligned}$$

BAYES RULE APPLIED TO LINEAR REGRESSION

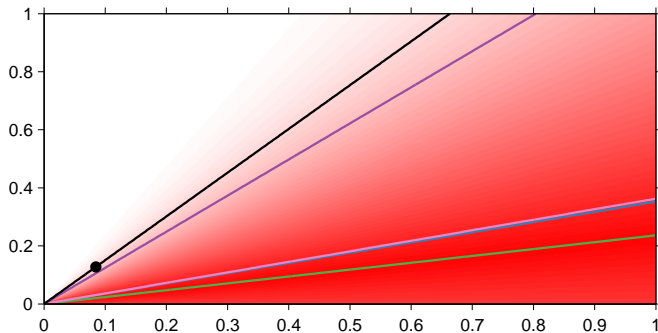
$$\begin{aligned} p(m | y, x) &= \frac{p(y | m, x) p(m)}{\int p(y | m, x) p(m) \, dm} \\ &\propto p(y | m, x) p(m) \\ &\propto \left(\prod_i \frac{1}{2\pi\sigma_\varepsilon} e^{-(y_i - mx_i)^2 / (2\sigma_\varepsilon^2)} \right) \frac{1}{2\pi} e^{-m^2/2} \\ &\propto e^{-\sum_i (y_i - mx_i)^2 / (2\sigma_\varepsilon^2) - m^2/2} \\ &\propto e^{-\frac{\sum_i x_i^2 + \sigma_\varepsilon^2}{2\sigma_\varepsilon^2} \left(m - \frac{\sum_i x_i y_i}{\sum_i x_i^2 + \sigma_\varepsilon^2} \right)^2} \end{aligned}$$

BAYES RULE APPLIED TO LINEAR REGRESSION

$$\begin{aligned} p(m | y, x) &= \frac{p(y | m, x) p(m)}{\int p(y | m, x) p(m) dm} \\ &\propto p(y | m, x) p(m) \\ &\propto \left(\prod_i \frac{1}{2\pi\sigma_\epsilon} e^{-(y_i - mx_i)^2 / (2\sigma_\epsilon^2)} \right) \frac{1}{2\pi} e^{-m^2/2} \\ &\propto e^{-\sum_i (y_i - mx_i)^2 / (2\sigma_\epsilon^2) - m^2/2} \\ &\propto e^{-\frac{\sum_i x_i^2 + \sigma_\epsilon^2}{2\sigma_\epsilon^2} \left(m - \frac{\sum_i x_i y_i}{\sum_i x_i^2 + \sigma_\epsilon^2} \right)^2} \end{aligned}$$

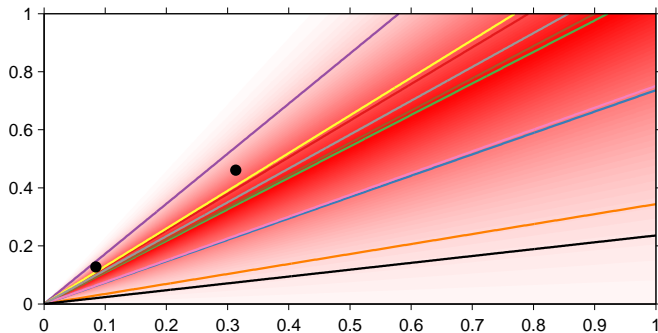
This is a normal distribution: $m | y, x \sim \mathcal{N}\left(\frac{\sum_i x_i y_i}{\sum_i x_i^2 + \sigma_\epsilon^2}, \frac{\sigma_\epsilon^2}{\sum_i x_i^2 + \sigma_\epsilon^2}\right)$

BAYESIAN LINEAR REGRESSION



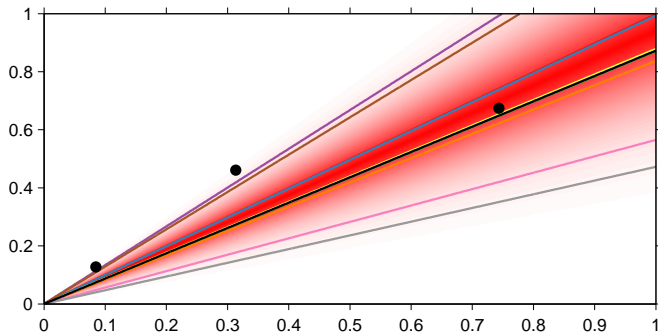
$$m | y, x \sim \mathcal{N}\left(\frac{\sum_i x_i y_i}{\sum_i x_i^2 + \sigma_\epsilon^2}, \frac{\sigma_\epsilon^2}{\sum_i x_i^2 + \sigma_\epsilon^2}\right)$$

BAYESIAN LINEAR REGRESSION



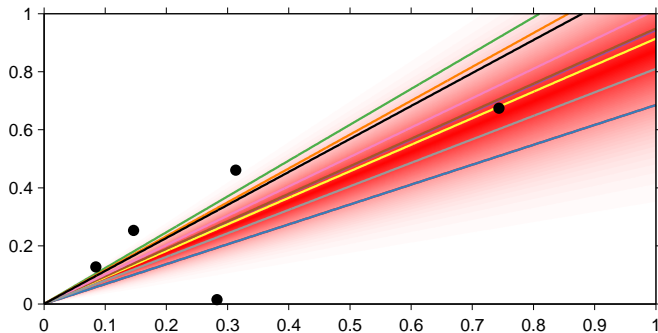
$$m | y, x \sim \mathcal{N}\left(\frac{\sum_i x_i y_i}{\sum_i x_i^2 + \sigma_\epsilon^2}, \frac{\sigma_\epsilon^2}{\sum_i x_i^2 + \sigma_\epsilon^2}\right)$$

BAYESIAN LINEAR REGRESSION



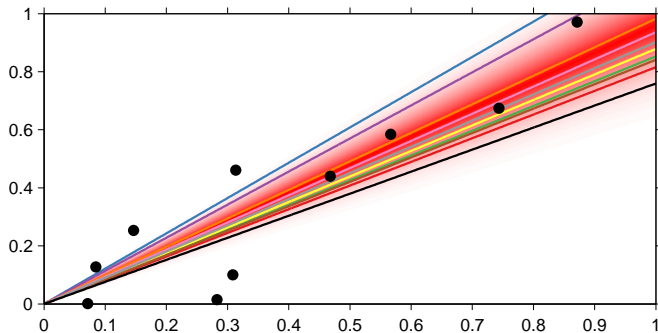
$$m | y, x \sim \mathcal{N}\left(\frac{\sum_i x_i y_i}{\sum_i x_i^2 + \sigma_\epsilon^2}, \frac{\sigma_\epsilon^2}{\sum_i x_i^2 + \sigma_\epsilon^2}\right)$$

BAYESIAN LINEAR REGRESSION



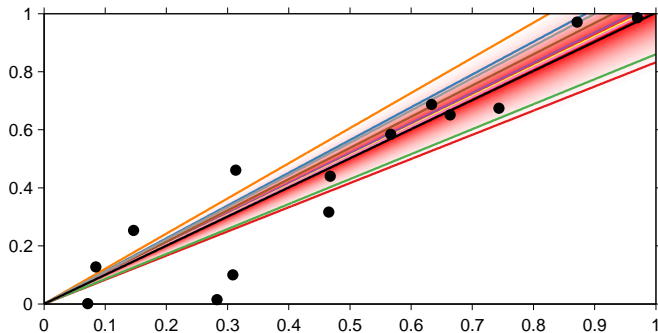
$$m | y, x \sim \mathcal{N}\left(\frac{\sum_i x_i y_i}{\sum_i x_i^2 + \sigma_\epsilon^2}, \frac{\sigma_\epsilon^2}{\sum_i x_i^2 + \sigma_\epsilon^2}\right)$$

BAYESIAN LINEAR REGRESSION



$$m | y, x \sim \mathcal{N}\left(\frac{\sum_i x_i y_i}{\sum_i x_i^2 + \sigma_\epsilon^2}, \frac{\sigma_\epsilon^2}{\sum_i x_i^2 + \sigma_\epsilon^2}\right)$$

BAYESIAN LINEAR REGRESSION



$$m | y, x \sim \mathcal{N}\left(\frac{\sum_i x_i y_i}{\sum_i x_i^2 + \sigma_\epsilon^2}, \frac{\sigma_\epsilon^2}{\sum_i x_i^2 + \sigma_\epsilon^2}\right)$$

BAYESIAN LINEAR REGRESSION REWRITTEN

- ▶ We defined a Bayesian linear regression model by specifying priors on m and ε_i

BAYESIAN LINEAR REGRESSION REWRITTEN

- ▶ We defined a Bayesian linear regression model by specifying priors on m and ε_i
 - ▶ $m \sim \mathcal{N}(0, 1)$, $\varepsilon_i \sim_{\text{iid}} \mathcal{N}(0, \sigma_\varepsilon^2)$

BAYESIAN LINEAR REGRESSION REWRITTEN

- ▶ We defined a Bayesian linear regression model by specifying priors on m and ε_i
 - ▶ $m \sim \mathcal{N}(0, 1)$, $\varepsilon_i \sim_{\text{iid}} \mathcal{N}(0, \sigma_\varepsilon^2)$
 - ▶ $y_i | m, \varepsilon_i = mx_i + \varepsilon_i$

BAYESIAN LINEAR REGRESSION REWRITTEN

- ▶ We defined a Bayesian linear regression model by specifying priors on m and ε_i
 - ▶ $m \sim \mathcal{N}(0, 1)$, $\varepsilon_i \sim_{\text{iid}} \mathcal{N}(0, \sigma_\varepsilon^2)$
 - ▶ $y_i | m, \varepsilon_i = mx_i + \varepsilon_i$
- ▶ This implicitly defined a joint prior on $\{y_i : i = 1, \dots, n\}$

BAYESIAN LINEAR REGRESSION REWRITTEN

- ▶ We defined a Bayesian linear regression model by specifying priors on m and ε_i
 - ▶ $m \sim \mathcal{N}(0, 1)$, $\varepsilon_i \sim_{\text{iid}} \mathcal{N}(0, \sigma_\varepsilon^2)$
 - ▶ $y_i | m, \varepsilon_i = mx_i + \varepsilon_i$
- ▶ This implicitly defined a joint prior on $\{y_i : i = 1, \dots, n\}$
 - ▶ $y_i \sim \mathcal{N}(0, x_i^2 + \sigma_\varepsilon^2)$ (sum of two normals)

BAYESIAN LINEAR REGRESSION REWRITTEN

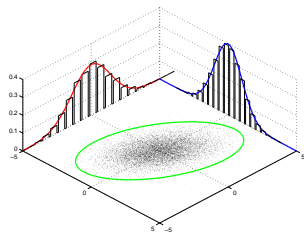
- ▶ We defined a Bayesian linear regression model by specifying priors on m and ε_i
 - ▶ $m \sim \mathcal{N}(0, 1)$, $\varepsilon_i \sim_{\text{iid}} \mathcal{N}(0, \sigma_\varepsilon^2)$
 - ▶ $y_i \mid m, \varepsilon_i = mx_i + \varepsilon_i$
- ▶ This implicitly defined a joint prior on $\{y_i : i = 1, \dots, n\}$
 - ▶ $y_i \sim \mathcal{N}(0, x_i^2 + \sigma_\varepsilon^2)$ (sum of two normals)
 - ▶ $\text{cov}(y_i, y_j) = x_i x_j \quad \forall i \neq j$

BAYESIAN LINEAR REGRESSION REWRITTEN

- ▶ We defined a Bayesian linear regression model by specifying priors on m and ε_i
 - ▶ $m \sim \mathcal{N}(0, 1)$, $\varepsilon_i \sim_{\text{iid}} \mathcal{N}(0, \sigma_\varepsilon^2)$
 - ▶ $y_i | m, \varepsilon_i = mx_i + \varepsilon_i$
- ▶ This implicitly defined a joint prior on $\{y_i : i = 1, \dots, n\}$
 - ▶ $y_i \sim \mathcal{N}(0, x_i^2 + \sigma_\varepsilon^2)$ (sum of two normals)
 - ▶ $\text{cov}(y_i, y_j) = x_i x_j \quad \forall i \neq j$

$\{y_i : i = 1, \dots, n\}$ has a multivariate normal distribution

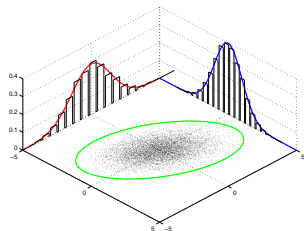
- ▶ $y \sim \mathcal{N}(\mathbf{0}, \mathbf{K})$
- ▶ where $k_{ij} = x_i x_j + \delta_{ij} \sigma_\varepsilon^2$



BAYESIAN LINEAR REGRESSION REWRITTEN

$\{y_i : i = 1, \dots, n\}$ has a multivariate normal distribution

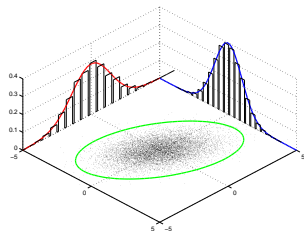
- ▶ $y \sim \mathcal{N}(\mathbf{0}, \mathbf{K})$
- ▶ where $k_{ij} = x_i x_j + \delta_{ij} \sigma_\epsilon^2$



BAYESIAN LINEAR REGRESSION REWRITTEN

$\{y_i : i = 1, \dots, n\}$ has a multivariate normal distribution

- ▶ $y \sim \mathcal{N}(\mathbf{0}, \mathbf{K})$
- ▶ where $k_{ij} = x_i x_j + \delta_{ij} \sigma_\epsilon^2$



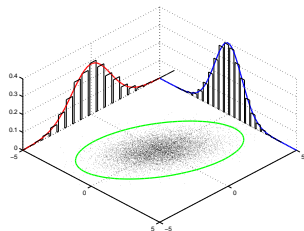
The covariance can be written as a function of pairs of x_i

$$k_{ij} = k(x_i, x_j) = x_i x_j + \delta_{x_i=x_j} \sigma_\epsilon^2$$

BAYESIAN LINEAR REGRESSION REWRITTEN

$\{y_i : i = 1, \dots, n\}$ has a multivariate normal distribution

- ▶ $y \sim \mathcal{N}(\mathbf{0}, \mathbf{K})$
- ▶ where $k_{ij} = x_i x_j + \delta_{ij} \sigma_\epsilon^2$



The covariance can be written as a function of pairs of x_i

$$k_{ij} = k(x_i, x_j) = x_i x_j + \delta_{x_i=x_j} \sigma_\epsilon^2$$

Do we define a valid probability distribution if $\{x_i\} = \mathbb{R}$?

GAUSSIAN PROCESSES

- ▶ A Gaussian process is collection of random variables, any finite number of which have a joint Gaussian distribution

GAUSSIAN PROCESSES

- ▶ A Gaussian process is collection of random variables, any finite number of which have a joint Gaussian distribution
- ▶ We can write this collection of random variables as $\{f(x) : x \in \mathcal{X}\}$ i.e. a function f evaluated at inputs x

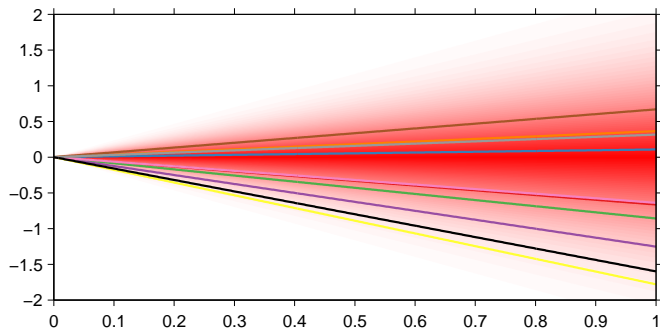
GAUSSIAN PROCESSES

- ▶ A Gaussian process is collection of random variables, any finite number of which have a joint Gaussian distribution
- ▶ We can write this collection of random variables as $\{f(x) : x \in \mathcal{X}\}$ i.e. a function f evaluated at inputs x
- ▶ A GP is completely specified by
 - ▶ Mean function, $\mu(x) = \mathbb{E}(f(x))$
 - ▶ Covariance / kernel function, $k(x, x') = \text{Cov}(f(x), f(x'))$
 - ▶ Denoted $f \sim \text{GP}(\mu, k)$

GAUSSIAN PROCESSES

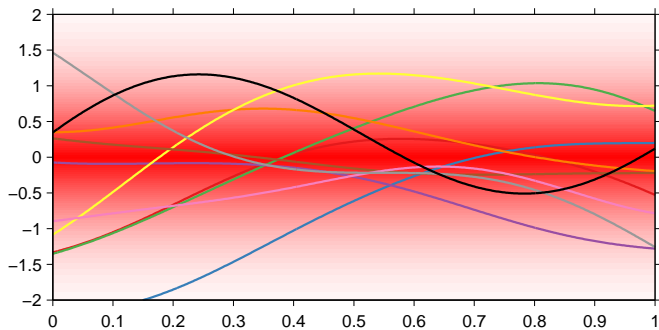
- ▶ A Gaussian process is collection of random variables, any finite number of which have a joint Gaussian distribution
- ▶ We can write this collection of random variables as $\{f(x) : x \in \mathcal{X}\}$ i.e. a function f evaluated at inputs x
- ▶ A GP is completely specified by
 - ▶ Mean function, $\mu(x) = \mathbb{E}(f(x))$
 - ▶ Covariance / kernel function, $k(x, x') = \text{Cov}(f(x), f(x'))$
 - ▶ Denoted $f \sim \text{GP}(\mu, k)$
- ▶ Can be thought of as a probability distribution on functions

LINEAR KERNEL = LINEAR FUNCTIONS



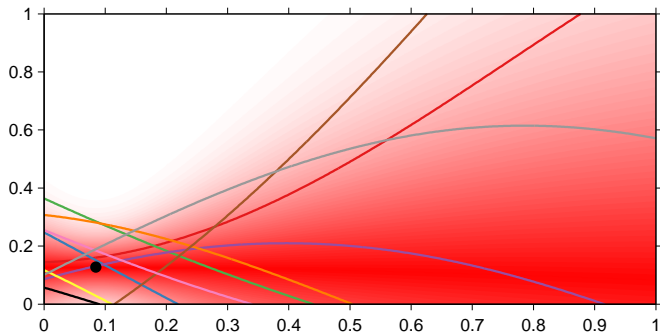
$$k(x, x') = xx'$$

WHAT ABOUT OTHER KERNELS?

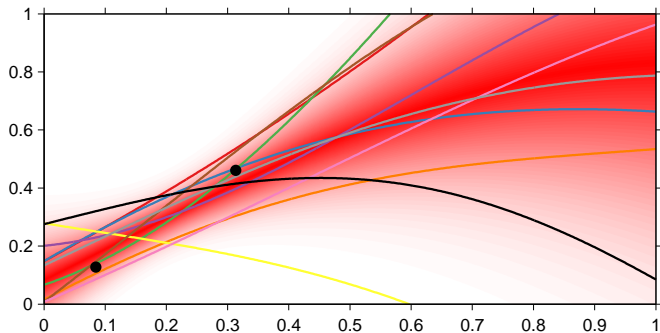


$$k(x, x') = e^{-(x-x')^2}$$

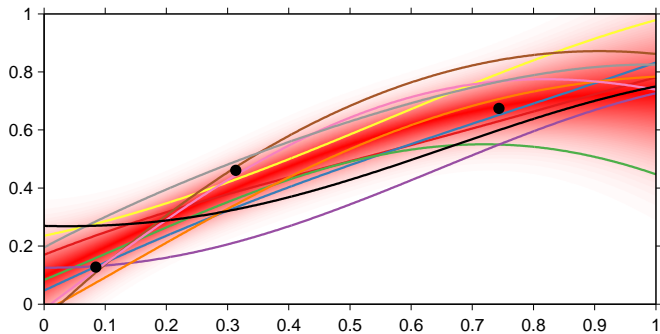
BAYESIAN NON-LINEAR REGRESSION



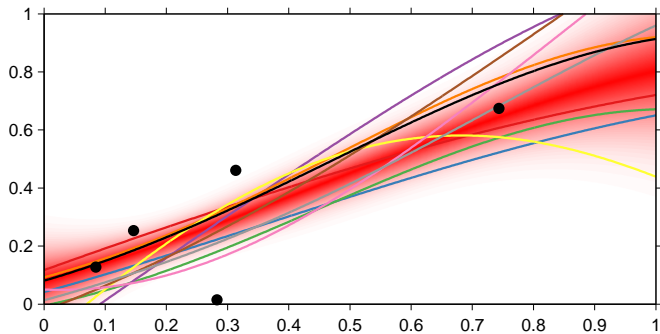
BAYESIAN NON-LINEAR REGRESSION



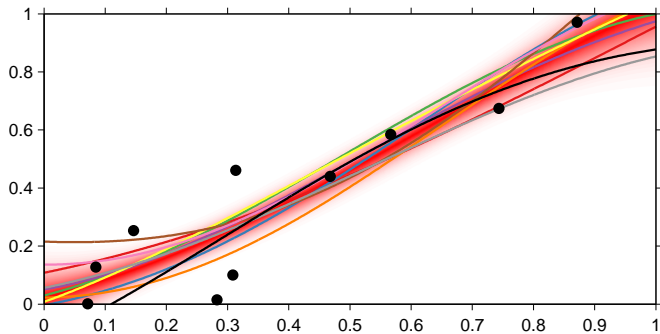
BAYESIAN NON-LINEAR REGRESSION



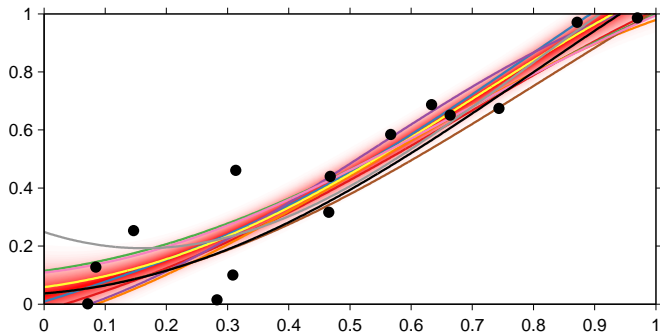
BAYESIAN NON-LINEAR REGRESSION



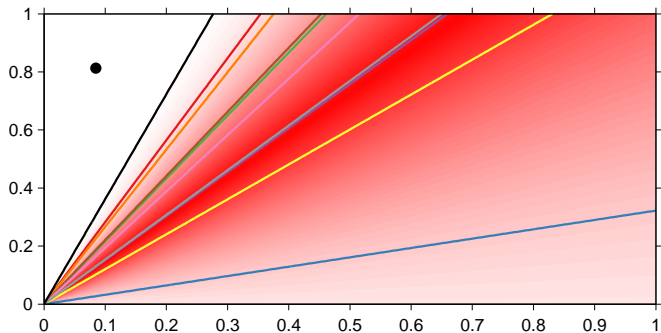
BAYESIAN NON-LINEAR REGRESSION



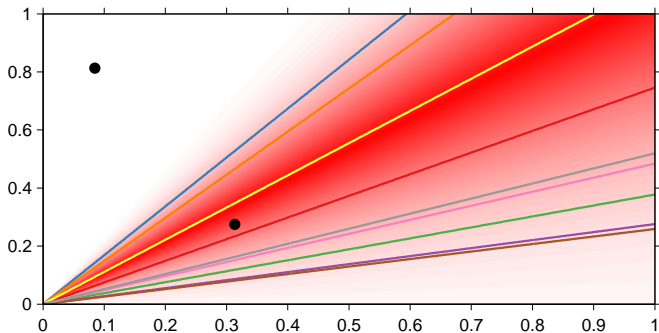
BAYESIAN NON-LINEAR REGRESSION



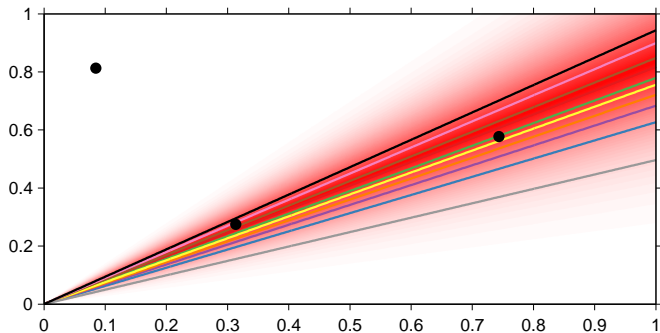
BAYESIAN LINEAR REGRESSION GONE WRONG



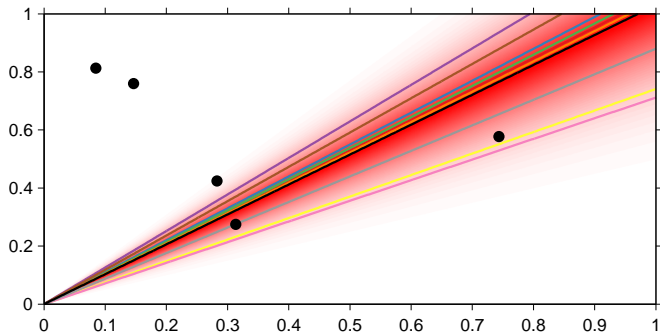
BAYESIAN LINEAR REGRESSION GONE WRONG



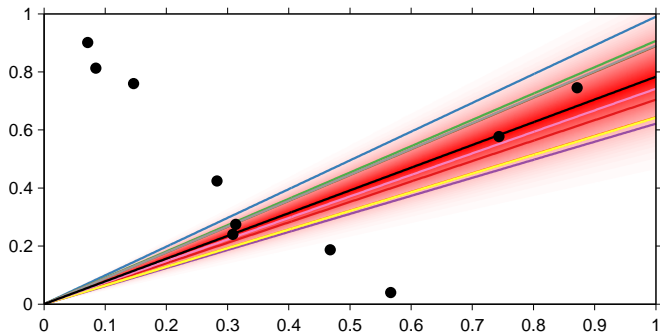
BAYESIAN LINEAR REGRESSION GONE WRONG



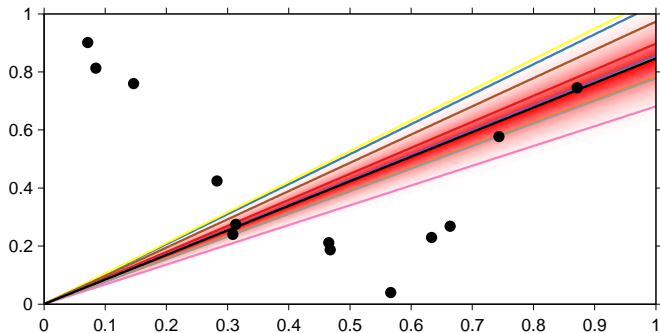
BAYESIAN LINEAR REGRESSION GONE WRONG



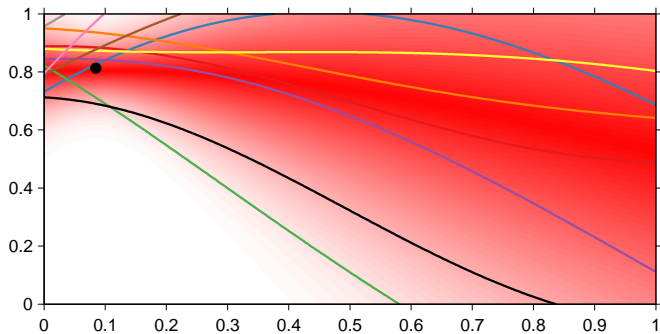
BAYESIAN LINEAR REGRESSION GONE WRONG



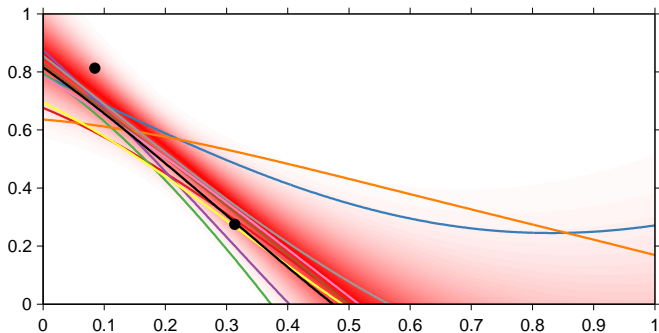
BAYESIAN LINEAR REGRESSION GONE WRONG



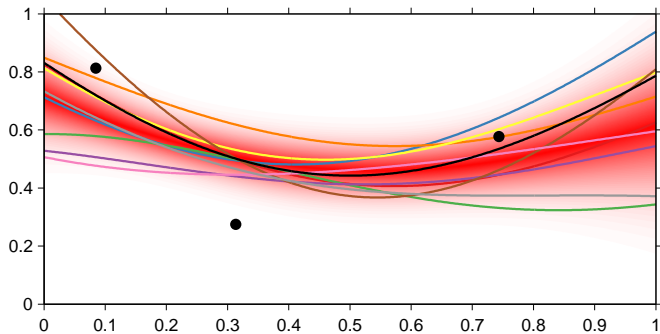
NON-LINEARITY TO THE RESCUE



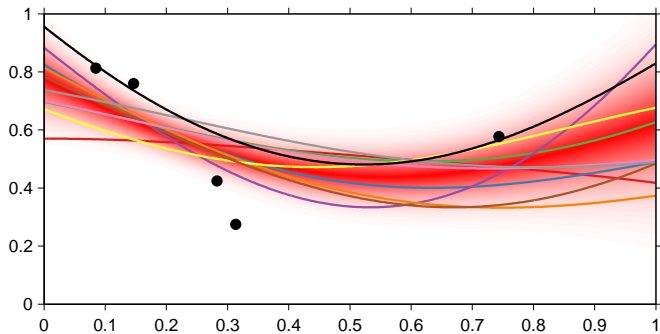
NON-LINEARITY TO THE RESCUE



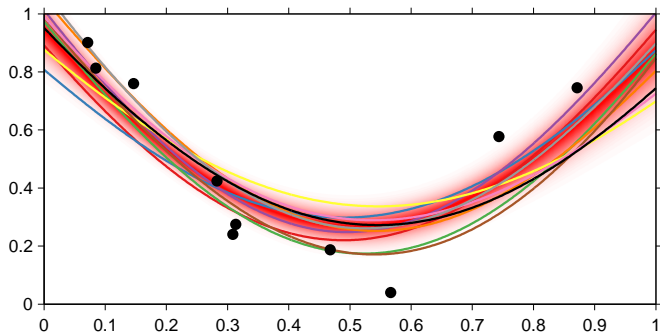
NON-LINEARITY TO THE RESCUE



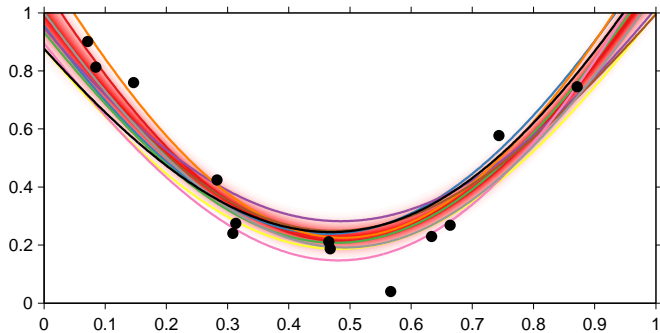
NON-LINEARITY TO THE RESCUE



NON-LINEARITY TO THE RESCUE



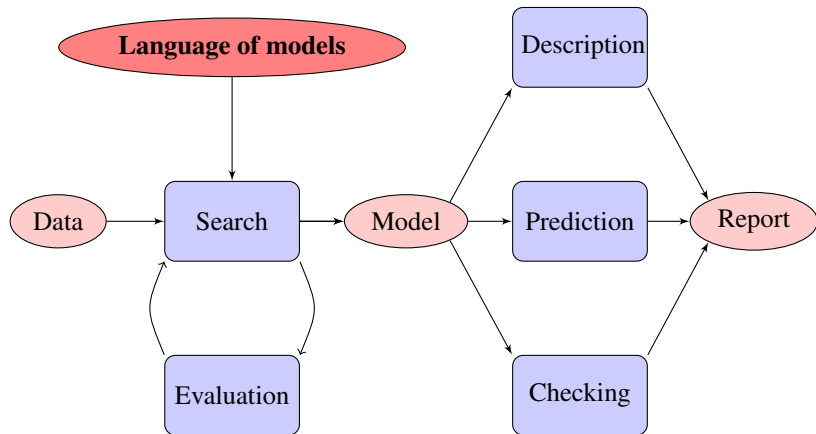
NON-LINEARITY TO THE RESCUE



WHICH KERNEL FUNCTION?

How far can we go with kernel functions?

DEFINING A LANGUAGE OF MODELS



THE ATOMS OF OUR LANGUAGE

Five base kernels



Squared
exp. (SE)



Periodic
(PER)



Linear
(LIN)



Constant
(C)



White
noise (WN)

Encoding for the following types of functions



Smooth
functions



Periodic
functions



Linear
functions



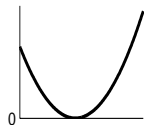
Constant
functions



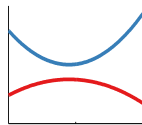
Gaussian
noise

THE COMPOSITION RULES OF OUR LANGUAGE

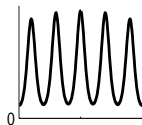
- ▶ Two main operations: addition, multiplication



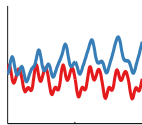
$\text{LIN} \times \text{LIN}$



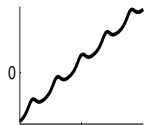
quadratic
functions



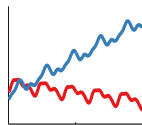
$\text{SE} \times \text{PER}$



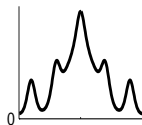
locally
periodic



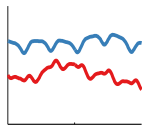
$\text{LIN} + \text{PER}$



periodic plus
linear trend



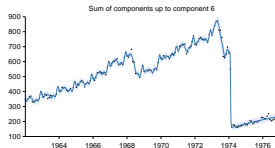
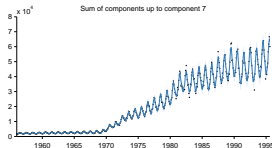
$\text{SE} + \text{PER}$



periodic plus
smooth trend

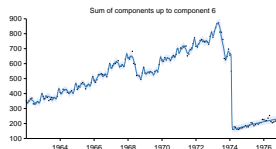
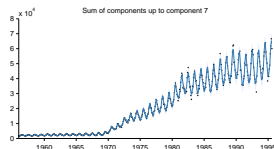
MODELING CHANGEPOINTS

Time series data often exhibit changepoints:



MODELING CHANGEPOINTS

Time series data often exhibit changepoints:



We can model this by assuming $f_1(x) \sim \text{GP}(0, k_1)$ and $f_2(x) \sim \text{GP}(0, k_2)$ and then defining

$$f(x) = (1 - \sigma(x))f_1(x) + \sigma(x)f_2(x)$$

where σ is a sigmoid function between 0 and 1.

MODELING CHANGEPOINTS

We can model this by assuming $f_1(x) \sim \text{GP}(0, k_1)$ and $f_2(x) \sim \text{GP}(0, k_2)$ and then defining

$$f(x) = (1 - \sigma(x))f_1(x) + \sigma(x)f_2(x)$$

where σ is a sigmoid function between 0 and 1.

Then $f \sim \text{GP}(0, k)$, where

$$k(x, x') = (1 - \sigma(x))k_1(x, x')(1 - \sigma(x')) + \sigma(x)k_2(x, x')\sigma(x')$$

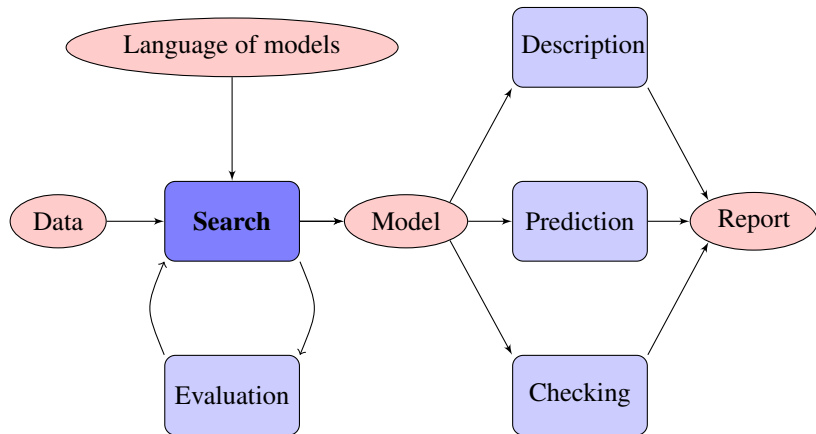
We define the changepoint operator $k = \text{CP}(k_1, k_2)$.

AN EXPRESSIVE LANGUAGE OF MODELS

Regression model	Kernel
GP smoothing	$SE + WN$
Linear regression	$C + LIN + WN$
Multiple kernel learning	$\sum SE + WN$
Trend, cyclical, irregular	$\sum SE + \sum PER + WN$
Fourier decomposition	$C + \sum \cos + WN$
Sparse spectrum GPs	$\sum \cos + WN$
Spectral mixture	$\sum SE \times \cos + WN$
Changepoints	e.g. $CP(SE, SE) + WN$
Heteroscedasticity	e.g. $SE + LIN \times WN$

Note: \cos is a special case of our version of PER

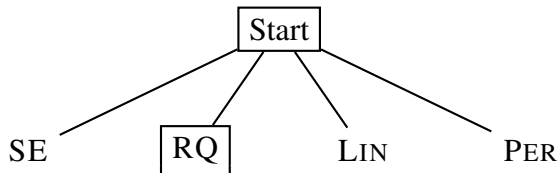
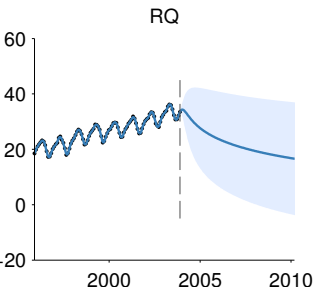
DISCOVERING A GOOD MODEL VIA SEARCH



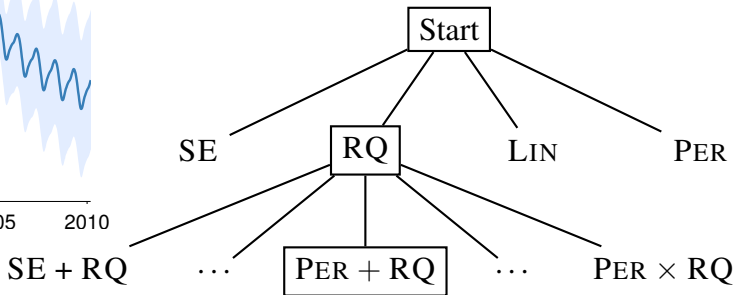
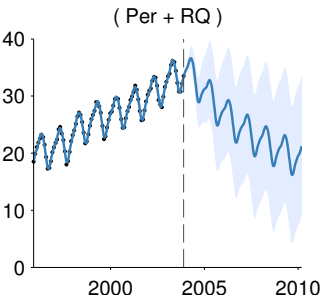
DISCOVERING A GOOD MODEL VIA SEARCH

- ▶ Language defined as the arbitrary composition of five base kernels (WN, C, LIN, SE, PER) via three operators (+, \times , CP).
- ▶ The space spanned by this language is open-ended and can have a high branching factor requiring a judicious search
- ▶ We propose a greedy search for its simplicity and similarity to human model-building

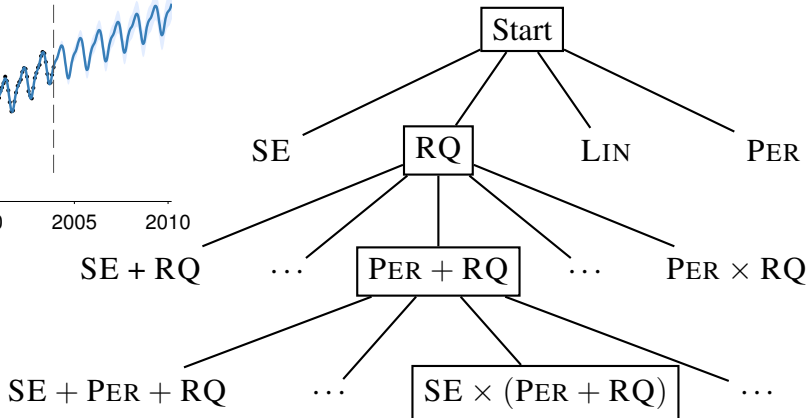
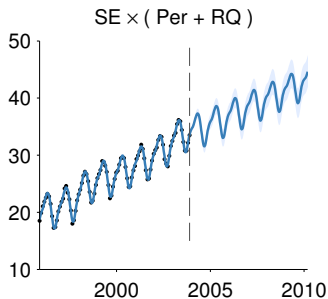
EXAMPLE: MAUNA LOA KEELING CURVE



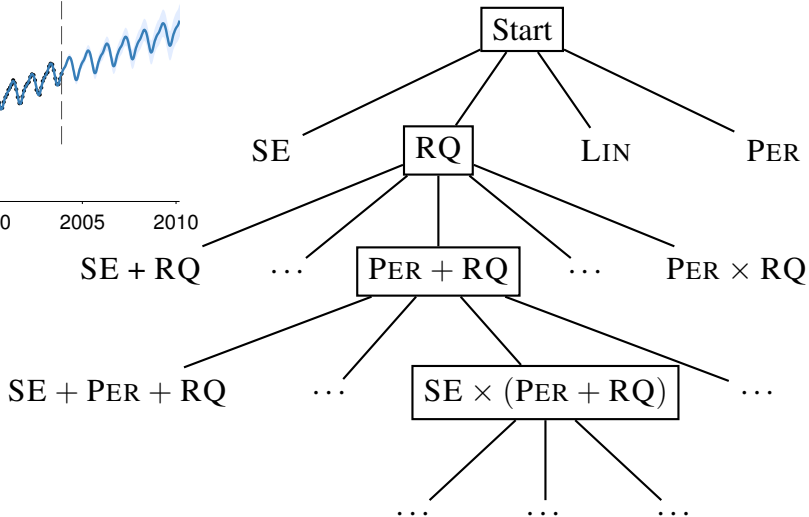
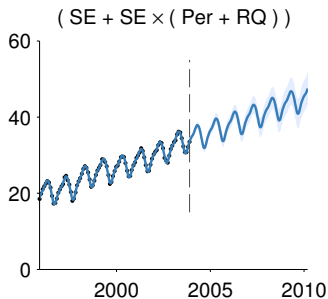
EXAMPLE: MAUNA LOA KEELING CURVE



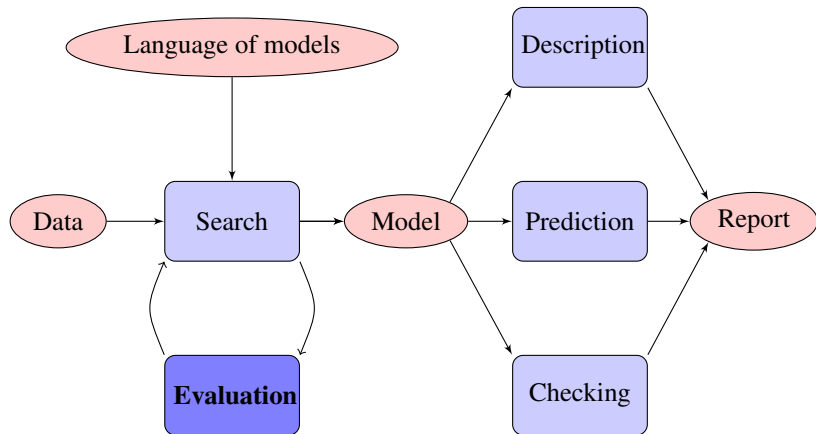
EXAMPLE: MAUNA LOA KEELING CURVE



EXAMPLE: MAUNA LOA KEELING CURVE



MODEL EVALUATION



BAYESIAN MODEL SELECTION

Suppose we have a collection of models $\{M_i\}$ and some data D

BAYESIAN MODEL SELECTION

Suppose we have a collection of models $\{M_i\}$ and some data D

Bayes rule tells us

$$p(M_i | D) = \frac{p(D | M_i)p(M_i)}{p(D)}$$

BAYESIAN MODEL SELECTION

Suppose we have a collection of models $\{M_i\}$ and some data D

Bayes rule tells us

$$p(M_i | D) = \frac{p(D | M_i)p(M_i)}{p(D)}$$

If $p(M_i)$ is equal for all i (prior ignorance) then

$$p(M_i | D) \propto p(D | M_i) = \int p(D | \theta_i, M_i)p(\theta_i | M_i)d\theta_i$$

BAYESIAN MODEL SELECTION

Suppose we have a collection of models $\{M_i\}$ and some data D

Bayes rule tells us

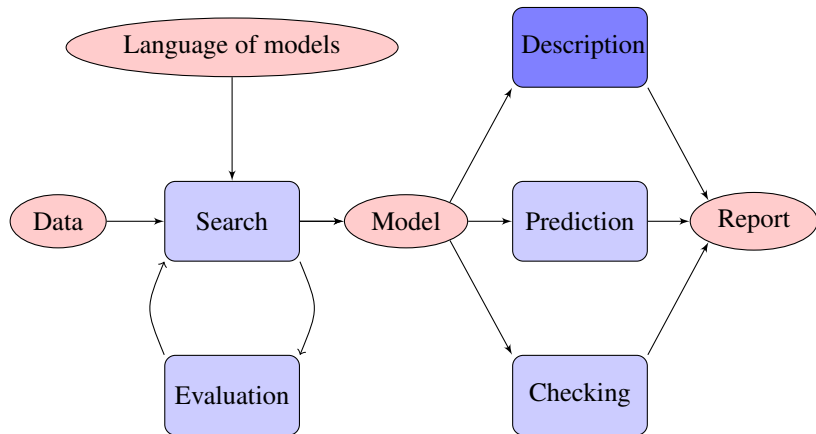
$$p(M_i | D) = \frac{p(D | M_i)p(M_i)}{p(D)}$$

If $p(M_i)$ is equal for all i (prior ignorance) then

$$p(M_i | D) \propto p(D | M_i) = \int p(D | \theta_i, M_i)p(\theta_i | M_i)d\theta_i$$

i.e. The most likely model has the highest **marginal likelihood**

AUTOMATIC TRANSLATION OF MODELS



AUTOMATIC TRANSLATION OF MODELS

- ▶ Search can produce **arbitrarily complicated models** from open-ended language but two main properties allow description to be automated
- ▶ Kernels can be **decomposed** into a **sum of products**
 - ▶ A sum of kernels corresponds to a sum of functions
 - ▶ Therefore, we can describe each product of kernels separately
- ▶ Each kernel in a product modifies a model in a **consistent** way
 - ▶ Each kernel roughly corresponds to an adjective

SUM OF PRODUCTS NORMAL FORM

Suppose the search finds the following kernel

$$SE \times (WN \times LIN + CP(C, PER))$$

SUM OF PRODUCTS NORMAL FORM

Suppose the search finds the following kernel

$$SE \times (WN \times LIN + CP(C, PER))$$

The changepoint can be converted into a sum of products

$$SE \times (WN \times LIN + C \times \sigma + PER \times \bar{\sigma})$$

SUM OF PRODUCTS NORMAL FORM

Suppose the search finds the following kernel

$$\text{SE} \times (\text{WN} \times \text{LIN} + \text{CP}(\text{C}, \text{PER}))$$

The changepoint can be converted into a sum of products

$$\text{SE} \times (\text{WN} \times \text{LIN} + \text{C} \times \sigma + \text{PER} \times \bar{\sigma})$$

Multiplication can be distributed over addition

$$\text{SE} \times \text{WN} \times \text{LIN} + \text{SE} \times \text{C} \times \sigma + \text{SE} \times \text{PER} \times \bar{\sigma}$$

SUM OF PRODUCTS NORMAL FORM

Suppose the search finds the following kernel

$$SE \times (WN \times LIN + CP(C, PER))$$

The changepoint can be converted into a sum of products

$$SE \times (WN \times LIN + C \times \sigma + PER \times \bar{\sigma})$$

Multiplication can be distributed over addition

$$SE \times WN \times LIN + SE \times C \times \sigma + SE \times PER \times \bar{\sigma}$$

Simplification rules are applied

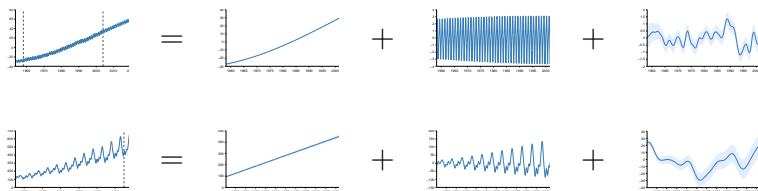
$$WN \times LIN + SE \times \sigma + SE \times PER \times \bar{\sigma}$$

SUMS OF KERNELS ARE SUMS OF FUNCTIONS

If $f_1 \sim \text{GP}(0, k_1)$ and independently $f_2 \sim \text{GP}(0, k_2)$ then

$$f_1 + f_2 \sim \text{GP}(0, k_1 + k_2)$$

e.g.

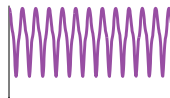
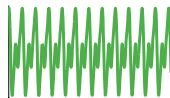
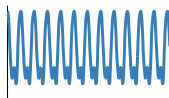
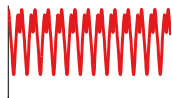


We can therefore describe each component separately

PRODUCTS OF KERNELS

PER
periodic function

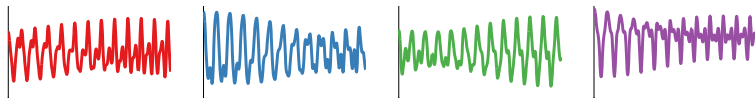
On their own, each kernel is described by a standard noun phrase



PRODUCTS OF KERNELS - SE

$$\underbrace{\text{SE}}_{\text{approximately}} \times \underbrace{\text{PER}}_{\text{periodic function}}$$

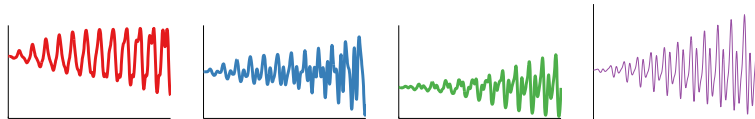
Multiplication by SE removes long range correlations from a model since $\text{SE}(x, x')$ decreases monotonically to 0 as $|x - x'|$ increases.



PRODUCTS OF KERNELS - LIN

$\underbrace{\text{SE}}$ \times $\underbrace{\text{PER}}$ \times $\underbrace{\text{LIN}}$
approximately periodic function with linearly growing amplitude

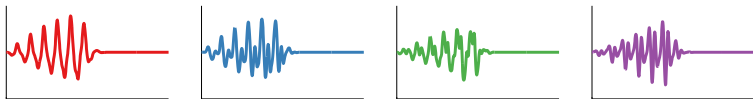
Multiplication by LIN is equivalent to multiplying the function being modeled by a linear function. If $f(x) \sim \text{GP}(0, k)$, then $xf(x) \sim \text{GP}(0, k \times \text{LIN})$. This causes the standard deviation of the model to vary linearly without affecting the correlation.



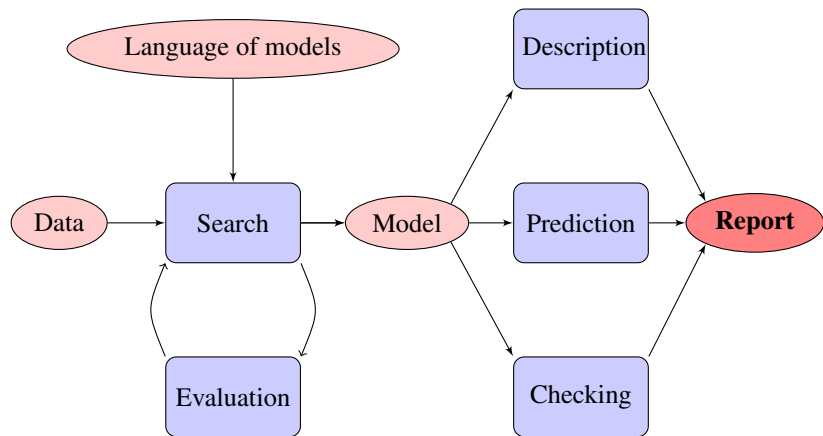
PRODUCTS OF KERNELS - CHANGEPOINTS

$\underbrace{\text{SE}}$ \times $\underbrace{\text{PER}}$ \times $\underbrace{\text{LIN}}$ \times $\underbrace{\sigma}$
approximately periodic function with linearly growing amplitude until 1700

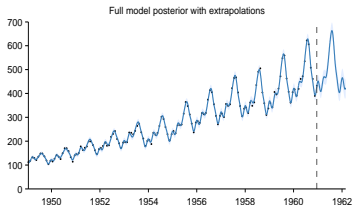
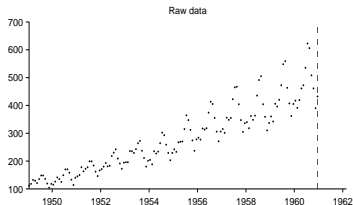
Multiplication by σ is equivalent to multiplying the function being modeled by a sigmoid.



AUTOMATICALLY GENERATED REPORTS



EXAMPLE: AIRLINE PASSENGER VOLUME

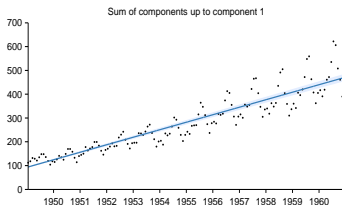
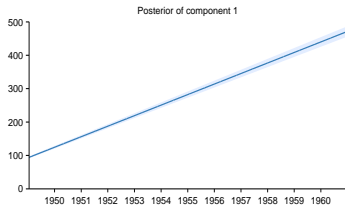


Four additive components have been identified in the data

- ▶ A linearly increasing function.
- ▶ An approximately periodic function with a period of 1.0 years and with linearly increasing amplitude.
- ▶ A smooth function.
- ▶ Uncorrelated noise with linearly increasing standard deviation.

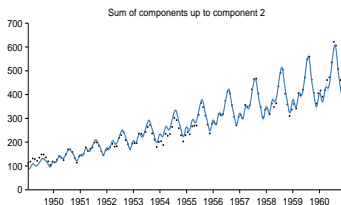
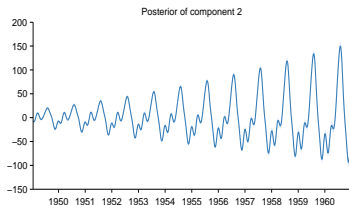
EXAMPLE: AIRLINE PASSENGER VOLUME

This component is linearly increasing.



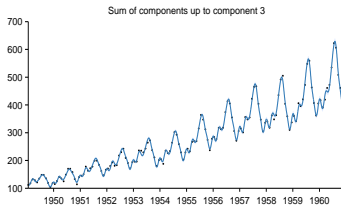
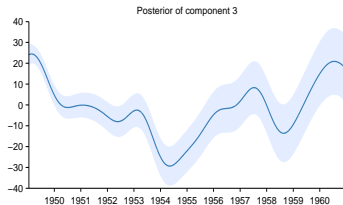
EXAMPLE: AIRLINE PASSENGER VOLUME

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 linearly. The shape of this function within each period has a typical lengthscale of 6.0 weeks.



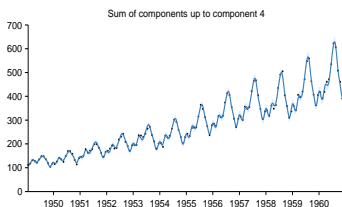
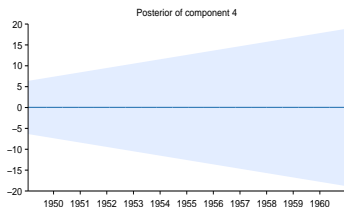
EXAMPLE: AIRLINE PASSENGER VOLUME

This component is a smooth function with a typical lengthscale of 8.1 months.



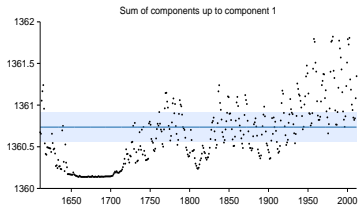
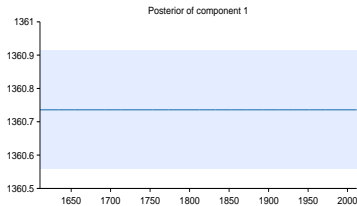
EXAMPLE: AIRLINE PASSENGER VOLUME

This component models uncorrelated noise. The standard deviation of the noise increases linearly.



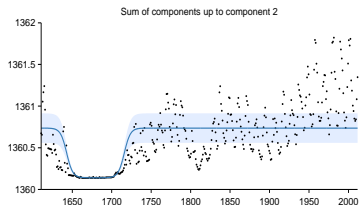
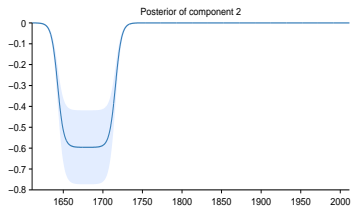
EXAMPLE: SOLAR IRRADIANCE

This component is constant.



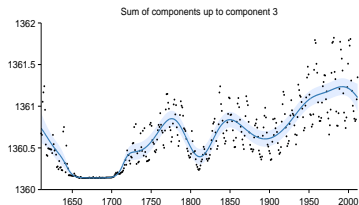
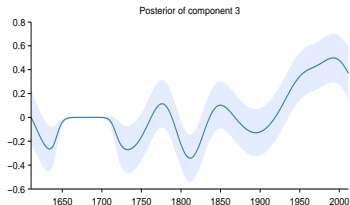
EXAMPLE: SOLAR IRRADIANCE

This component is constant. This component applies from 1643 until 1716.



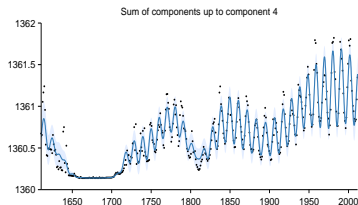
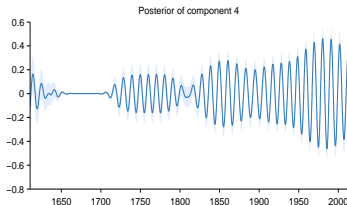
EXAMPLE: SOLAR IRRADIANCE

This component is a smooth function with a typical lengthscale of 23.1 years. This component applies until 1643 and from 1716 onwards.



EXAMPLE: SOLAR IRRADIANCE

This component is approximately periodic with a period of 10.8 years. Across periods the shape of this function varies smoothly with a typical lengthscale of 36.9 years. The shape of this function within each period is very smooth and resembles a sinusoid. This component applies until 1643 and from 1716 onwards.



EXAMPLE: FULL REPORT

See pdf

SUMMARY

- ▶ We have briefly introduced Bayesian statistics and Gaussian processes
- ▶ Described the beginnings of an automatic statistician for regression
 - ▶ Defines an open-ended language of models
 - ▶ Searches greedily through this space
 - ▶ Produces detailed reports describing patterns in data
- ▶ We believe this line of research has the potential to make powerful statistical model-building techniques accessible to non-experts

THANKS

Thanks

APPENDIX

SUM AND PRODUCT RULES OF PROBABILITY

Sum rule :

$$\mathbb{P}(X = x) = \sum_y \mathbb{P}(X = x, Y = y)$$

Product rule :

$$\mathbb{P}(X = x, Y = y) = \mathbb{P}(X = x) \mathbb{P}(Y = y | X = x)$$

SUM AND PRODUCT RULES OF PROBABILITY

Sum rule :
$$p(x) = \int p(x, y) dy$$

Product rule :
$$p(x, y) = p(x) p(y | x)$$

DERIVING BAYES RULE

$$p(\theta, D | M) = p(\theta, D | M)$$

DERIVING BAYES RULE

$$p(\theta, D | M) = p(\theta, D | M)$$

$$p(\theta | D, M)p(D | M) = p(D | \theta, M)p(\theta | M) \quad (\text{product rule})$$

DERIVING BAYES RULE

$$p(\theta, D | M) = p(\theta, D | M)$$

$$p(\theta | D, M)p(D | M) = p(D | \theta, M)p(\theta | M) \quad (\text{product rule})$$

$$p(\theta | D, M) = \frac{p(D | \theta, M)p(\theta | M)}{p(D | M)} \quad (\text{divide})$$

DERIVING BAYES RULE

$$p(\theta, D | M) = p(\theta, D | M)$$

$$p(\theta | D, M)p(D | M) = p(D | \theta, M)p(\theta | M) \quad (\text{product rule})$$

$$p(\theta | D, M) = \frac{p(D | \theta, M)p(\theta | M)}{p(D | M)} \quad (\text{divide})$$

$$p(\theta | D, M) = \frac{p(D | \theta, M)p(\theta | M)}{\int p(D, \theta | M)d\theta} \quad (\text{sum rule})$$

DERIVING BAYES RULE

$$p(\theta, D | M) = p(\theta, D | M)$$

$$p(\theta | D, M)p(D | M) = p(D | \theta, M)p(\theta | M) \quad (\text{product rule})$$

$$p(\theta | D, M) = \frac{p(D | \theta, M)p(\theta | M)}{p(D | M)} \quad (\text{divide})$$

$$p(\theta | D, M) = \frac{p(D | \theta, M)p(\theta | M)}{\int p(D, \theta | M)d\theta} \quad (\text{sum rule})$$

$$p(\theta | D, M) = \frac{p(D | \theta, M)p(\theta | M)}{\int p(D | \theta, M)p(\theta | M)d\theta} \quad (\text{product rule})$$

DEFINING A LANGUAGE OF REGRESSION MODELS

- ▶ Regression consists of **learning a function** $f : \mathcal{X} \rightarrow \mathcal{Y}$ from inputs to outputs from example input / output pairs
- ▶ Language should include **simple parametric forms** . . .
 - ▶ e.g. Linear functions, Polynomials, Exponential functions
- ▶ . . . as well as functions specified by **high level properties**
 - ▶ e.g. Smoothness, Periodicity
- ▶ Inference should be **tractable for all models** in language

A LANGUAGE OF GAUSSIAN PROCESS KERNELS

- ▶ It is common practice to use a zero mean function since the mean can be marginalised out
 - ▶ Suppose, $f(x) | a \sim \text{GP}(a \times \mu(x), k(x, x'))$ where $a \sim \mathcal{N}(0, 1)$
 - ▶ Then equivalently, $f(x) \sim \text{GP}(0, \mu(x)\mu(x') + k(x, x'))$
- ▶ We therefore define a language of GP regression models by specifying a **language of kernels**

MODEL EVALUATION

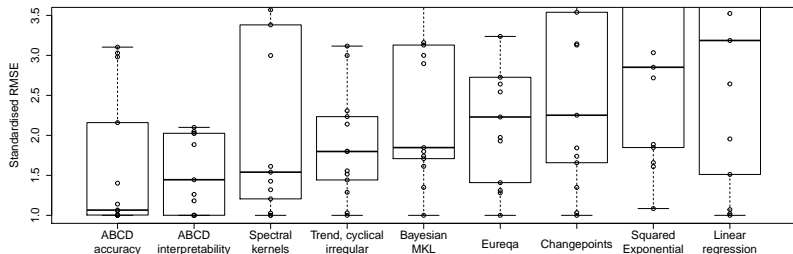
- ▶ After proposing a new model its kernel parameters are optimised by conjugate gradients
- ▶ We evaluate each optimised model, M , using the **marginal likelihood** which can be computed analytically for GPs
- ▶ We **penalise** the marginal likelihood for the **optimised kernel parameters** using the Bayesian Information Criterion (BIC):

$$-0.5 \times \text{BIC}(M) = \log p(D | M) - \frac{p}{2} \log n$$

where p is the number of kernel parameters, D represents the data, and n is the number of data points.

GOOD PREDICTIVE PERFORMANCE AS WELL

Standardised RMSE over 13 data sets

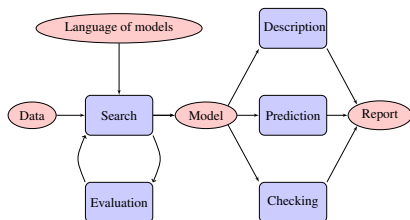


- ▶ Tweaks can be made to the algorithm to improve accuracy or interpretability of models produced. . .
- ▶ . . . but both methods are highly competitive at extrapolation (shown above) and interpolation

GOALS OF THE AUTOMATIC STATISTICIAN PROJECT

- ▶ Provide a set of tools for understanding data that require minimal expert input
- ▶ Uncover challenging research problems in e.g.
 - ▶ Automated inference
 - ▶ Model construction and comparison
 - ▶ Data visualisation and interpretation
- ▶ Advance the field of machine learning in general

INGREDIENTS OF AN AUTOMATIC STATISTICIAN



- ▶ **An open-ended language of models**
 - ▶ Expressive enough to capture real-world phenomena...
 - ▶ ...and the techniques used by human statisticians
- ▶ **A search procedure**
 - ▶ To efficiently explore the language of models
- ▶ **A principled method of evaluating models**
 - ▶ Trading off complexity and fit to data
- ▶ **A procedure to automatically explain the models**
 - ▶ Making the assumptions of the models explicit...
 - ▶ ...in a way that is intelligible to non-experts

CHALLENGES

- ▶ Interpretability / accuracy
- ▶ Increasing the expressivity of language
 - ▶ e.g. Monotonicity, positive functions, symmetries
- ▶ Computational complexity of searching through a huge space of models
- ▶ Extending the automatic reports to multidimensional datasets
 - ▶ Search and descriptions naturally extend to multiple dimensions, but automatically generating relevant visual summaries harder

CURRENT AND FUTURE DIRECTIONS

- ▶ Automatic statistician for:
 - ▶ Multivariate nonlinear regression
 - ▶ Classification
 - ▶ Completing and interpreting tables and databases

- ▶ Probabilistic programming
 - ▶ Probabilistic models are expressed in a general (Turing complete) programming language (e.g. Church/Venture/Anglican)
 - ▶ A universal inference engine can then be used to infer unobserved variables given observed data
 - ▶ This can be used to implement search over the model space in an automated statistician