

Module 3: Bayesian Nonparametrics

Gaussian Processes for Regression

STAT/BIOSTAT 527, University of Washington

Emily Fox

April 21st-22nd, 2014

©Emily Fox 2014

1

Maximum Likelihood Estimation

- Model:

$$y = f(x) + \epsilon \quad \text{where } \epsilon \sim N(0, \sigma^2)$$

$$f(x) = \sum_{m=1}^M \beta_m h_m(x)$$

- Equivalently,

$$p(y | x, \beta, \sigma^2) = N(y | f(x), \sigma^2)$$

- For our training data (independent obs)

$$p(y | X, \beta, \sigma^2) = \prod_{i=1}^n N(y_i | f(x_i), \sigma^2)$$

max log-like
(instead of min L)

©Emily Fox 2014

2

Maximum Likelihood Estimation

$$p(y | X, \beta, \sigma^2) = \prod_i N(y_i | \beta^T h(x_i), \sigma^2)$$

- Taking the log

$$\log p(y | X, \beta, \sigma^2) = \sum_i \underbrace{-\frac{1}{2}(y_i - \beta^T h(x_i))^2}_{\text{const. wrt } \beta} - \frac{n}{2} \log 2\pi - \frac{n}{2} \log \sigma^2$$

$$\mathcal{N}(x | \mu, \Sigma) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

- Equivalent objective to RSS (Gaussian log-like loss = L_2 loss)

- Taking the gradient and setting to zero, we have already shown

$$\hat{\beta}^{ML} = (H^T H)^{-1} H^T y$$

point estimate $\leftarrow H = \begin{pmatrix} h_1(x_1) & \dots & h_n(x_1) \\ \vdots & & \vdots \\ h_1(x_n) & \dots & h_n(x_n) \end{pmatrix}$

©Emily Fox 2014

3

Bayesian Linear Regression

- More generally, consider a conjugate prior on the basis expansion coefficients:

$$\text{prior } p(\beta) = N(\beta | \mu_0, \Sigma_0) \quad \leftarrow \text{prior mean + cov.}$$

- Combining this with the Gaussian likelihood function, and using standard Gaussian identities, gives posterior

$$\text{where } p(\beta | y) = N(\beta | \mu_n, \Sigma_n) \quad \leftarrow \text{posterior (updated) mean + cov. after obs. } y$$

$$\mu_n = \Sigma_n^{-1} (\Sigma_0^{-1} \mu_0 + \sigma^{-2} H^T y)$$

$$\Sigma_n^{-1} = \Sigma_0^{-1} + \sigma^{-2} H^T H$$

} incorporates obs. y w/ prior

prior $\leftarrow \Sigma_0^{-1} \mu_0$ data dependent $\leftarrow \sigma^{-2} H^T H$

©Emily Fox 2014

4

Predictive Distribution

- Predict y^* at new locations x^* by integrating over parameters β

$$p(y^* | y) = \int p(y^* | \beta) p(\beta | y) d\beta$$

posterior:
 $p(\beta | y) = N(\beta | \mu_n, \Sigma_n)$

$y^* = h(x^*)^T \beta + \epsilon$
 $\beta | y \sim N(\mu_n, \Sigma_n)$
 $\epsilon \sim N(0, \sigma^2)$

likelihood
 $p(y | x, \beta, \sigma^2) = N(y | f(x), \sigma^2)$

$= N(\mu_n^*(x^*), \Sigma_n^*(x^*))$

Σ_n fun of our obs. locations x_1, \dots, x_n

$\mu_n^*(x^*) = E[y^* | y] = \mu_n^T h(x^*)$

$\Sigma_n^*(x^*) = \text{cov}(y^* | y) = h(x^*)^T \text{cov}(\beta | y) h(x^*) + \sigma^2$
 $= h(x^*)^T \Sigma_n h(x^*) + \sigma^2$

var of β var of obs.

©Emily Fox 2014

5

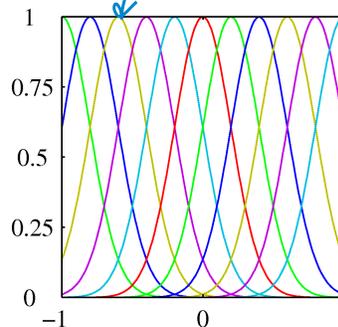
Example: Gaussian Basis Expansion

- Gaussian basis functions:

$$h_j(x) = \exp\left\{-\frac{(x - \mu_j)^2}{2s^2}\right\}$$

← a choice of a basis exp. with $M=9$

- These are local; a small change in x only affects nearby basis functions. Parameters control location and scale (width)

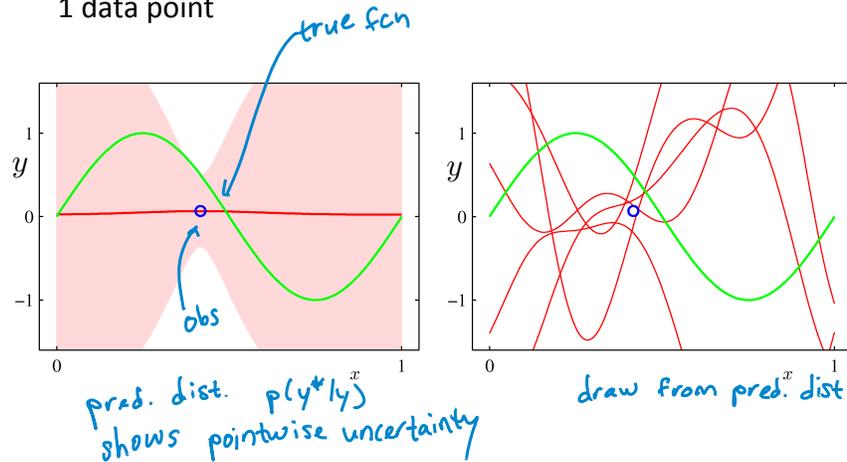


©Emily Fox 2014

6

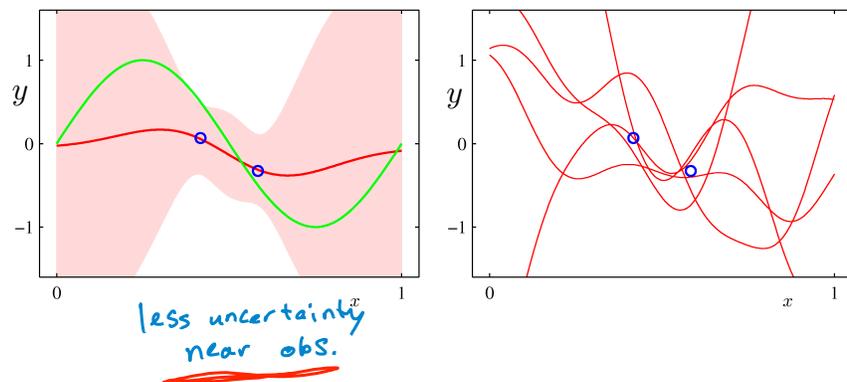
Example: Gaussian Basis Expansion

- Example: Sinusoidal data, 9 Gaussian basis functions, 1 data point



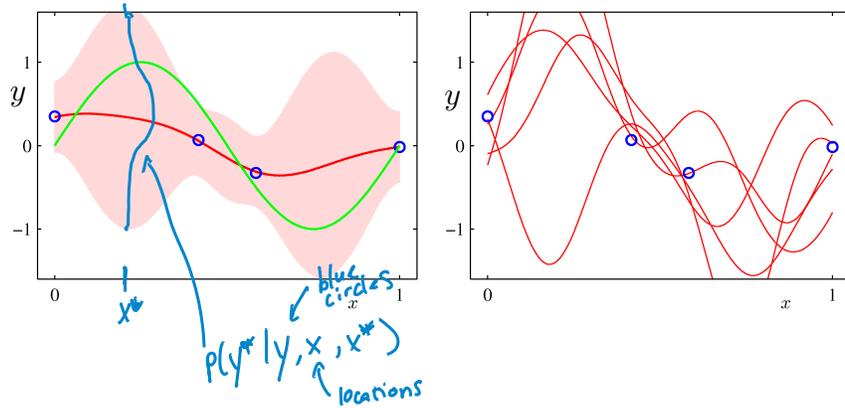
Example: Gaussian Basis Expansion

- Example: Sinusoidal data, 9 Gaussian basis functions, 2 data points



Example: Gaussian Basis Expansion

- Example: Sinusoidal data, 9 Gaussian basis functions, 4 data points

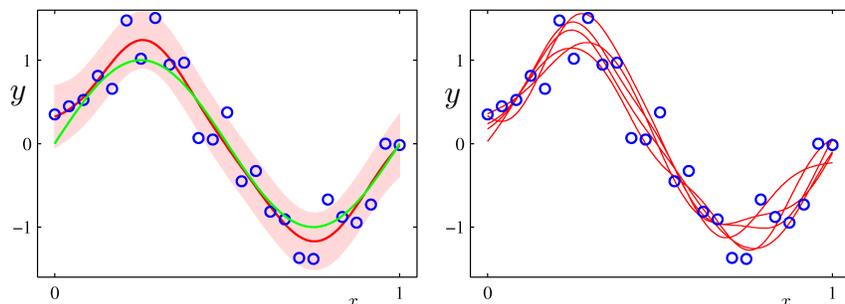


©Emily Fox 2014

9

Example: Gaussian Basis Expansion

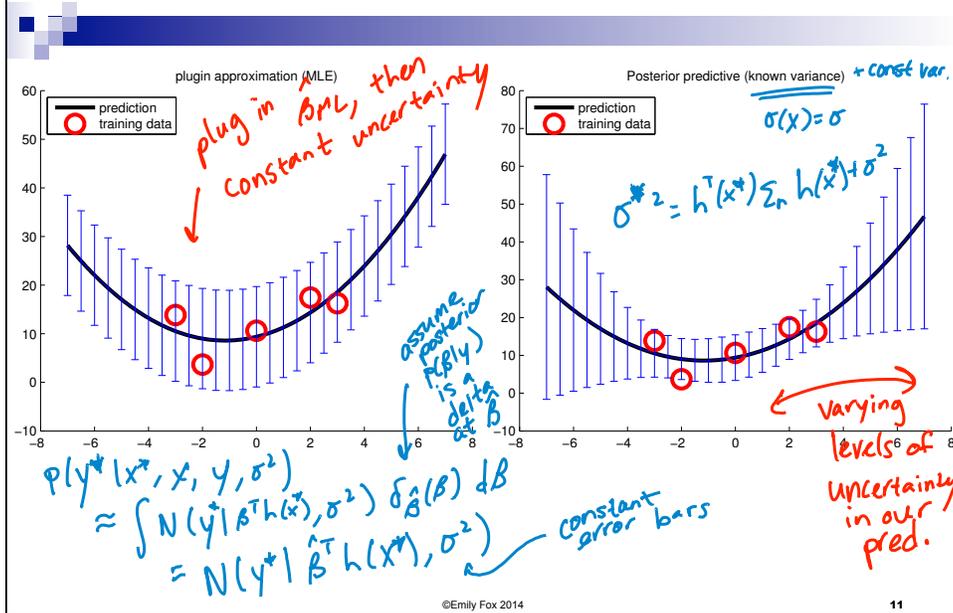
- Example: Sinusoidal data, 9 Gaussian basis functions, 25 data points



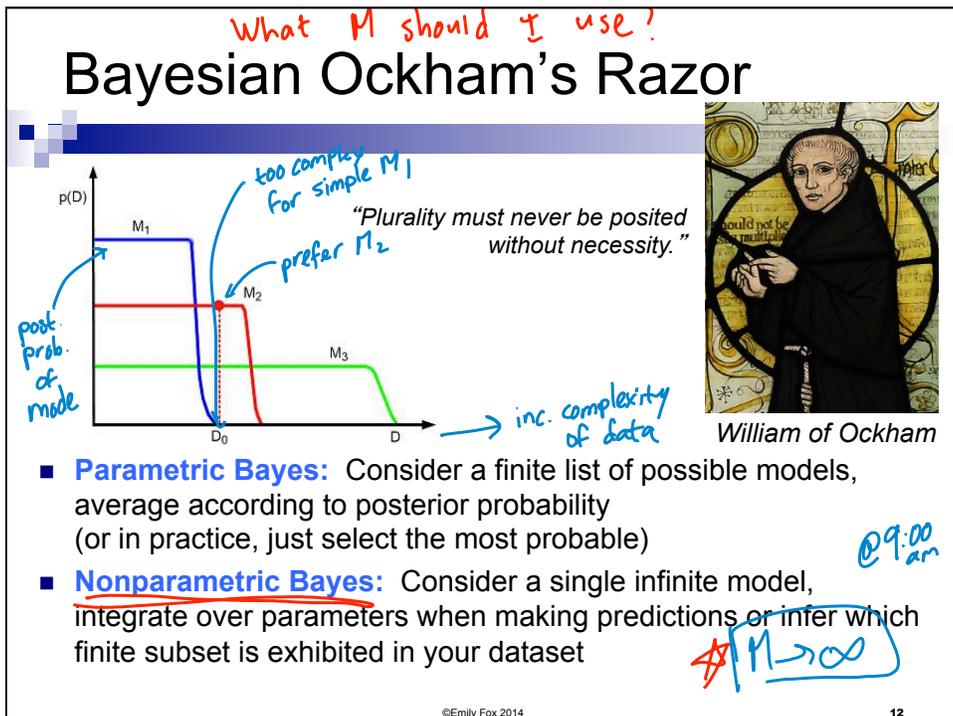
©Emily Fox 2014

10

Estimation vs. Predictive Distributions



Bayesian Ockham's Razor



Going Infinite...

model: $f(x) = \phi(x)^T \beta$

Change of notation:

$$h_u(x) \rightarrow \phi_j(x)$$

- Nonparametric Gaussian regression: Would like to let the number of basis functions $M \rightarrow \infty$ ← "features"

- Prior: $N(\beta | 0, \alpha^{-1} I_M)$ ← special case of Gaussian prior

- Distribution on f : $f = \Phi \beta$ ← linear comb. of Gaussians β_j
← $(f(x_1), \dots, f(x_n))^T$
← $(\beta_1, \dots, \beta_n)^T$
→ f is Gauss.

$$p(f) = N(f | 0, \alpha^{-1} \Phi \Phi^T)$$

$$E[f] = \Phi E[\beta] = 0$$

$$\text{cov}(f) = \Phi E[\beta \beta^T] = \frac{1}{\alpha} \Phi \Phi^T$$

- Gaussian process models replace explicit basis function representation with a direct specification in terms of a positive definite kernel function

Mercer Kernel Functions

- Distributions are of the form $p(f) = N(f | 0, \alpha^{-1} \Phi \Phi^T) = N(f | 0, K)$ ← $n \times n$ matrix
↑ $n \times M$ $M \times n$

where the **Gram matrix** K is defined as

$$K_{ij} = K(x_i, x_j) = \phi(x_i)^T \phi(x_j) \quad \text{dim } M$$

kernel ϕ_n →

- K is a **Mercer kernel** if the Gram matrix is positive definite for any n and any x_1, \dots, x_n

Note: K is $n \times n$ matrix regardless of M (dim of basis)

Example of the "kernel trick"

Mercer's Theorem

- If K is positive definite, we can compute the eigendecomp:

$$K = U^T \Lambda U$$

- Then $K_{ij} = (\Lambda^{1/2} U_{\cdot i})^T (\Lambda^{1/2} U_{\cdot j})$

- Define $\phi(x) = \Lambda^{1/2} U_{\cdot i}$ so that

$$K_{ij} = \phi(x_i)^T \phi(x_j) \quad M \text{ dim}$$

- If a kernel is Mercer, there exists a function $\phi : \mathcal{X} \rightarrow \mathbb{R}^M$ s.t.

$$K(x, x') = \phi(x)^T \phi(x')$$

\uparrow
M might be infinite

Hard to show in general given $K(x, x')$

©Emily Fox 2014

15

Example Mercer Kernels

- Example #1: (non-stationary) **polynomial kernel**

$$\kappa(x, x') = (\gamma x^T x' + r)^p$$

- For $p=2, \gamma = r = 1,$

$$(1 + x^T x')^2 = (1 + x_1 x'_1 + x_2 x'_2)^2 = 1 + 2x_1 x'_1 + 2x_2 x'_2 + (x_1 x'_1)^2 + (x_2 x'_2)^2 + 2x_1 x'_1 x_2 x'_2$$

- This can be written as $\phi(x)^T \phi(x')$, with

$$\phi(x) = [1, \sqrt{2} x_1, \sqrt{2} x_2, x_1^2, x_2^2, \sqrt{2} x_1 x_2]$$

- Equivalent to working in a 6-dimensional feature space
- For general p , basis contains all terms up to degree p

$M=6$

- Example #2: **Gaussian kernel**

$$\kappa(x, x') = \exp\left(-\frac{1}{2}(x - x')^T \Sigma^{-1}(x - x')\right)$$

- Feature map lives in an infinite-dimensional space

©Emily Fox 2014

16

Gaussian Processes

~~fixed β~~ \leftarrow prior on β
 work $f(x)$ \leftarrow prior on f

- Dispense of parametric view (prior on β) and consider prior on functions themselves (prior on f)

- Seems hard, but we have shown that it is feasible when we look at a finite set of values x_1, \dots, x_n

$$p(f) = N(f \mid 0, K) \quad \leftarrow \text{dist on } \begin{pmatrix} f(x_1) \\ \vdots \\ f(x_n) \end{pmatrix}$$

- Defined by a *Mercer kernel*
- More generally, a Gaussian process provides a distribution over functions

dist on $f(\cdot)$

Gaussian Processes

- Distribution on functions

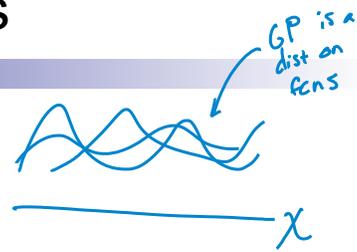
□ $f \sim \text{GP}(m, K)$

- m : mean function
- K : covariance function

\iff iff $\forall n$ and any x_1, \dots, x_n

□ $p(f(x_1), \dots, f(x_n)) \sim N_n(\mu, K)$

- $\mu = [m(x_1), \dots, m(x_n)]$
- $K_{ij} = K(x_i, x_j)$

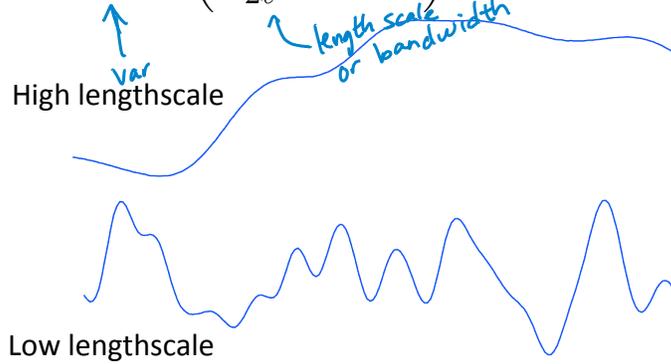


- Idea: If x_i, x_j are similar according to the kernel, then $f(x_i)$ is similar to $f(x_j)$

k: covariance function

Example: squared exp (SE) or Gauss. kernel

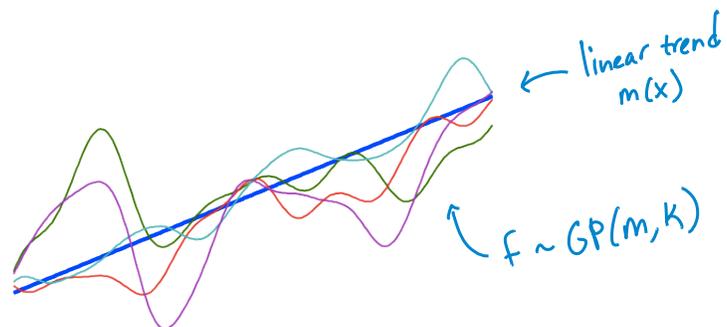
$$\kappa(x, x') = \sigma_f^2 \exp\left(-\frac{1}{2\ell^2}(x - x')^2\right)$$



©Emily Fox 2014

19

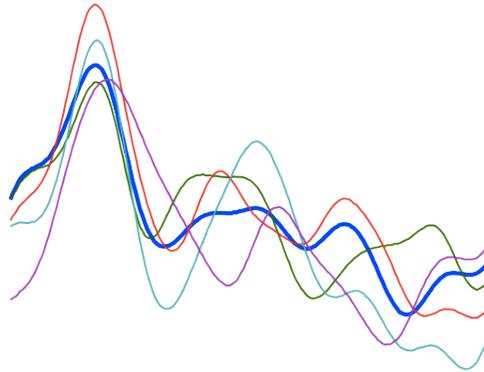
m: mean function



©Emily Fox 2014

20

m: mean function



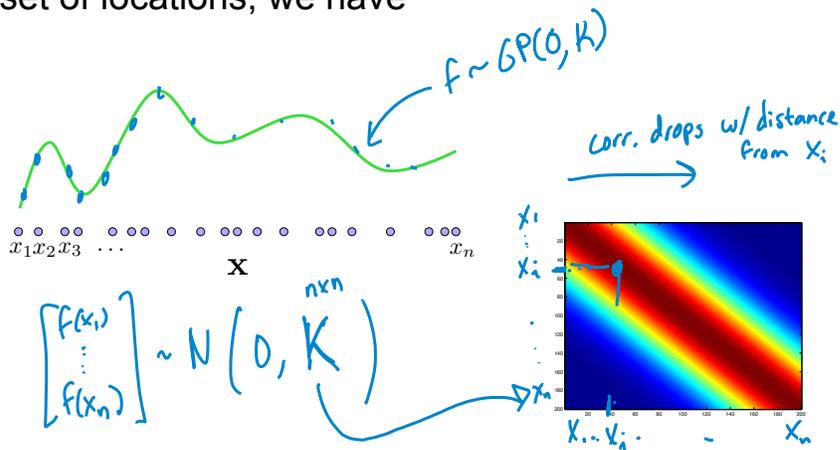
Can define more complicated $m(x)$, but typically people choose $m(x)=0$

©Emily Fox 2014

21

Induced Multivariate Gaussian

- Evaluating the GP-distributed function at any set of locations, we have

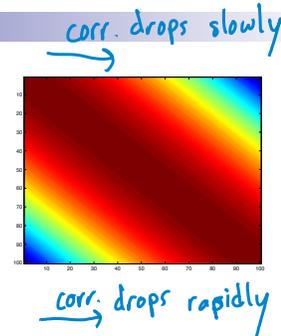
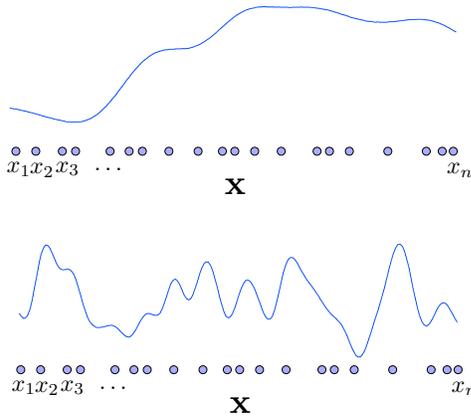


©Emily Fox 2014

22

Induced Multivariate Gaussian

■ Comparing length-scales:

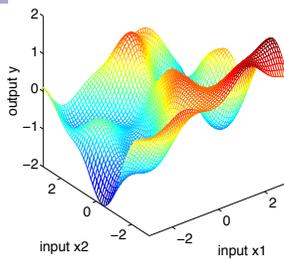


©Emily Fox 2014

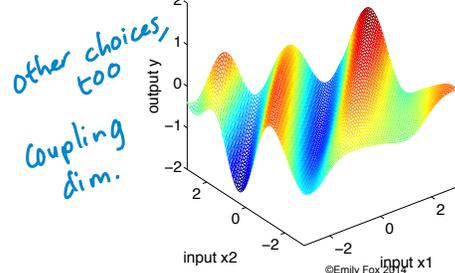
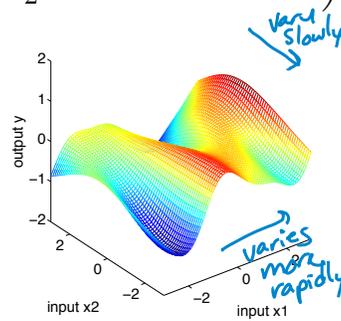
23

2D Gaussian Processes

$$\kappa(x_p, x_q) = \sigma_f^2 \exp\left(-\frac{1}{2}(x_p - x_q)^T M (x_p - x_q)\right)$$



$M = \ell^{-2} I$
 $k(x, x') = e^{-\frac{\|x-x'\|^2}{2\ell^2}}$
 radial basis fn (RBF)



diff. length scales to diff. dims

$$M = \begin{bmatrix} \frac{1}{\ell_1^2} & 0 \\ 0 & \frac{1}{\ell_2^2} \end{bmatrix}$$

$\ell_1 = 1 \quad \ell_2 = 3$

©Emily Fox 2014

24

GPs for Regression

(up to this point, just background on GPs... now have "data")

- Start with noise-free scenario: directly observe the function
- Training data $\mathcal{D} = \{(x_i, f_i), i = 1, \dots, n\}$
- Test data locations $X^* \rightarrow$ predict f^*

Jointly, we have *by defn of GP, if $f(\cdot) \sim \text{GP}(\mu, K)$ then*

$$\begin{pmatrix} f_1 \\ \vdots \\ f_n \\ f^* \end{pmatrix} \sim N \left(\begin{pmatrix} \mu \\ \mu_* \end{pmatrix}, \begin{pmatrix} K & K_* \\ K_*^T & K_{**} \end{pmatrix} \right)$$

Annotations:
 - f_1, \dots, f_n are circled in blue.
 - f^* is circled in blue with "cond. on this" written below it.
 - K is labeled $K(X, X)$.
 - K_* is labeled $K(X, X^*)$.
 - K_{**} is labeled $K(X^*, X^*)$.
 - μ_* is labeled $(m(x_1^*), \dots, m(x_n^*))$.
 - K_{**} is also labeled "training obs."

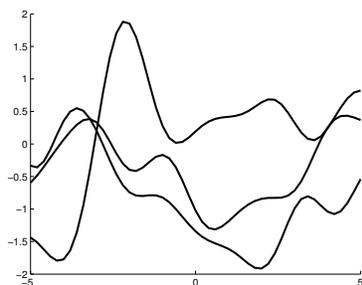
Therefore,

$$p(f^* | X^*, X, f) = N(f^* | \mu_* + K_*^T K^{-1} (f - \mu), K_{**} - K_*^T K^{-1} K_*)$$

©Emily Fox 2014

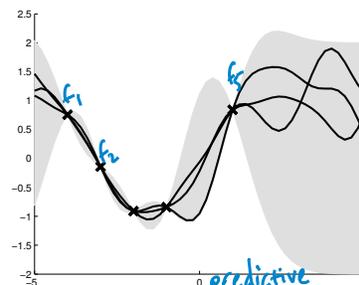
25

1D Noise-Free Example



Samples from Prior

$$\kappa(x, x') = \sigma_f^2 \exp\left(-\frac{1}{2\ell^2}(x - x')^2\right)$$



Posterior Given 5 Noise-Free Observations

- Interpolator, where uncertainty increases with distance
- Useful as a computationally cheap proxy for a complex simulator
 - Examine effect of simulator params on GP predictions instead of doing expensive runs of the simulator

©Emily Fox 2014

26

GPs for Regression

- Noisy scenario: observe a noisy version of underlying function

$$y = f(x) + \epsilon \quad \epsilon \sim N(0, \sigma_y^2)$$

- Not required to interpolate, just come "close" to observed data

$$\text{cov}(y|X) = \text{cov}(f) + \text{cov}(\epsilon) = K + \sigma_y^2 I_n \triangleq K_y$$

- Training data $\mathcal{D} = \{(x_i, y_i), i = 1, \dots, n\}$

- Test data locations $X^* \rightarrow$ predict f^*

- Jointly, we have $\begin{pmatrix} y \\ f^* \end{pmatrix} \sim N\left(0, \begin{pmatrix} K_y & K_* \\ K_y^T & K_{**} \end{pmatrix}\right)$

- Therefore, $p(f^* | X^*, X, y) = N(f^* | K_*^T K_y^{-1} y, K_{**} - K_*^T K_y^{-1} K_*)$

©Emily Fox 2014

27

GPs for Regression

$$p(f^* | X^*, X, y) = N(K_*^T K_y^{-1} y, K_{**} - K_*^T K_y^{-1} K_*)$$

- For a single point x^*

$$p(f^* | X^*, X, y) = N(k_*^T K_y^{-1} y, k_{**} - k_*^T K_y^{-1} k_*)$$

so

$$\bar{f}^* = k_*^T K_y^{-1} y = \sum_{i=1}^n \alpha_i K(x_i, x^*)$$

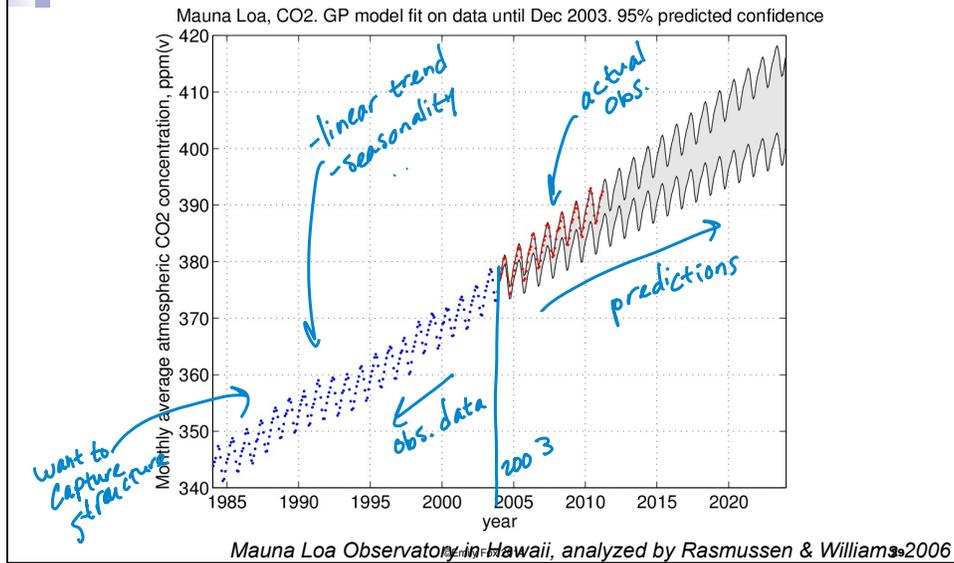
↑ predictive mean

will see this later

©Emily Fox 2014

28

CO2 Concentration Over Time



Mixing Kernels for CO2 GP Analysis

Smooth global trend

$$\kappa_1(x, x') = \theta_1^2 \exp\left(-\frac{(x - x')^2}{2\theta_2^2}\right)$$

Seasonal periodicity

$$\kappa_2(x, x') = \theta_3^2 \exp\left(-\frac{(x - x')^2}{2\theta_4^2} - \frac{2 \sin^2(\pi(x - x'))}{\theta_5^2}\right)$$

Medium term irregularities

$$\kappa_3(x, x') = \theta_6^2 \left(1 + \frac{(x - x')^2}{2\theta_8\theta_7^2}\right)^{-\theta_8}$$

Correlated Observation Noise

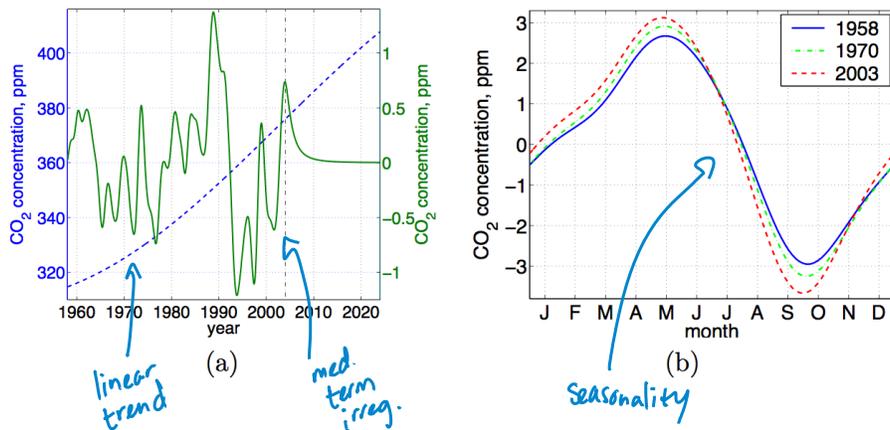
$$\kappa_4(x_p, x_q) = \theta_9^2 \exp\left(-\frac{(x_p - x_q)^2}{2\theta_{10}^2}\right) + \theta_{11}^2 \delta_{pq}$$

Fun fact:
- sum of kernels
is a kernel

- product of
kernels is
a kernel

$$k = k_1 + k_2 + k_3 + k_4$$

CO2 Concentration Over Time



hyperparams θ_j are optimized.

Mauna Loa Observatory in Hawaii, analyzed by Rasmussen & Williams, 2006

Estimating Hyperparameters

- How should we choose the kernel parameters?

- Example: squared exponential kernel parameterization

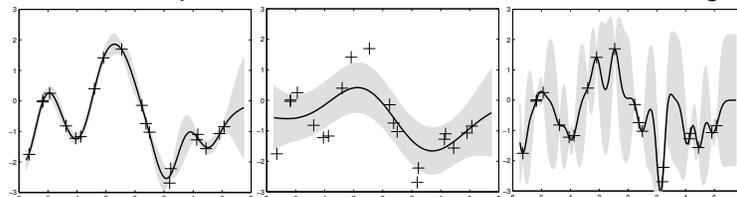
$$\kappa(x, x') = \sigma_f^2 \exp\left(\frac{-1}{2}(x_p - x_q)^T M (x'_p - x'_q)\right) + \sigma_y^2 \delta_{pq}$$

- Hyperparameters $\theta = \{M, \sigma_f^2, \sigma_y^2\}$

- As we saw before, can choose

$$M = \ell^{-2} I \quad M = \text{diag}(\ell_1^{-2}, \dots, \ell_d^{-2}) \quad M = \Lambda \Lambda' + \text{diag}(\ell_1^{-2}, \dots, \ell_d^{-2}) \dots$$

- As in other nonparametric methods, choice can have large effect



©Emily Fox 2014

32

Estimating Hyperparameters

Options:

- #1: Define a grid of possible values and use cross validation
can be slow...
- #2: Full Bayesian analysis: Place prior on hyperparameters and integrate over these as well in making predictions
some challenges in practice
- #3: Maximize the marginal likelihood

$$p(y | X, \theta) = \int p(y | f, X)p(f | X, \theta)df$$

$$\log p(y | X, \theta) =$$

Acknowledgements

*Many figures courtesy Kevin Murphy's textbook
[Machine Learning: A Probabilistic Perspective](#),
and Chris Bishop's textbook
[Pattern Recognition and Machine Learning](#)*

*Slides based on parts of the lecture notes of Erik Sudderth for
"Applied Bayesian Nonparametrics" at Brown University*