

## Module 3: Bayesian Nonparametrics

# Gaussian Processes for Regression

STAT/BIOSTAT 527, University of Washington

Emily Fox

April 21<sup>st</sup>-22<sup>nd</sup>, 2014

©Emily Fox 2014

1

## Recap of regression so far

- Recall our regression setting

$$f(x) = E[Y | x]$$

- How to estimate from finite training set?

Restrict to  
model class

- Example = linear basis expansion

□ Standard linear

□ Polynomial

□ Splines

□ ...

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \epsilon$$

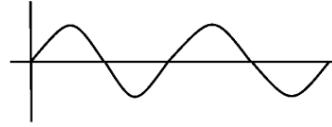
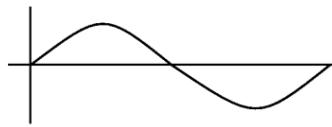
$y = \sum_j \beta_j x^j + \epsilon$

} good locally,  
but not  
globally

©Emily Fox 2014

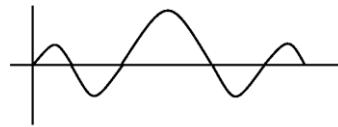
2

## Other Important Basis Expansions



•  
•

Fourier Basis



Wavelet Basis

not looking at these in this class

©Emily Fox 2014

3

## Recap of regression so far

- Recall our regression setting

$$f(x) = E[Y | x]$$

- How to estimate from finite training set?

Restrict to  
model class

- Example = linear basis expansion

Overfitting as model  
complexity grows

- Penalized linear basis expansions

(regularized LS)

□ Ridge

□ Lasso

□ Smoothing splines

□ Penalized regression splines

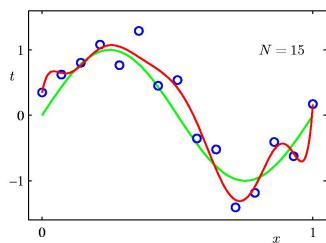
equivalent to searching over all functions  
subject to a smoothness constraints

©Emily Fox 2014

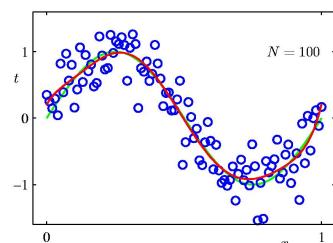
4

# Overfitting

9<sup>th</sup> Order Polynomial



$n = 15$



$n = 100$

- model complexity is relative to sample size
- can consider more complex forms with more data

©Emily Fox 2014

5

## Recap of regression so far

- Recall our regression setting

$$f(x) = E[Y | x]$$

- How to estimate from finite training set?

- Example = linear basis expansion
  - ↓ Overfitting as model complexity grows
- Penalized linear basis expansions

Restrict to  
model class

Local nbhd  
methods

- Example = kernel regression

kNN regression  
local averages  
Nadaraya-Watson  
locally weighted poly.

©Emily Fox 2014

6

## Again: Linear Basis Expansion

- Instead of just considering input variables  $x$  (potentially mult.), augment/replace with transformations = “input features”

In this lecture, we'll focus on these forms

- Linear basis expansions maintain linear form in terms of these transformations

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

trans.

- What transformations should we use?

- $h_m(x) = x_m \rightarrow$  linear model
- $h_m(x) = x_j^2, h_m(x) = x_j x_k \rightarrow$  polynomial reg.
- $h_m(x) = I(L_m \leq x_k \leq U_m) \rightarrow$  piecewise constant
- ...

©Emily Fox 2014

7

## Making Predictions

- So far, our focus has been on  $L_2$  loss:

$$\min_{\beta} \text{RSS}(\beta) + \lambda \|\beta\|$$
$$\sum_i (y_i - \hat{f}(x_i))^2 \quad \hat{f}(x) = \beta^T h(x)$$

plus penalty

- Here, we assumed  $y = f(x) + \epsilon$  with  $E[\epsilon] = 0 \quad \text{var}(\epsilon) = \sigma^2$

- Now, let's assume a distributional form and log-likelihood loss

$$y \sim N(0, \sigma^2) \Rightarrow p(y | f(x), \sigma^2) = N(f(x), \sigma^2)$$

First, recall some facts about Gaussians...

©Emily Fox 2014

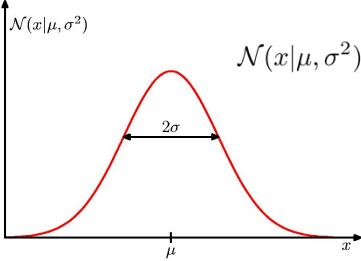
8

# Quick Review of Gaussians

- Univariate and multivariate Gaussians

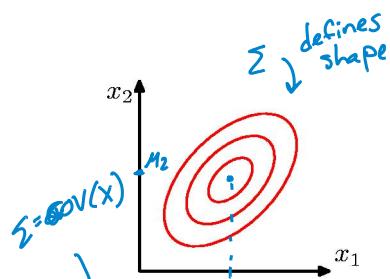
$$\mathcal{N}(x|\mu, \sigma^2)$$

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



$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

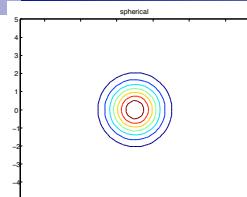
$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^D/2} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right\}$$



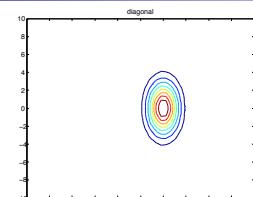
©Emily Fox 2014

9

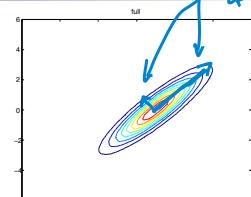
## Two-Dimensional Gaussians



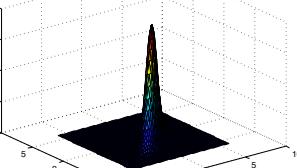
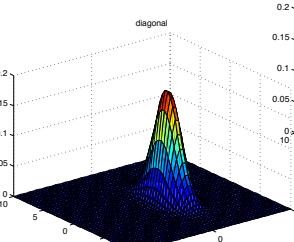
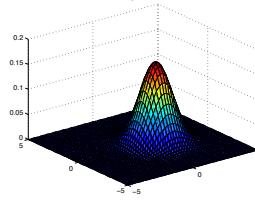
$$\boldsymbol{\Sigma} = \sigma^2 \mathbf{I}$$



$$\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_d^2 \end{bmatrix}$$



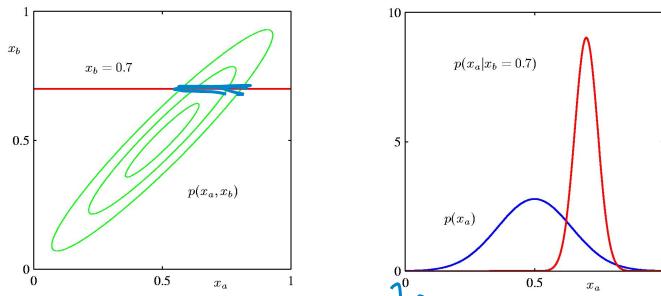
$$\begin{aligned} \text{eigvecs of } \boldsymbol{\Sigma} \\ \boldsymbol{\Sigma} \text{ general} \end{aligned}$$



©Emily Fox 2014

10

## Conditional & Marginal Distributions



$$\begin{pmatrix} x_a \\ x_b \end{pmatrix} \sim N \left( \begin{pmatrix} \mu_a \\ \mu_b \end{pmatrix}, \begin{bmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{bmatrix} \right)$$

Marg:  $x_a \sim N(\mu_a, \Sigma_{aa})$

Cond:  $x_a | x_b \sim N(\mu_a + \Sigma_{ab}\Sigma_{bb}^{-1}(x_b - \mu_b), \Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba})$

©Emily Fox 2014

11

## 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)  $(x_1, y_1), \dots, (x_n, y_n)$

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

©Emily Fox 2014

12

# Maximum Likelihood Estimation

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

- Taking the log

$$\log p(y | X, \beta, \sigma^2) \approx \sum_i -\frac{1}{2} (y_i - \beta^T h(x_i))^2 - \frac{n}{2} \log 2\pi - \frac{n}{2} \log \sigma^2$$

const. wrt  $\beta$

- Equivalent objective to RSS (*Gaussian log-like loss =  $L_2$  loss*)
- Taking the gradient and setting to zero, we have already shown

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

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

©Emily Fox 2014

13

## A Bayesian Formulation

*Look at penalized regression*

- Consider a model with likelihood

$$y_i | \beta \sim N(\beta_0 + x_i^T \beta, \sigma^2)$$

and prior

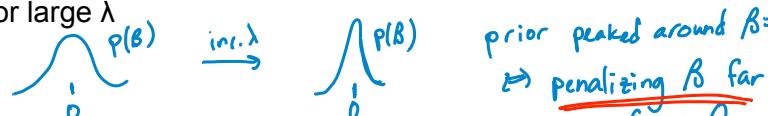
$$\beta \sim N\left(0, \frac{\sigma^2}{\lambda} I_p\right)$$

if  $\epsilon \sim N(0, \sigma^2)$

prior places penalty

$\beta_j \sim N(0, \frac{\sigma^2}{\lambda})$

- For large  $\lambda$



- The posterior is

$$\beta | y \sim N\left(\hat{\beta}^{ridge}, \sigma^2(X^T X + \lambda I)^{-1}\right)$$

$\hat{\beta}^{MAP} = \hat{\beta}^{ridge}$

works against overfitting of MLE ↑ easy to show  $\text{var}(\hat{\beta}^{ridge})$

©Emily Fox 2014

14

# Bayesian Linear Regression

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

$$p(\beta) = N(\beta | \mu_0, \Sigma_0)$$

*prior mean + cov.*

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

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

*posterior (updated)  
mean + cov.  
after obs. y*

where

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

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

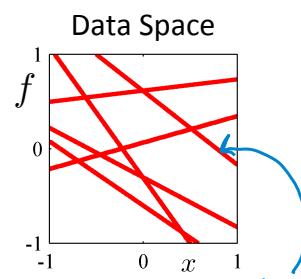
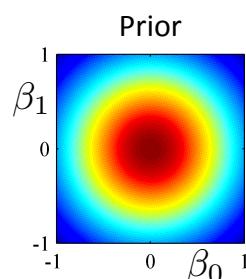
©Emily Fox 2014

15

## Example: Standard Linear Basis

0 data points observed

$$y = \beta_0 + \beta_1 x + \epsilon$$



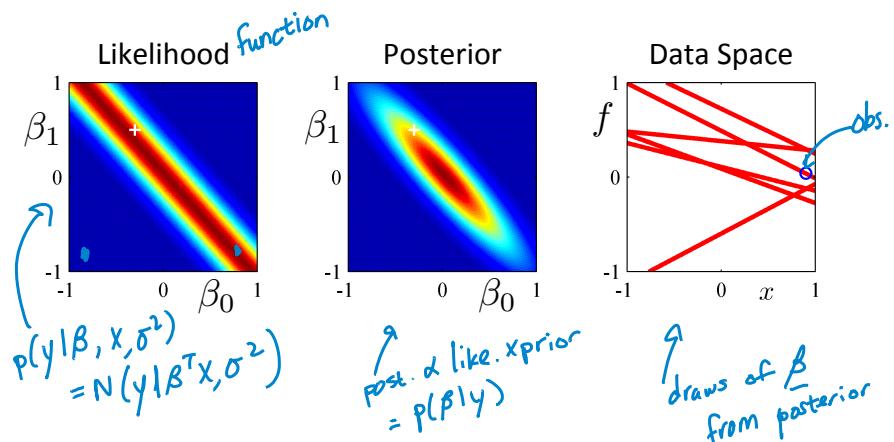
*draw  $\beta_0, \beta_1$  from prior  
and set  $f = \beta_0 + \beta_1 x$*

©Emily Fox 2014

16

## Example: Standard Linear Basis

1 data point observed

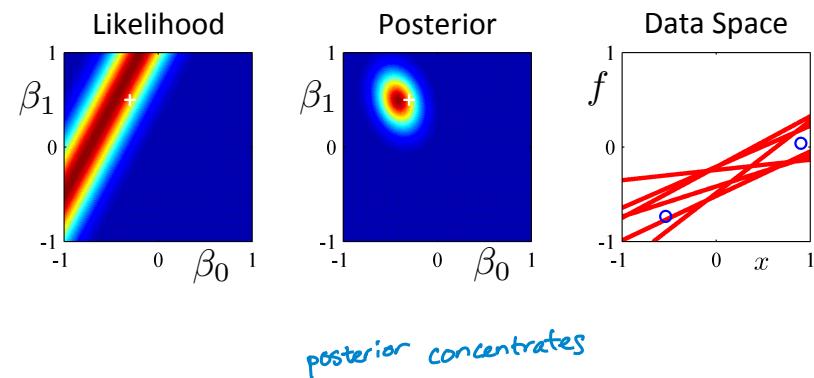


©Emily Fox 2014

17

## Example: Standard Linear Basis

2 data points observed

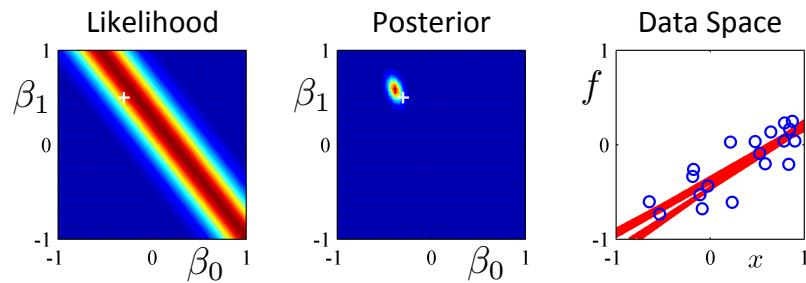


©Emily Fox 2014

18

## Example: Standard Linear Basis

20 data points observed



©Emily Fox 2014

19

## Predictive Distribution

- Predict  $y^*$  at new locations  $x^*$  by integrating over parameters  $\beta$

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

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

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

$y^* = h(x^*)^\top \beta + \epsilon$

$\beta | y \sim N(\mu_n, \Sigma_n)$

$\epsilon \sim N(0, \sigma^2)$

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

$\sum_n$  fan of  
our obs.  
locations  
 $x_1, \dots, x_n$

$$\begin{aligned} \mu_n^*(x^*) &= E[y^* | y] = \mu_n^\top h(x^*) \\ \Sigma_n^*(x^*) &= \text{cov}(y^* | y) = h(x^*)^\top \text{cov}(\beta | y) h(x^*) + \sigma^2 \\ &= h(x^*)^\top \Sigma_n h(x^*) + \sigma^2 \end{aligned}$$

$\var_{\text{of } \beta}$        $\var_{\text{of obs.}}$

©Emily Fox 2014

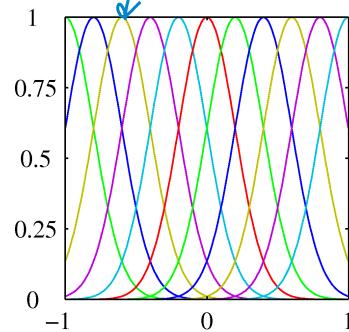
20

## Example: Gaussian Basis Expansion

- Gaussian basis functions:

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

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

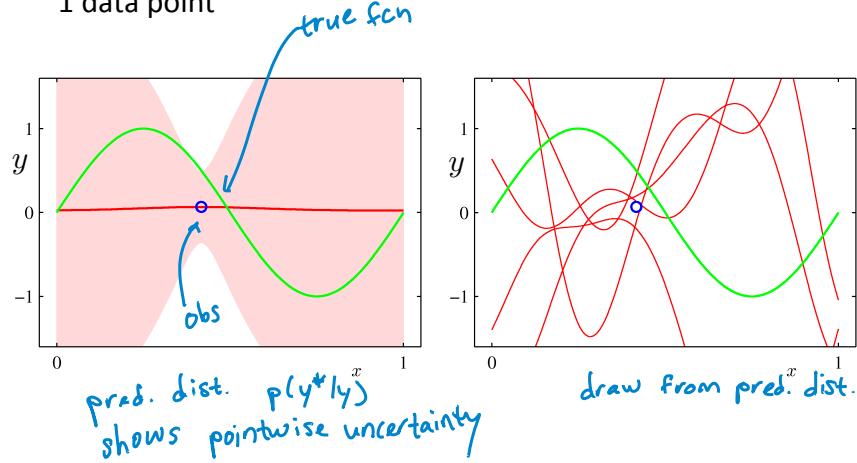


©Emily Fox 2014

21

## Example: Gaussian Basis Expansion

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

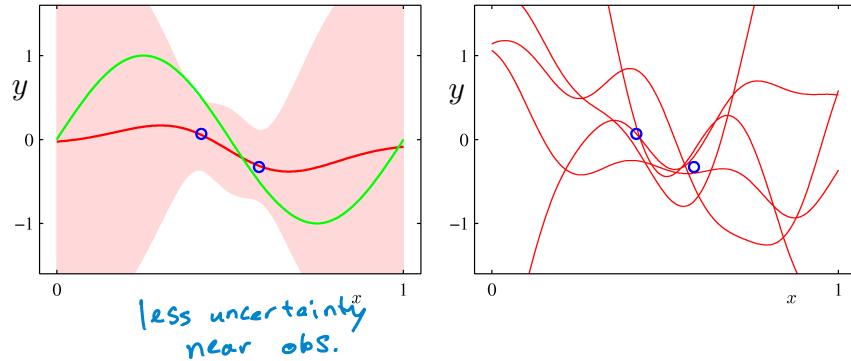


©Emily Fox 2014

22

## Example: Gaussian Basis Expansion

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

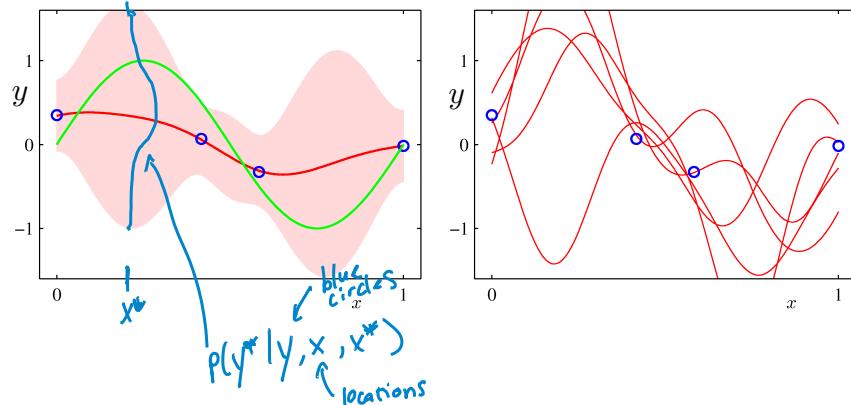


©Emily Fox 2014

23

## Example: Gaussian Basis Expansion

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

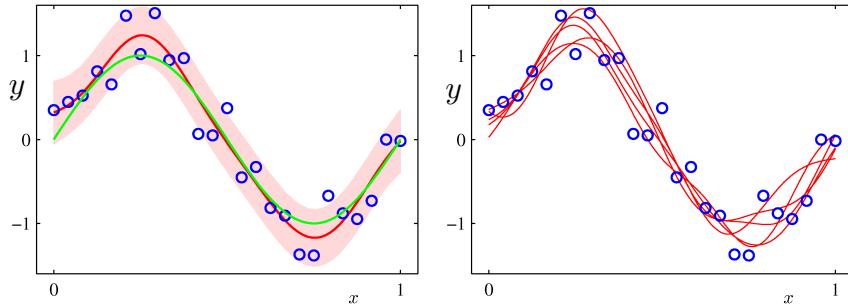


©Emily Fox 2014

24

## Example: Gaussian Basis Expansion

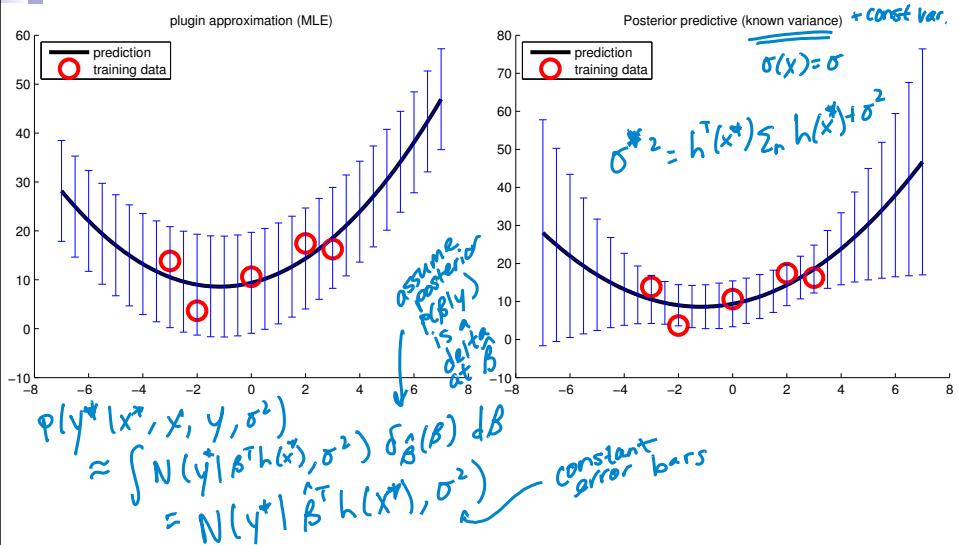
- Example: Sinusoidal data, 9 Gaussian basis functions, 25 data points



©Emily Fox 2014

25

## Estimation vs. Predictive Distributions



©Emily Fox 2014

26

# What # of basis fns should we use?

## Bayesian Model Selection

- Assume some  $M$  possible models
  - Model  $M_m$ ,  $m=1, \dots, M$  has parameters  $\theta_m$  and prior  $p(\theta_m | M_m)$
  - Prior over models  $p(M_m)$

- Model posterior

$$p(M_m | Z) \propto p(M_m) p(Z | M_m)$$

$$\propto p(M_m) \int p(Z | \theta_m, M_m) p(\theta_m | M_m) d\theta_m$$

eg mean + cov Gauss or  $\{\theta_j\}$

- Compare models:

$$\frac{p(M_m | Z)}{p(M_\ell | Z)} = \frac{p(M_m) p(Z | M_m)}{p(M_\ell) p(Z | M_\ell)} > 1$$

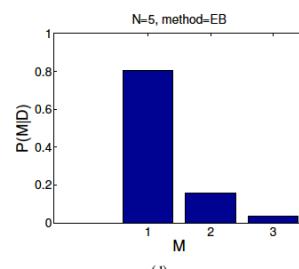
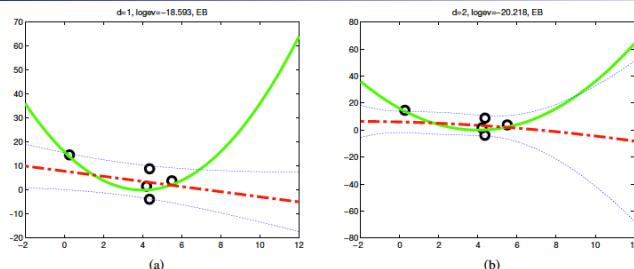
Often, uniform prior

Bayes factor "marginal likelihood"

©Emily Fox 2014

27

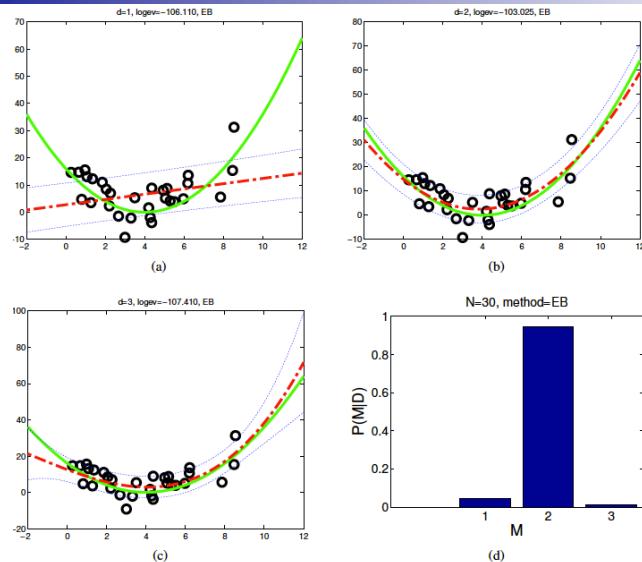
## BMS Example ( $n=5$ )



©Emily Fox 2014

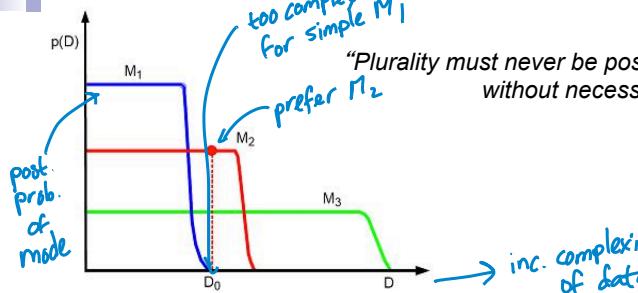
28

## BMS Example (n=30)



29

## Bayesian Ockham's Razor



- **Parametric Bayes:** Consider a finite list of possible models, average according to posterior probability (or in practice, just select the most probable)
- **Nonparametric Bayes:** Consider a single infinite model, integrate over parameters when making predictions or infer which finite subset is exhibited in your dataset

@9:00 am

$M \rightarrow \infty$

©Emily Fox 2014

30

# 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*