

Module 2: Splines and Kernel Methods

Inference for Linear Smoothers

STAT/BIOSTAT 527, University of Washington

Emily Fox

April 23rd, 2013

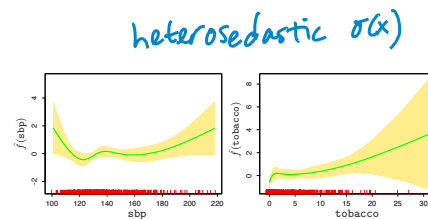
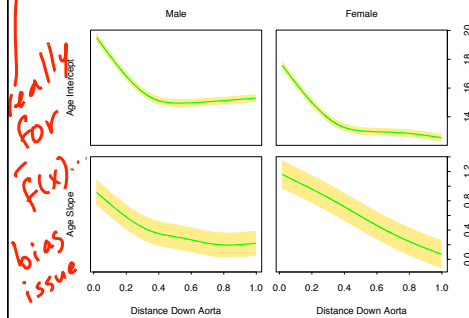
©Emily Fox 2013

1

Confidence Bands

- So far we have focused on point estimation: $\hat{f}(x)$
- Often, we want to define a **confidence interval** for which $f(x)$ is in this interval with some pre-specified probability
- Looking over all x , we refer to these as **confidence bands**

homoscedastic $\sigma(x) = \sigma$



From Hastie, Tibshirani, Friedman book

©Emily Fox 2013

2

CIs for Linear Smoothers

- For linear smoothers, and assuming constant variance $\sigma(x) = \sigma$

$$\hat{f}(x) = \sum_{i=1}^n \ell_i(x) y_i \quad \begin{array}{l} \nearrow \bar{f}(x) = \sum_{i=1}^n \ell_i(x) f(x_i) \\ \searrow \text{Var}(\hat{f}(x)) = \sigma^2 \|\ell(x)\|^2 \end{array}$$

- Consider confidence band of the form

$$CI(x) = \hat{f}(x) \pm c \hat{\sigma} \|\ell(x)\| \quad \begin{array}{l} a \leq x \leq b \\ \uparrow c > 0 \\ \uparrow \text{est. of } \sigma \end{array}$$

- Using this, let's solve for c

©Emily Fox 2013

3

CIs for Linear Smoothers

- Based on approach of Sun and Loader (1994)

- Case #2: Assume σ unknown use est. $\hat{\sigma}$?

- Case #3: Assume $\sigma(x)$ non-constant

$$\text{var}(\hat{f}(x)) = \sum_i \sigma^2(x_i) \ell_i^2(x)$$

$$CI(x) = \hat{f}(x) \pm c \sqrt{\sum_i \sigma^2(x_i) \ell_i^2(x)} \quad *$$

- If $\hat{\sigma}(x)$ varies slowly with x, then (Faraway and Sun 1995)

$$\begin{array}{l} \sigma(x_i) \approx \sigma(x) \text{ for those } x \text{ w/ } \ell_i(x) \text{ large} \\ \Rightarrow CI(x) = \hat{f}(x) \pm c \hat{\sigma}(x) \|\ell(x)\| \quad * \end{array}$$

©Emily Fox 2013

4

Variance Estimation

- In most cases σ is unknown and must be estimated
- For linear smoothers, consider the following estimator

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n (y_i - \hat{f}(x_i))^2}{n - 2\nu + \tilde{\nu}}$$

$\gamma = \text{tr}(L) \quad \tilde{\gamma} = \text{tr}(L^T L) = \sum_i \|\ell(x_i)\|^2$

- If target function is sufficiently smooth, $\nu = o(n)$, $\tilde{\nu} = o(n)$
- Then $\hat{\sigma}^2$ is a consistent estimator of σ^2

©Emily Fox 2013

5

Variance Estimation

- Proof outline:

- Recall that

$$Y - \hat{f} = Y - LY = (I - L)Y \triangleq \Lambda^{1/2} Y$$

and

$$E[Y^T Q Y] = \text{tr}(QV) + \mu^T Q \mu$$

- Then,

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n (y_i - \hat{f}(x_i))^2}{n - 2\nu + \tilde{\nu}} = \frac{Y^T \Lambda Y}{\text{tr}(\Lambda)}$$

$$E[\hat{\sigma}^2] = \frac{\text{tr}(\Lambda \sigma^2) + f^T \Lambda f}{\text{tr}(\Lambda)} = \sigma^2 + \frac{f^T \Lambda f}{n - 2\nu + \tilde{\nu}}$$

small for large n assuming f smooth

- Therefore, bias $\rightarrow 0$ for large n if f is smooth.
- Likewise for variance.

©Emily Fox 2013

6

Alternative Estimator

- Estimator:

$$\hat{\sigma}^2 = \frac{1}{2(n-1)} \sum_{i=1}^{n-1} (y_{i+1} - y_i)^2$$

- Motivation:

$$y_{i+1} - y_i = [f(x_{i+1}) - f(x_i)] + [\epsilon_{i+1} - \epsilon_i]$$

$$E[(y_{i+1} - y_i)^2] \approx E[\epsilon_{i+1}^2] + E[\epsilon_i^2] = 2\sigma^2$$

$$\Rightarrow E[\hat{\sigma}^2] \approx \sigma^2$$

- Estimator will be inflated *ignores $f(x_{i+1}) - f(x_i)$*
- Other estimators exist, too. See Wakefield or Wasserman.

©Emily Fox 2013

7

Heteroscedasticity



- The point estimate $\hat{f}(x)$ is relatively insensitive to heterosced., but confidence bands need to account for non-constant variance

- Re-examine model $y_i = f(x_i) + \sigma(x_i)\epsilon_i$ *$E[\epsilon] = 0$ var(1)*

- Define *redefine obs:*

$$Z_i = \log(y_i - f(x_i))^2 \quad \delta_i = \log \epsilon_i^2$$

- Then,

$$Z_i = \log(\sigma^2(x_i)) + \delta_i$$

- Algorithm:

- Estimate $f(x)$ using a nonparametric method w/ constant var to get $\hat{f}(x)$
- Define $Z_i = \log(y_i - \hat{f}(x_i))^2$
- Regress Z_i 's on x_i 's to get estimate $\hat{g}(x)$ of $\log \sigma^2(x)$ *new obs.*

$$\Rightarrow \hat{\sigma}^2(x) = e^{\hat{g}(x)}$$

$$Z_i = g(x_i) + \delta_i$$

$\log(\sigma^2(x_i))$

©Emily Fox 2013

8

Heteroscedasticity

- Drawbacks:

- Taking log of a very small residual leads to a large outlier
- A more statistically rigorous approach is to jointly estimate f, g

- Alternative = Generalized linear models

©Emily Fox 2013

9

Module 3: Bayesian Nonparametrics

Gaussian Processes

STAT/BIOSTAT 527, University of Washington

Emily Fox

April 23rd, 2013

©Emily Fox 2013

10

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

- ☐ ...

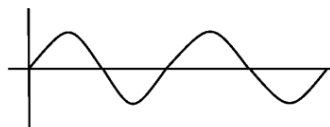
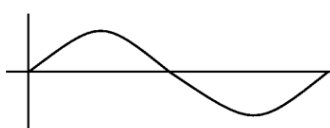
$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \epsilon$$
$$y = \sum \beta_j x^j + \epsilon$$

} good locally,
but not globally

©Emily Fox 2013

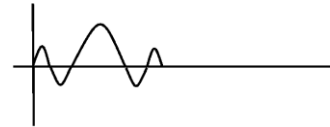
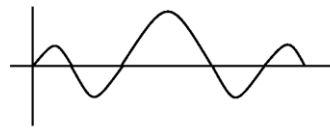
11

Other Important Basis Expansions



•
•
•

Fourier Basis



Wavelet Basis

not looking at these in this class

©Emily Fox 2013

12

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

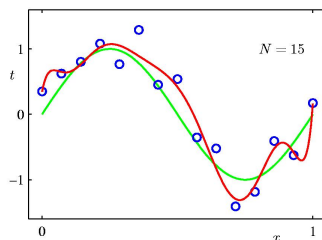
- Penalized linear basis expansions (regularized LS)
 - Ridge
 - Lasso
 - Smoothing splines
 - Penalized regression splines
- equivalent to searching over all fens s.t. smoothness constraints*

©Emily Fox 2013

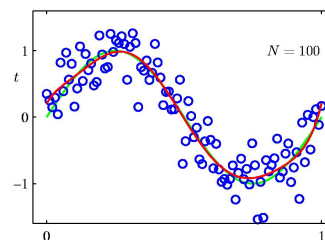
13

Overfitting

9th Order Polynomial



$n = 15$



$n = 100$

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

©Emily Fox 2013

14

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

- Example = kernel regression

k-NN regression
local averages
Nadaraya-Watson
local weighted poly.

©Emily Fox 2013

15

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

©Emily Fox 2013

16

Making Predictions

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

$$\min_{\beta} \text{RSS}(\beta) + \lambda \|\beta\|$$

$\sum_i (y_i - \hat{f}(x_i))^2$ $f(x) = \beta^T h(x)$

- Here, we assumed $y = f(x) + \epsilon$ with $E[\epsilon] = 0$, $\text{var}(\epsilon) = \sigma^2$
- Now, let's ^{further} assume a distributional form and log-likelihood loss

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

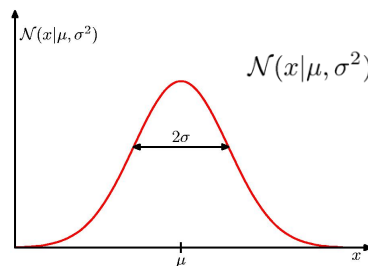
First, recall some facts about Gaussians...

©Emily Fox 2013

17

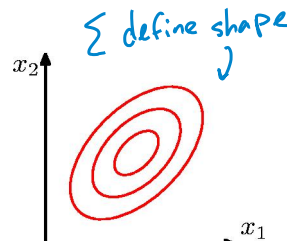
Quick Review of Gaussians

- Univariate and multivariate Gaussians



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

$\Sigma = \text{cov}(x)$

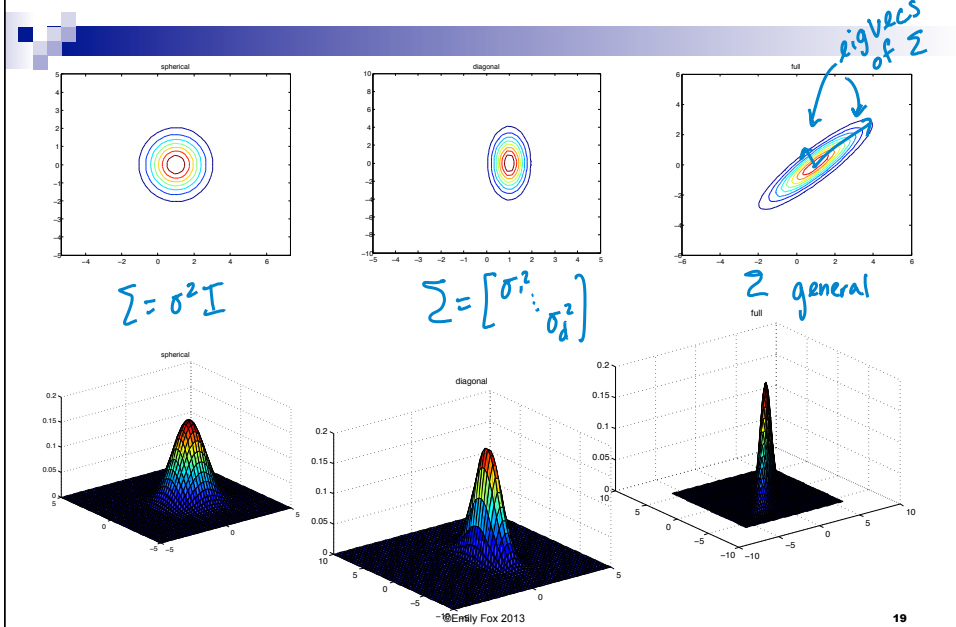


$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right\}$$

©Emily Fox 2013

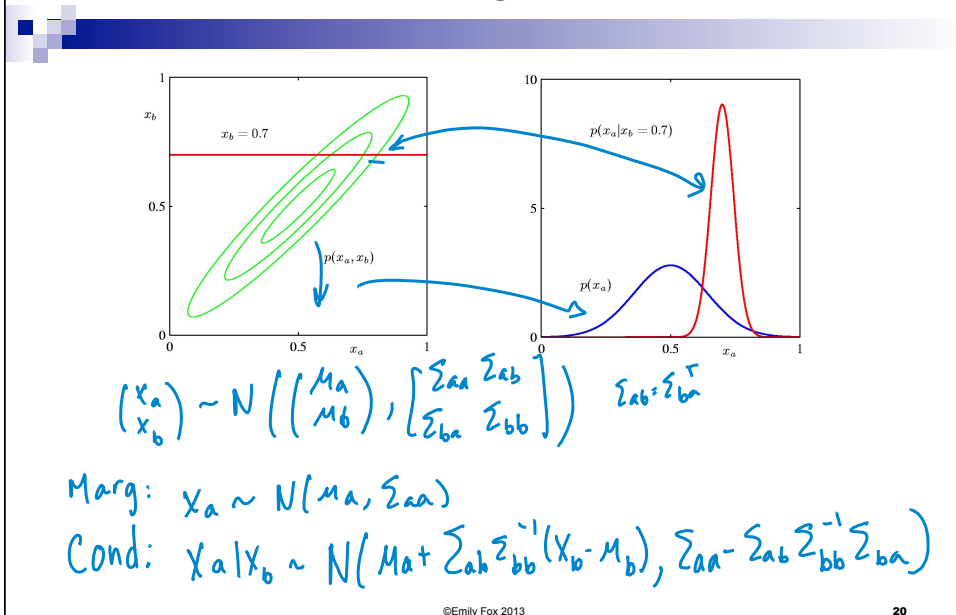
18

Two-Dimensional Gaussians



19

Conditional & Marginal Distributions



20

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 2013

21

Maximum Likelihood Estimation

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

- Taking the log

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

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

Const. wrt β

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

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

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

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

©Emily Fox 2013

22

A Bayesian Formulation of Ridge

- Consider a model with likelihood

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

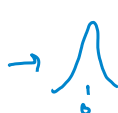
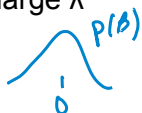
and prior

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

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

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

- For large λ



prior is peaked around $\beta=0$

\Rightarrow penalizing β far from 0

- The posterior is

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

works against overfitting of MLE

easy to show
 $\text{Var}(\hat{\beta}^{ridge})$

©Emily Fox 2013

23

Bayesian Linear Regression

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

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

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

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

posterior \propto likelihood \times prior

where

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

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

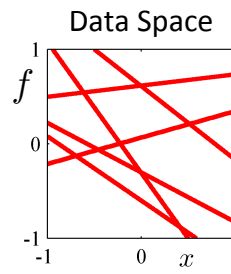
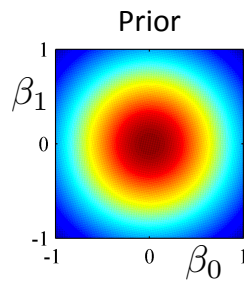
©Emily Fox 2013

24

Example: Standard Linear Basis

0 data points observed

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



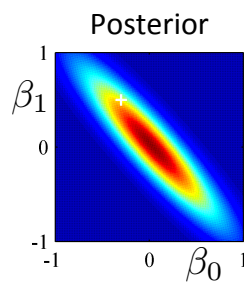
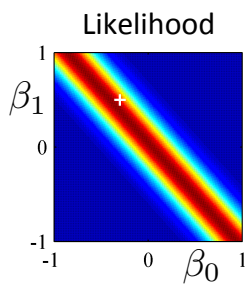
draw β_0, β_1 from prior
and set $f = \beta_0 + \beta_1 x$

©Emily Fox 2013

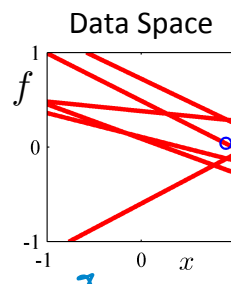
25

Example: Standard Linear Basis

1 data point observed



post. \propto like. \times prior



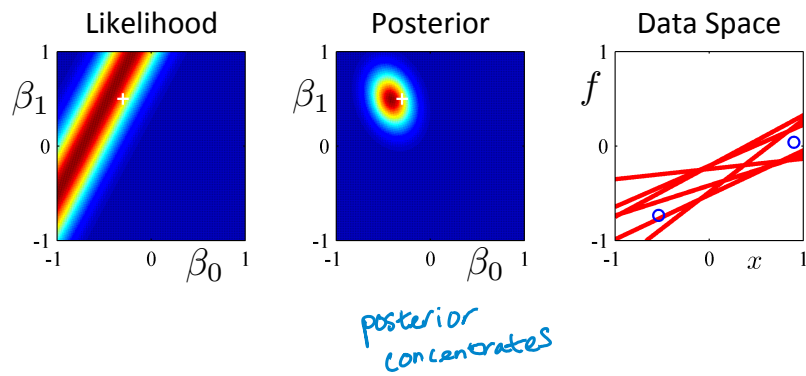
draws of β
from posterior

©Emily Fox 2013

26

Example: Standard Linear Basis

2 data points observed

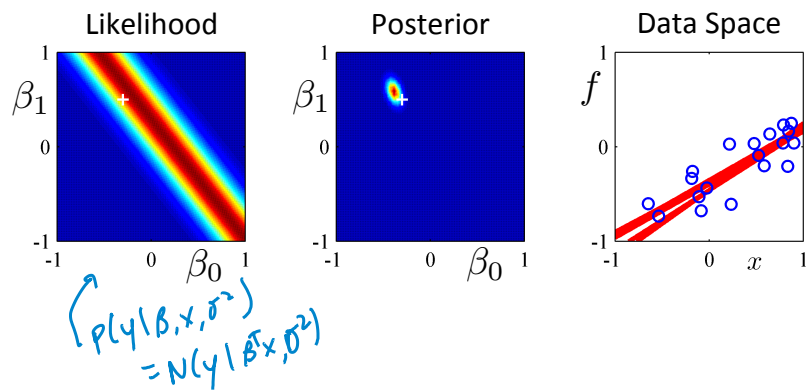


©Emily Fox 2013

27

Example: Standard Linear Basis

20 data points observed



©Emily Fox 2013

28

Predictive Distribution

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

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

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

$p(\beta | y) = N(\beta | \mu_n, \Sigma_n)$ (posterior)
 $p(y | x, \beta, \sigma^2) = N(y | f(x), \sigma^2)$
 $f(x) = \beta^T h(x)$

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

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

Var of our params β (points to Σ_n)
 Var of obs (points to σ^2)
 Var of locations x (points to $h(x^*)$)

©Emily Fox 2013

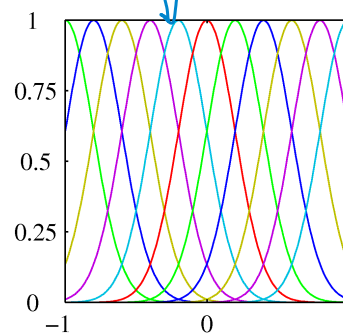
29

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)

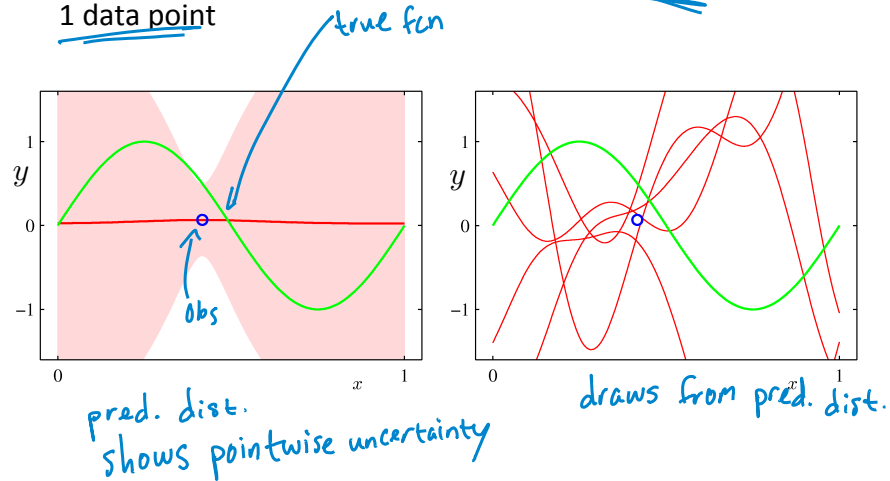


©Emily Fox 2013

30

Example: Gaussian Basis Expansion

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

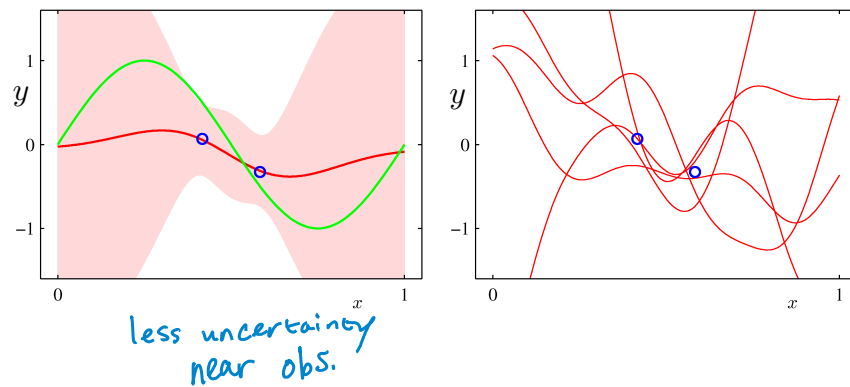


©Emily Fox 2013

31

Example: Gaussian Basis Expansion

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

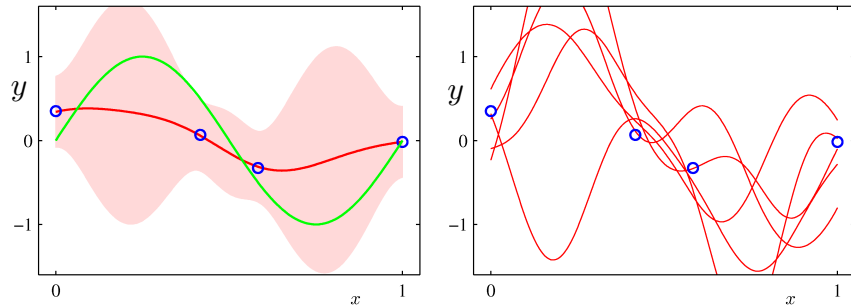


©Emily Fox 2013

32

Example: Gaussian Basis Expansion

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

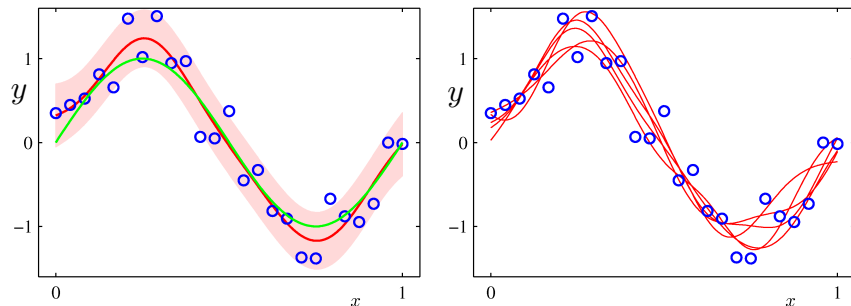


©Emily Fox 2013

33

Example: Gaussian Basis Expansion

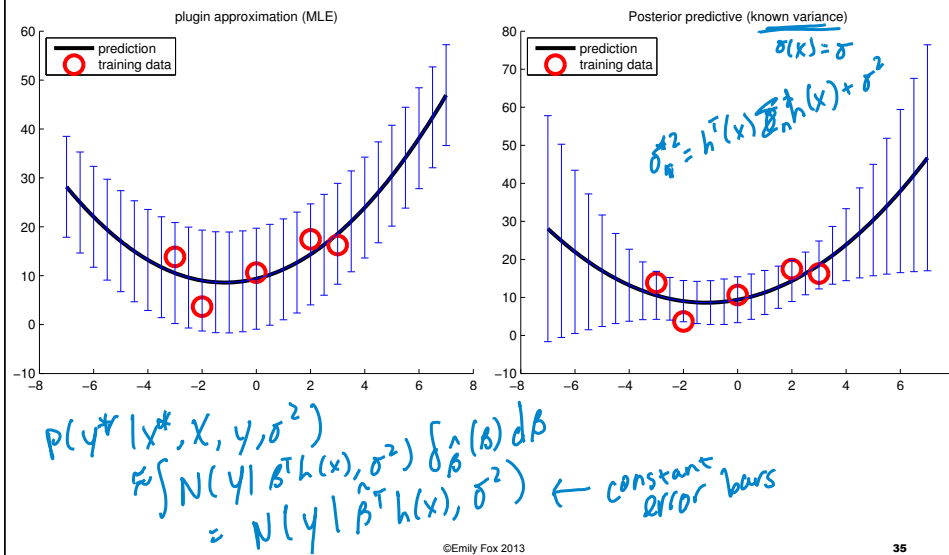
- Example: Sinusoidal data, 9 Gaussian basis functions, 25 data points



©Emily Fox 2013

34

Estimation vs. Predictive Distributions



What # of basis fns should we use? Bayesian Model Selection

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

- Model posterior

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

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

- Compare models:

post. odds

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

often, uniform prior

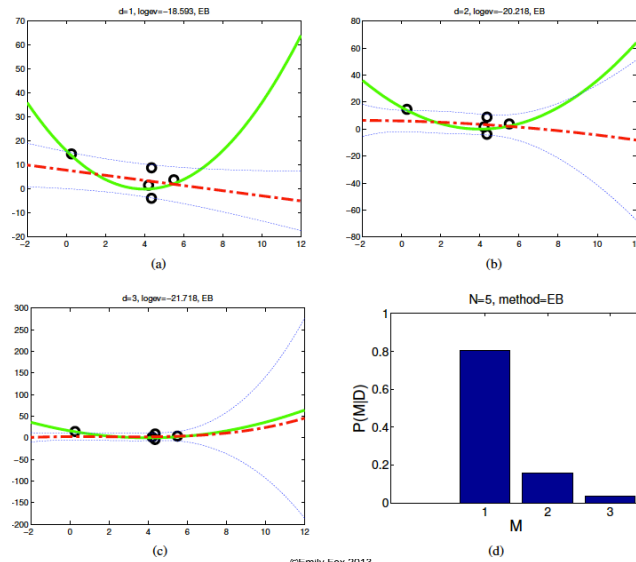
Bayes factor

©Emily Fox 2013

36

BMS Example (n=5)

uniform priors
on models

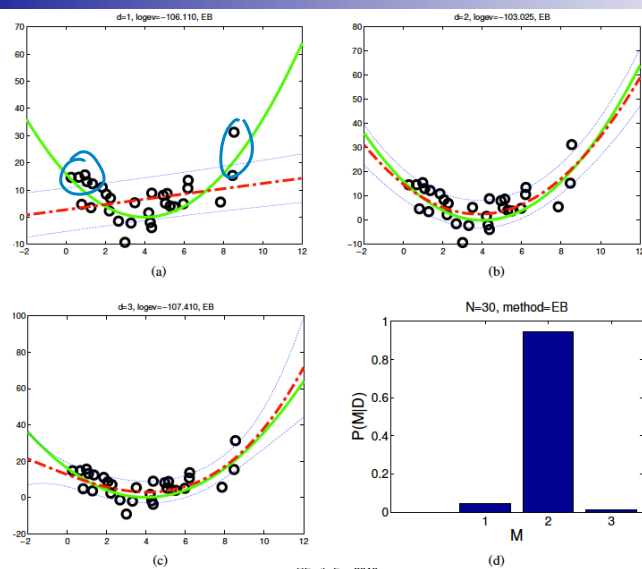


high
posterior
prob.
on simple
models

©Emily Fox 2013

37

BMS Example (n=30)

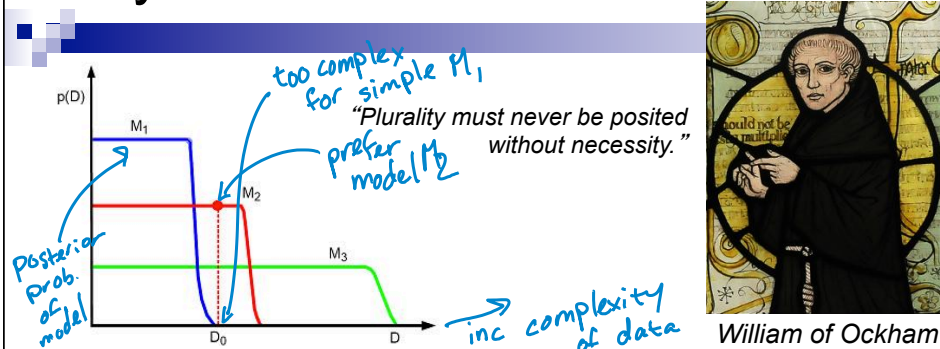


favors
slightly
more
complex
models

©Emily Fox 2013

38

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

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

Announcements

- Upcoming changes...
- Lectures:
 - Instead of lecture next Tuesday, Shirley will provide an examples section
 - Instead of recitation on Tuesday May 9, I will do a lecture on nonparametrics for generalized linear models (GLM)
- Homeworks:
 - Starting this Thursday, homeworks will be 2 weeks long
 - Provides extra flexibility on timing to accommodate project
 - Each homework (HW4 and HW5) will count the same as two 1-wk assignments
 - Should be slightly shorter than two 1-wk assignments