

# Designing Supercomputer Experiments

## Trinity Mathematical Society

Trinity College, March 2008

Robert B. Gramacy

Statistical Laboratory, University of Cambridge

<http://www.statslab.cam.ac.uk/~bobby>

[bobby@statslab.cam.ac.uk](mailto:bobby@statslab.cam.ac.uk)

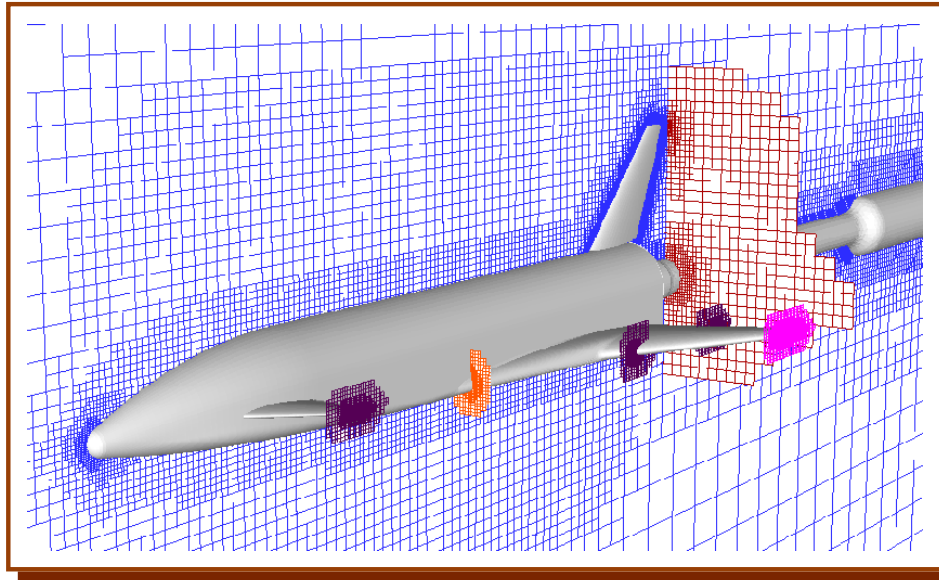
with Herbert K. H. Lee

Applied Math & Statistics Dept., UC Santa Cruz

# DESIGNING COMPUTER EXPERIMENTS

Real-world Application: CFD simulations of a proposed reusable NASA launch vehicle (Langley-Glide-Back Booster)

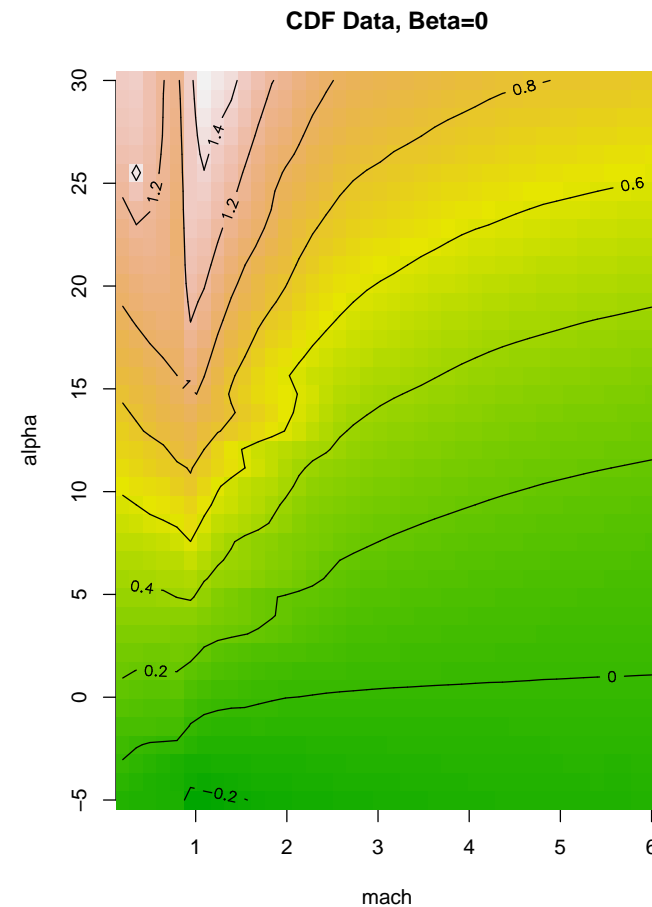
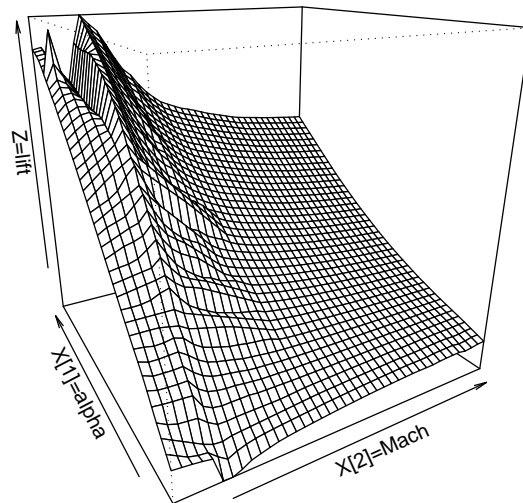
3 inputs:  
side slip angle  
Mach number  
angle of attack



6 outputs:  
lift  
drag  
pitching moment  
side-force  
yawing moment  
rolling moment

- integrate inviscid Euler equations over mesh of  $1.4 \times 10^6$  cells
- each input configuration takes 5-20 hours of CPU time

# MOTIVATION: EXAMPLE



- Exhibits typical characteristics of a computer experiment
  - not: stationary, linear, differentiable, continuous

# AUGMENTING THE STATE OF THE ART

---

Canonical model: Gaussian process (Santner, et al. 2003)

- conceptually simple non-parametric extension of LM
- flexible & non-linear

But:

- inference scales poorly with number of data points  $N$
- strictly stationary

# AUGMENTING THE STATE OF THE ART

---

Canonical model: Gaussian process (Santner, et al. 2003)

- conceptually simple non-parametric extension of LM
- flexible & non-linear

But:

- inference scales poorly with number of data points  $N$
- strictly stationary

Our idea: use trees to partition the input space like Bayesian treed LMs (Chipman et al., 2002) with GPs instead of LMs

- allows for modeling of non-stationary behavior
- predictive variance is location (region) dependent
- smaller covariance matrices ameliorate computational demands

# STANDARD LINEAR MODEL (LM)

---

LM for  $n$  inputs (covariates)  $\mathbf{x}$  of dimension  $p$ , and responses  $y$

$$y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + e_i, \quad \text{where } e_i \sim N(0, \sigma^2), \quad \text{for } i = 1, \dots, n$$

- Collect  $Y = (y_1, \dots, y_n)^\top$ , and  $\mathbf{X} = [\mathbf{x}_1^\top \cdots \mathbf{x}_n^\top]$  and

define the *likelihood*:  $Y \sim N_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I})$

Jeffrey's prior:  $(\boldsymbol{\beta}, \sigma^2) \propto 1/\sigma^2$

MLE:  $\hat{\boldsymbol{\beta}} \sim (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top Y$ , &  $s^2 \sim \frac{1}{n-p} (Y - \mathbf{X}\hat{\boldsymbol{\beta}})^\top (Y - \mathbf{X}\hat{\boldsymbol{\beta}})$

- *Bayes Rule*  $p(\boldsymbol{\theta}|Y) \propto p(Y|\boldsymbol{\theta})p(\boldsymbol{\theta})$  gives the posterior:

$$\boldsymbol{\beta}|\sigma^2, Y = N_p(\hat{\boldsymbol{\beta}}, (\mathbf{X}^\top \mathbf{X})^{-1}\sigma^2), \quad \& \quad \sigma^2|Y = \text{Inv-}\chi^2(n-p, s^2)$$

# GAUSSIAN PROCESS (GP)

---

The Gaussian Process is *also* a LM (of sorts)

- Replacing  $\mathbf{I}$  with  $\mathbf{K}$

$$Y|\boldsymbol{\beta}, \sigma^2 \sim N_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I}) \quad \Rightarrow \quad Z|\boldsymbol{\beta}, \sigma^2, \mathbf{K} \sim N_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{K})$$

- $\mathbf{K}_{i,j} = K(\mathbf{x}_i, \mathbf{x}_j)$  is, e.g., an isotropic correlation matrix

$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp \left\{ -\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{d} \right\} + g\delta_{j,j}$$

- with *range* (or length-scale) parameter  $d$
- and *nugget* (or measurement-error) parameter  $g$ 
  - Range  $d \rightarrow 0$  gives  $\mathbf{K} \rightarrow (1 + g)\mathbf{I}$ ; i.e., the LM

- Conditional on  $\mathbf{K}$  through  $d$  and  $g$ , we have MLEs

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{K}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{K}^{-1} Z,$$

and 
$$s^2 = \frac{1}{n - p} (Z - \mathbf{X} \hat{\boldsymbol{\beta}})^\top \mathbf{K}^{-1} (Z - \mathbf{X} \hat{\boldsymbol{\beta}})$$

- Numerical optimization required for  $d$  and  $g$

Bayesian inference is similar

$$\boldsymbol{\beta} | \sigma^2, d, g, Z \sim N_p(\hat{\boldsymbol{\beta}}, (\mathbf{X}^\top \mathbf{K}^{-1} \mathbf{X})^{-1} \sigma^2)$$

$$\sigma^2 | d, g, Z \sim \text{Inv-}\chi^2(n - p, s^2)$$

- Metropolis–Hastings is used to sample the posterior  $p(d, g | Z)$



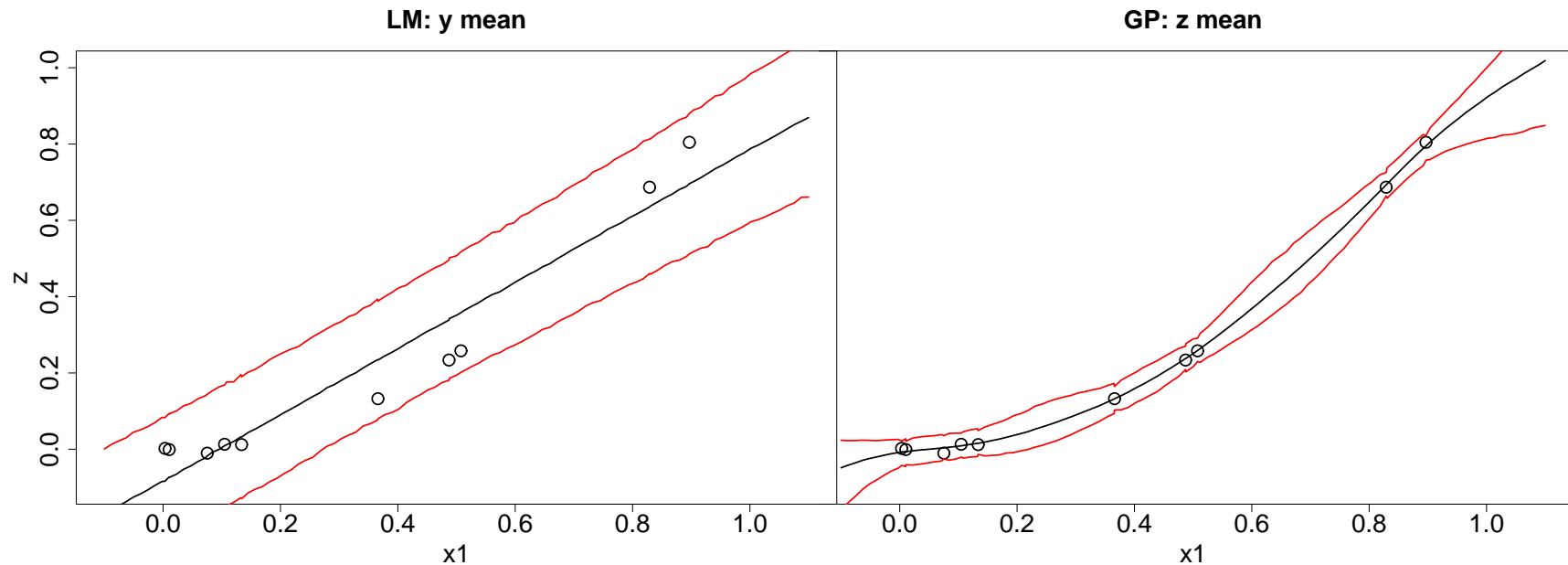
# GP PREDICTION (OR KRIGING)

The predicted value of  $z(\tilde{\mathbf{x}})$  is Normal with

mean  $\hat{y}(\tilde{\mathbf{x}}) = \tilde{\mathbf{x}}^\top \hat{\boldsymbol{\beta}} + \tilde{\mathbf{k}}(\tilde{\mathbf{x}})^\top \mathbf{K}^{-1}(\mathbf{Z} - \mathbf{X}\hat{\boldsymbol{\beta}})$

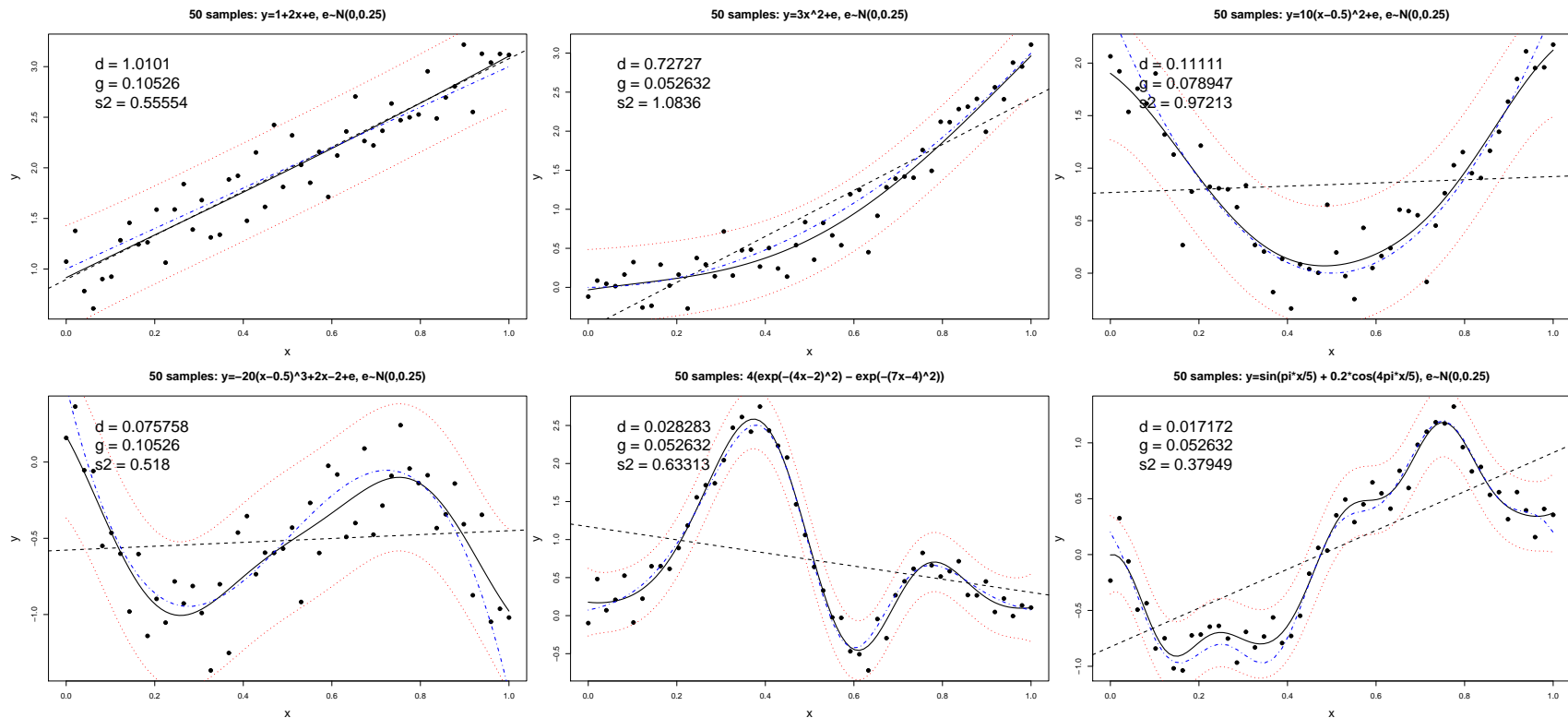
variance  $\hat{\sigma}(\tilde{\mathbf{x}})^2 = \sigma^2[\kappa(\tilde{\mathbf{x}}) - \mathbf{q}(\tilde{\mathbf{x}})^\top (\mathbf{K} + \mathbf{X}\mathbf{X}^\top)^{-1} \mathbf{q}(\tilde{\mathbf{x}})]$

where  $\kappa(\tilde{\mathbf{x}}) = K(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}) + \tilde{\mathbf{x}}^\top \tilde{\mathbf{x}}$ ,  $\mathbf{q}(\tilde{\mathbf{x}}) = \mathbf{k}(\tilde{\mathbf{x}}) + \mathbf{X}\tilde{\mathbf{x}}$ ,  $\mathbf{k}_j(\tilde{\mathbf{x}}) = K(\tilde{\mathbf{x}}, \mathbf{x}_j)$ , and  $\mathbf{x}_j^\top$  is the  $j^{\text{th}}$  column of  $\mathbf{X}$



# THE GP RANGE PARAMETER

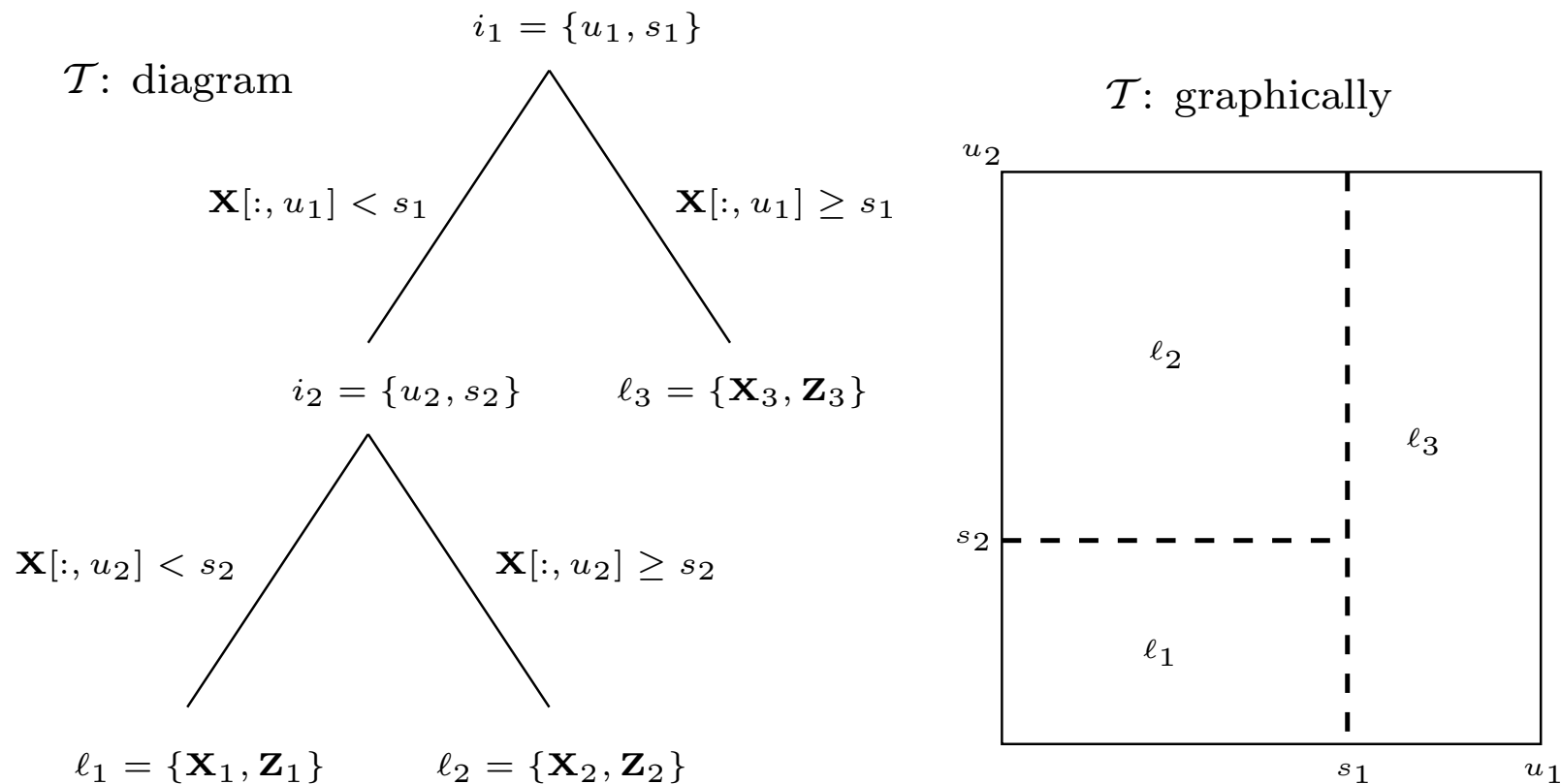
MLE range ( $d$ ) for various data:



- However, recall that  $d \rightarrow 0$  gives  $\mathbf{K} \rightarrow (1 + g)\mathbf{I}$

# TREE EXAMPLE

How a tree  $\mathcal{T}$  recursively partitions the input space:



- prior for node  $\eta$  at depth  $q_\eta \in \mathcal{T}$ :  $p_{\text{SPLIT}}(\eta, \mathcal{T}) = a(1 + q_\eta)^{-b}$

# HIERARCHICAL MODEL

---

Conditioning on tree  $\mathcal{T}$  we have  $R$  regions:  $\{r_\nu\}_{\nu=1}^R$

$$\begin{aligned}\mathbf{Z}_\nu | \boldsymbol{\beta}_\nu, \sigma_\nu^2, \mathbf{K}_\nu &\sim N(\mathbf{F}_\nu \boldsymbol{\beta}_\nu, \sigma_\nu^2 \mathbf{K}_\nu) & \sigma_\nu^2 &\sim IG(\alpha_\sigma/2, q_\sigma/2) \\ \boldsymbol{\beta}_\nu | \sigma_\nu^2, \tau_\nu^2, \mathbf{W}, \boldsymbol{\beta}_0 &\sim N(\boldsymbol{\beta}_0, \sigma_\nu^2 \tau_\nu^2 \mathbf{W}) & \tau_\nu^2 &\sim IG(\alpha_\tau/2, q_\tau/2) \\ \boldsymbol{\beta}_0 &\sim N(\boldsymbol{\mu}, \mathbf{B}) & \mathbf{W}^{-1} &\sim W((\rho \mathbf{V})^{-1}, \rho)\end{aligned}$$

with  $\mathbf{F}_\nu = (\mathbf{1}, \mathbf{X}_\nu)$ , and

$$K_\nu(\mathbf{x}_j, \mathbf{x}_k) = K_\nu^*(\mathbf{x}_j, \mathbf{x}_k) + g_\nu \delta_{j,k}$$

for nugget  $g$  and *true* correlation  $K^*$  is separable:

$$K_\nu^*(\mathbf{x}_j, \mathbf{x}_k | \mathbf{d}_\nu) = \exp\left\{-\sum_{i=1}^{m_X} |x_{ij} - x_{ik}|^{p_0} / d_{i\nu}\right\}$$

Conditional on a particular tree  $\mathcal{T}$

- Gibbs samples for all GP parameters  $\{\boldsymbol{\theta}_\nu\}_{\nu=1}^R | \mathcal{T}$ 
  - except  $K(\cdot, \cdot | \mathbf{d}_\nu, g_\nu)$  requires MH

Sample from the joint posterior of  $(\mathcal{T}, \boldsymbol{\theta})$

(Richardson & Green, 1997; Chipman et al., 2002)

- Average over  $\mathcal{T}$  with reversible-jump MCMC (RJ-MCMC)
  - Tree operations: *grow, prune, change, swap*

# BAYESIAN TREE *grow & prune*

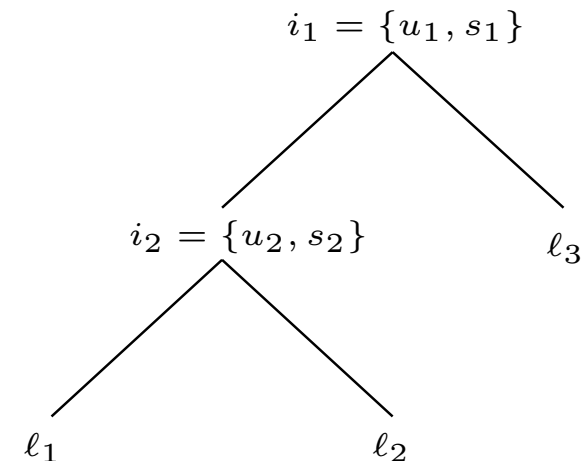
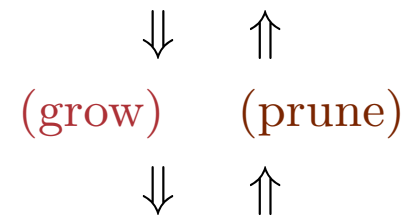
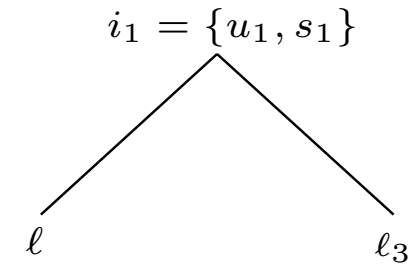
There is always a *leaf* to *grow*  $\mathcal{T}_t$ :

- randomly select a *leaf* node  $\eta \in \ell(\mathcal{T}_t)$
- randomly select a splitting rule  $(u, s)$
- create  $\mathcal{T}^*$  from  $\mathcal{T}_t$ 
  - split  $\eta$  in dimension  $u$  at location  $s$
  - $\eta \in \ell(\mathcal{T}_t)$  becomes  $(\eta_1, \eta_2) \in \ell(\mathcal{T}^*)$
- accept or reject the proposed  $\mathcal{T}^*$ 
  - via the MH acceptance ratio:

$$\alpha = \min \left\{ 1, \frac{p(\mathcal{T}^*|Y)q(\mathcal{T}|\mathcal{T}^*)}{p(\mathcal{T}_t|Y)q(\mathcal{T}^*|\mathcal{T})} \right\}$$

- set  $\mathcal{T}_{t+1} = \mathcal{T}^*$  w.p.  $\alpha$ ; or  $\mathcal{T}_{t+1} = \mathcal{T}_t$

The opposite of *grow* is *prune*  $\mathcal{T}_t$ :



# BAYESIAN TREE *change & swap*

Splits are moved with *change* operations

- randomly select *internal* node  $\eta \in i(\mathcal{T})$ 
  - $\eta$  has split point  $(u, s)$ , say
- propose a new  $\eta^* \in \mathcal{T}^*$ 
  - by modifying the split location  $(u, s^*)$

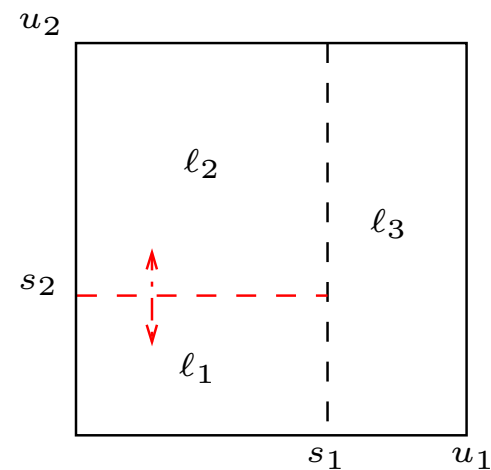
Or, switch the order of splits with *swap*

- randomly select a pair of *internal* nodes
  - $\eta_1, \eta_2 \in i(\mathcal{T})$  so that  $\mathcal{P}(\eta_2) = \eta_1$
  - swap their order, thus proposing  $\mathcal{T}^*$

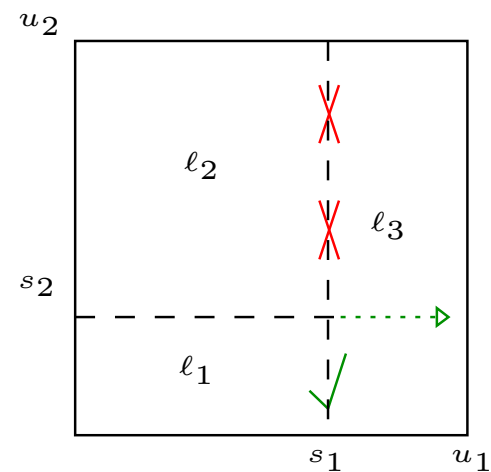
Accept or reject  $\mathcal{T}^*$  with probability

$$\alpha = \min \left\{ 1, \frac{p(\mathcal{T}^*|Y)}{p(\mathcal{T}_t|Y)} \right\}$$

(change)



(swap)



# PREDICTION (OR KRIGING)

The predicted value of  $Y(\mathbf{x})$  at  $\mathbf{x} \in r_\nu$  is Normal with

$$\begin{aligned} \text{mean} \quad & \hat{y}(\mathbf{x}) = \mathbf{f}^\top(\mathbf{x})\tilde{\boldsymbol{\beta}}_\nu + \mathbf{k}_\nu(\mathbf{x})^\top \mathbf{K}_\nu^{-1}(\mathbf{Z}_\nu - \mathbf{F}_\nu\tilde{\boldsymbol{\beta}}_\nu) \\ \text{variance} \quad & \hat{\sigma}(\mathbf{x})^2 = \sigma_\nu^2[\kappa(\mathbf{x}, \mathbf{x}) - \mathbf{q}_\nu^\top(\mathbf{x})\mathbf{C}_\nu^{-1}\mathbf{q}_\nu(\mathbf{x})] \end{aligned}$$

where  $\mathbf{k}_\nu(\mathbf{x})$  is a  $n_\nu$ -vector with  $\mathbf{k}_{\nu,j}(\mathbf{x}) = K_\nu(\mathbf{x}, \mathbf{x}_j)$ ,  $\mathbf{x}_j \in \mathbf{X}_\nu$

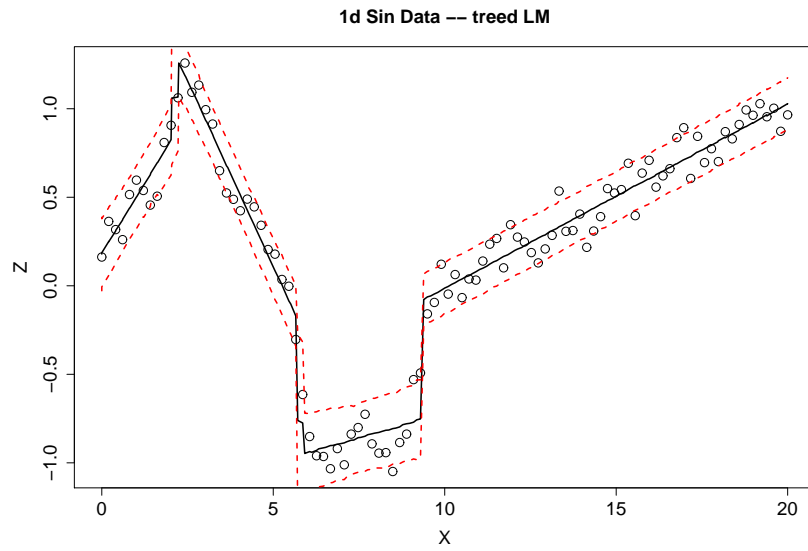
$$\begin{aligned} \mathbf{C}_\nu^{-1} &= (\mathbf{K}_\nu + \mathbf{F}_\nu\mathbf{W}\mathbf{F}_\nu^\top/\tau_\nu^2)^{-1} & \mathbf{q}_\nu(\mathbf{x}) &= \mathbf{k}_\nu(\mathbf{x}) + \tau_\nu^2\mathbf{F}_\nu\mathbf{W}_\nu\mathbf{f}(\mathbf{x}) \\ \kappa(\mathbf{x}, \mathbf{y}) &= K_\nu(\mathbf{x}, \mathbf{y}) + \tau_\nu^2\mathbf{f}^\top(\mathbf{x})\mathbf{W}\mathbf{f}(\mathbf{y}) & \mathbf{f}^\top(\mathbf{x}) &= (1, \mathbf{x}^\top) \end{aligned}$$

Expected reduction in squared error

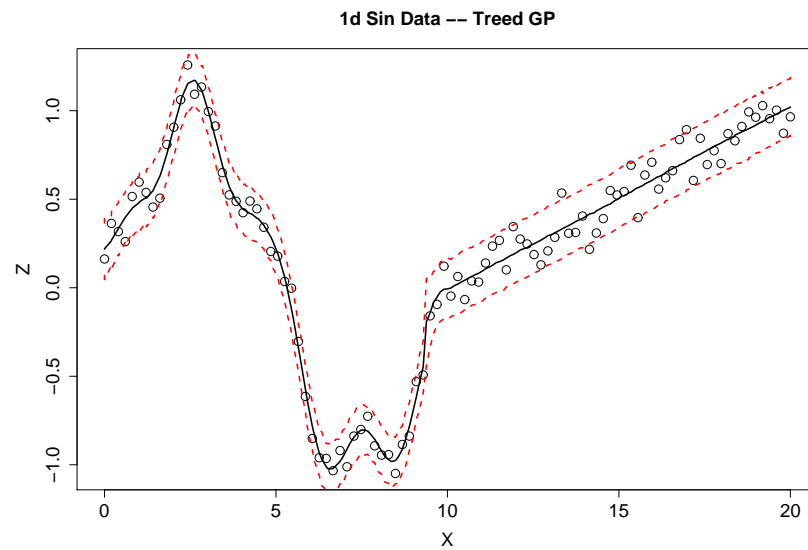
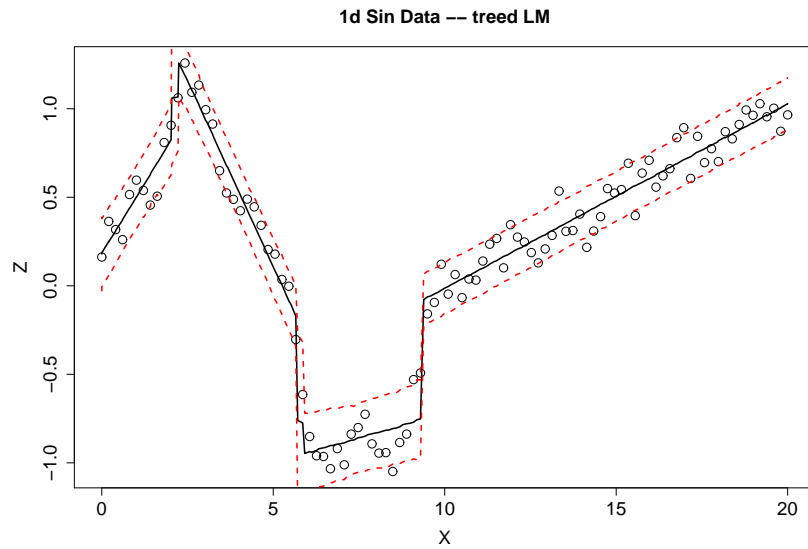
$$\Delta\hat{\sigma}_y^2(\tilde{\mathbf{x}}) = \hat{\sigma}_y^2 - \hat{\sigma}_y^2(\tilde{\mathbf{x}}) = \frac{[\mathbf{q}_N^\top(\mathbf{y})\mathbf{C}_N^{-1}\mathbf{q}_N(\tilde{\mathbf{x}}) - \kappa(\tilde{\mathbf{x}}, \mathbf{y})]^2}{\kappa(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}) - \mathbf{q}_N^\top(\tilde{\mathbf{x}})\mathbf{C}_N^{-1}\mathbf{q}_N(\tilde{\mathbf{x}})}$$



# TREED GP ON SINE DATA

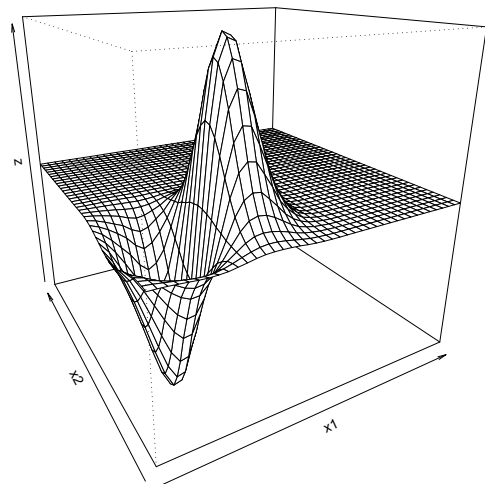


# TREED GP ON SINE DATA

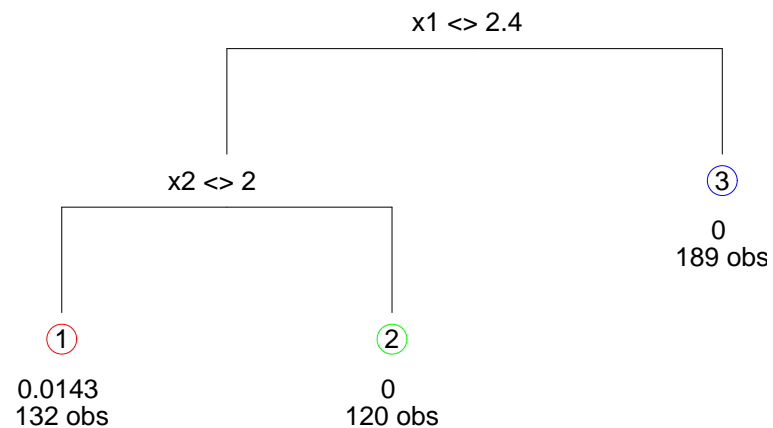
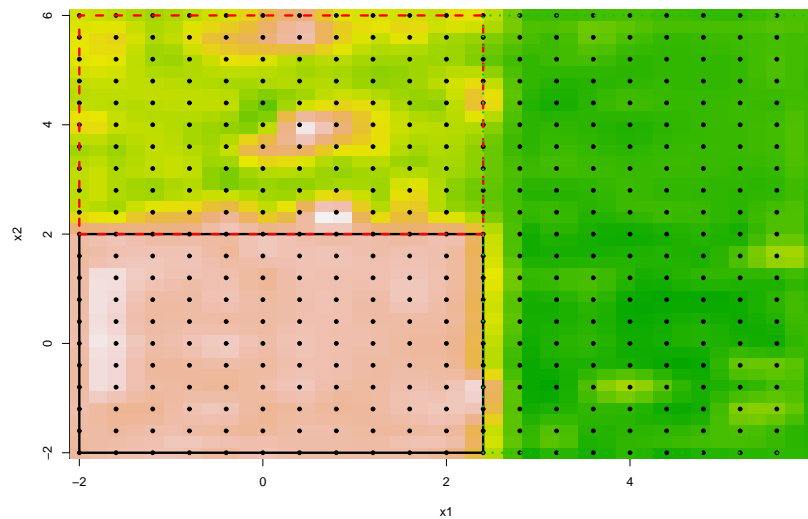


# TREED GP ON EXP EXAMPLE

B-TGPLM: z mean



B-TGPLM: z quantile diff (error)

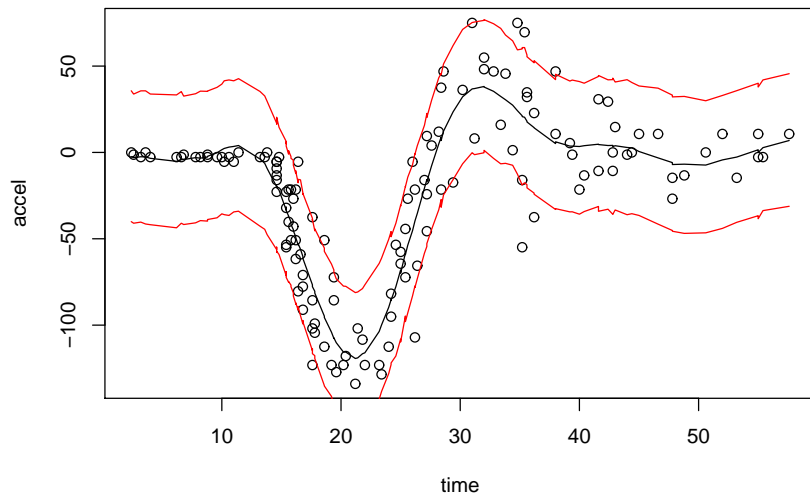


# MOTORCYCLE DATA

## Motorcycle data: (Silverman, 1985)

- non-stationary
- input-dependent noise

GP, accel mean and error

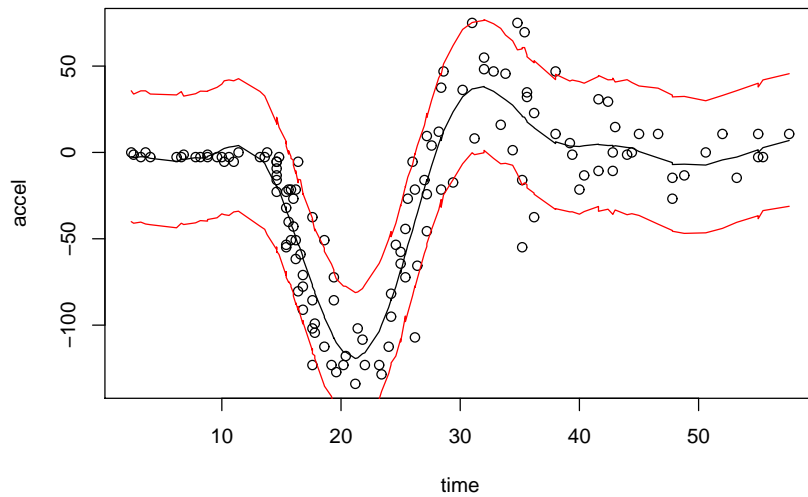


# MOTORCYCLE DATA

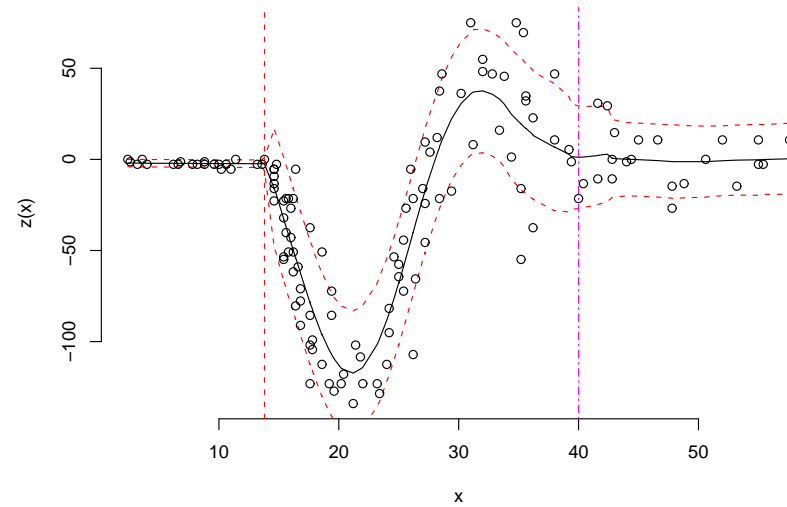
## Motorcycle data: (Silverman, 1985)

- non-stationary
- input-dependent noise

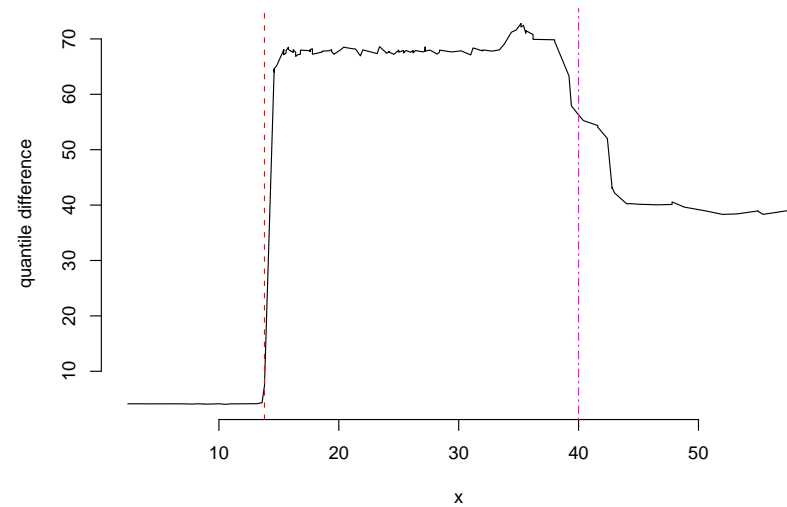
GP, accel mean and error



Estimated Surface



Estimated Error Spread (95th - 5th Quantile)



Bayesian adaptive sampling proceeds in trials ...

Current trial:

- Choose candidates:  $\tilde{\mathbf{X}}$
- estimate model parameters  $(\boldsymbol{\theta}, \mathcal{T})$  for data  $\{\mathbf{X}_i, \mathbf{z}_i\}_{i=1}^N$
- Order candidates  $\tilde{\mathbf{X}}$  by
  - **ALM**: maximize predictive error  $\hat{\sigma}^2(\tilde{\mathbf{x}})$
  - **ALC**: maximize average expected reduction in predictive error

$$\Delta\hat{\sigma}^2(\tilde{\mathbf{x}}) = \int_D \Delta\hat{\sigma}_y^2(\tilde{\mathbf{x}}) d\mathbf{y}$$

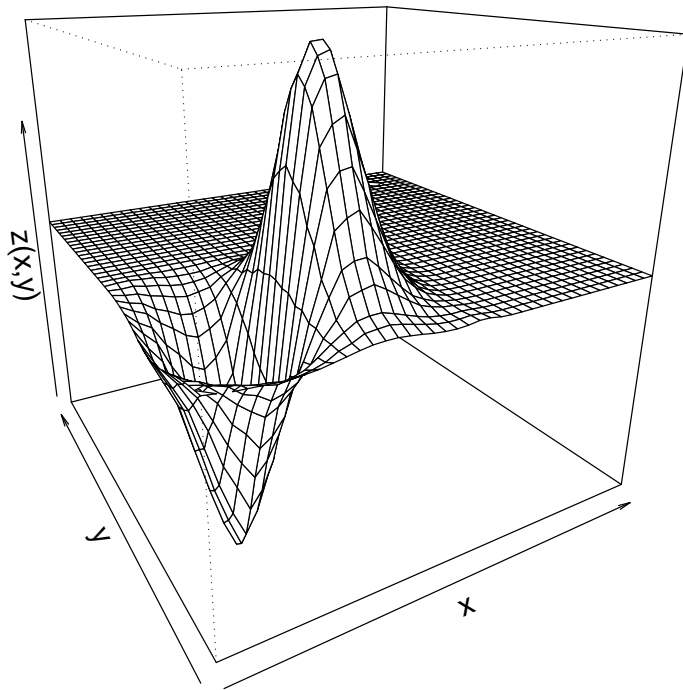
- (repeat) begin next trial

# ADAPTIVE SAMPLING COMPARISON

## 2d exponential data:

$$y(\mathbf{x}) = x_1 \exp(-x_1^2 - x_2^2)$$

2d Exp Data -- Treed GP



75 adaptive samples

model	as	rmse
btgp	alc	0.0035
btgp	alm	0.0037
bcart	alm	0.0090
bcart	alc	0.0093
btlm	alm	0.0099
btlm	alc	0.0115
bgp	alm	0.0352
bgp	alc	0.0493

Adaptive sampling on the Langley-Glide-Back Booster (LGBB):

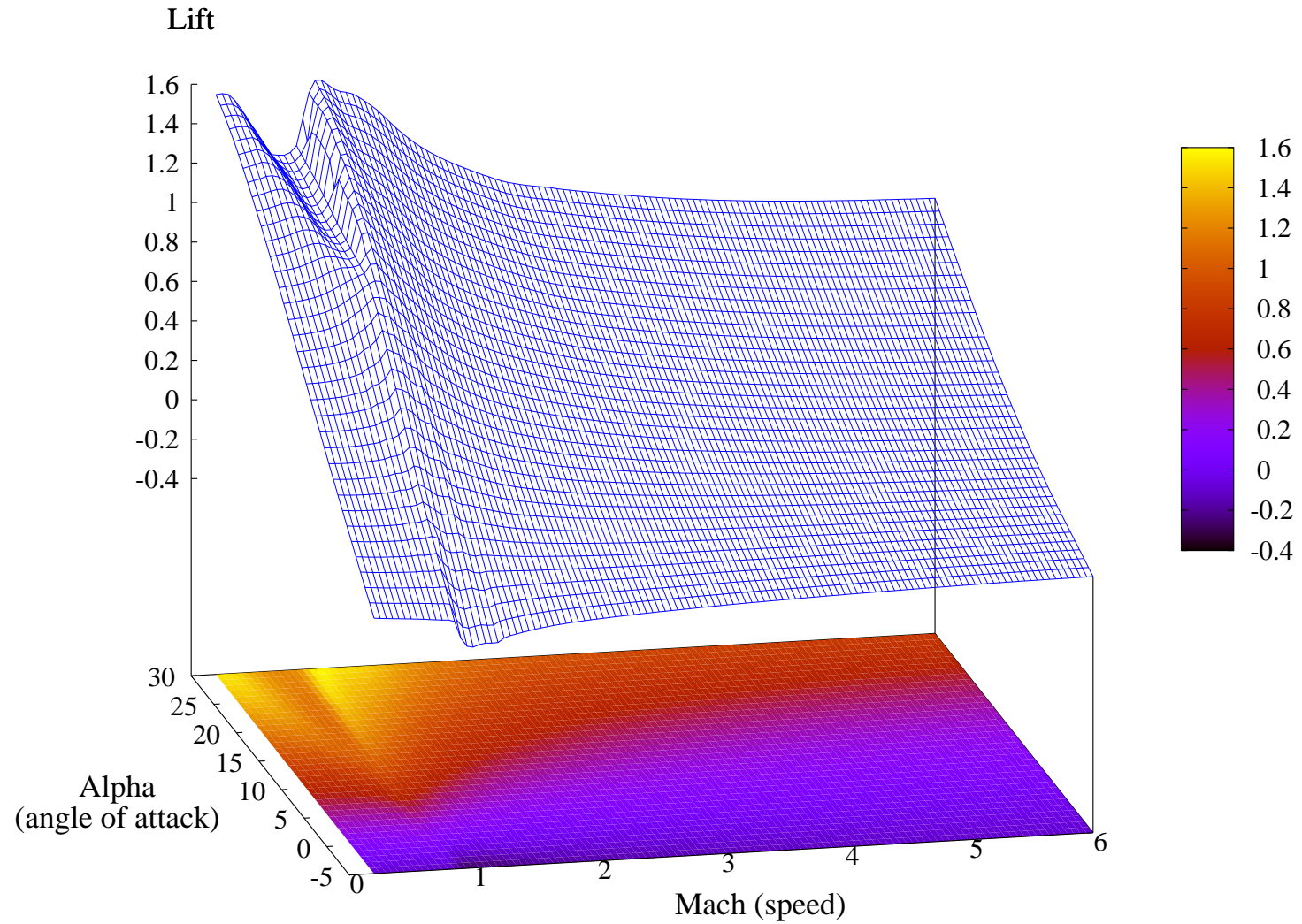
Interfacing with a NASA supercomputer:

- fit six independent treed GPs— one for each response (lift, drag, pitch, side-force, yaw, roll)
  - Six parallel treed GP modules
- design via expected reduction in predictive error (ALC)
  - pooled predictive quantiles across the six models

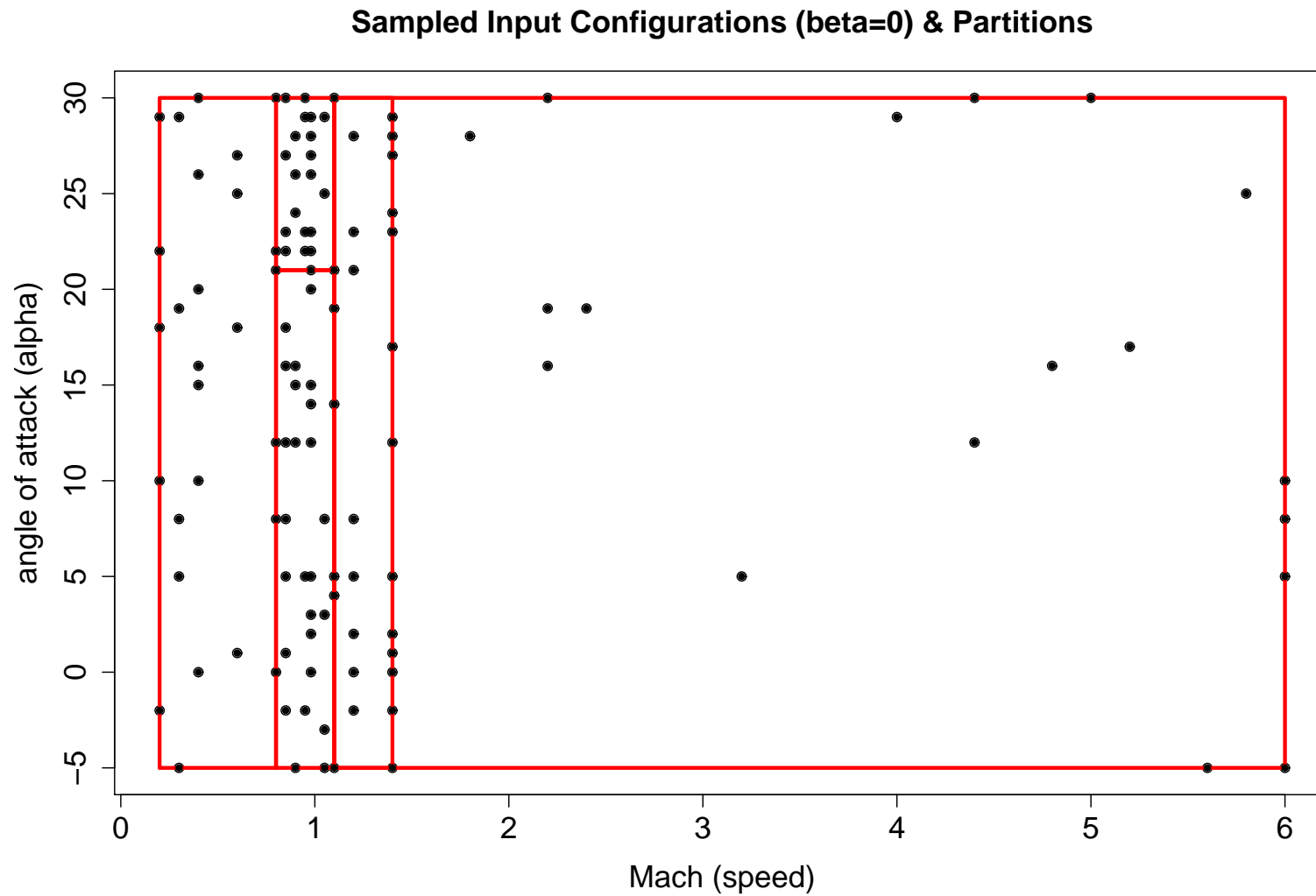


# ADAPTIVE SAMPLING ON LGBB: LIFT

Mean posterior predictive -- Lift  
fixing Beta (side slip angle) to zero

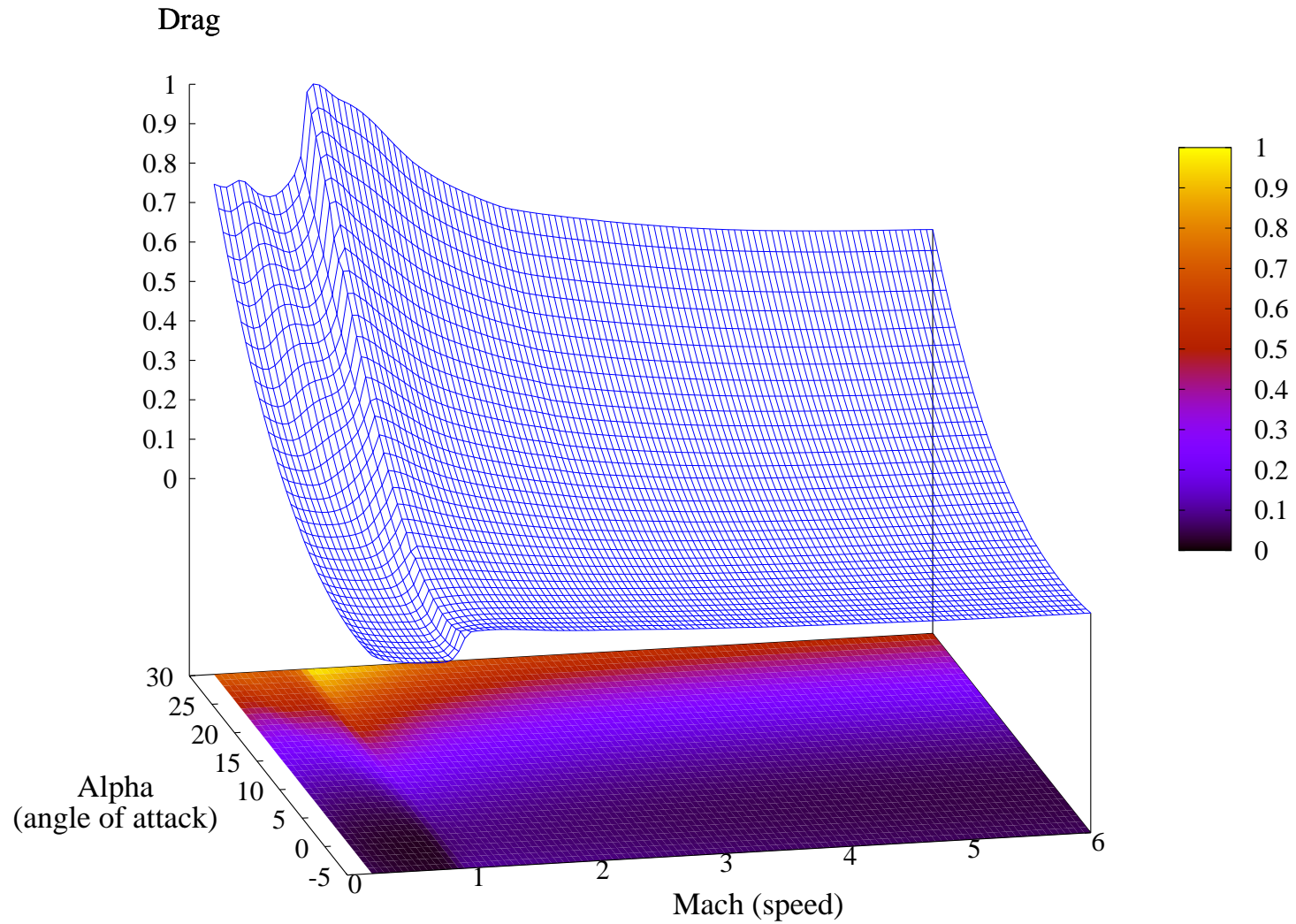


# ADAPTIVE SAMPLING ON LGBB: LIFT



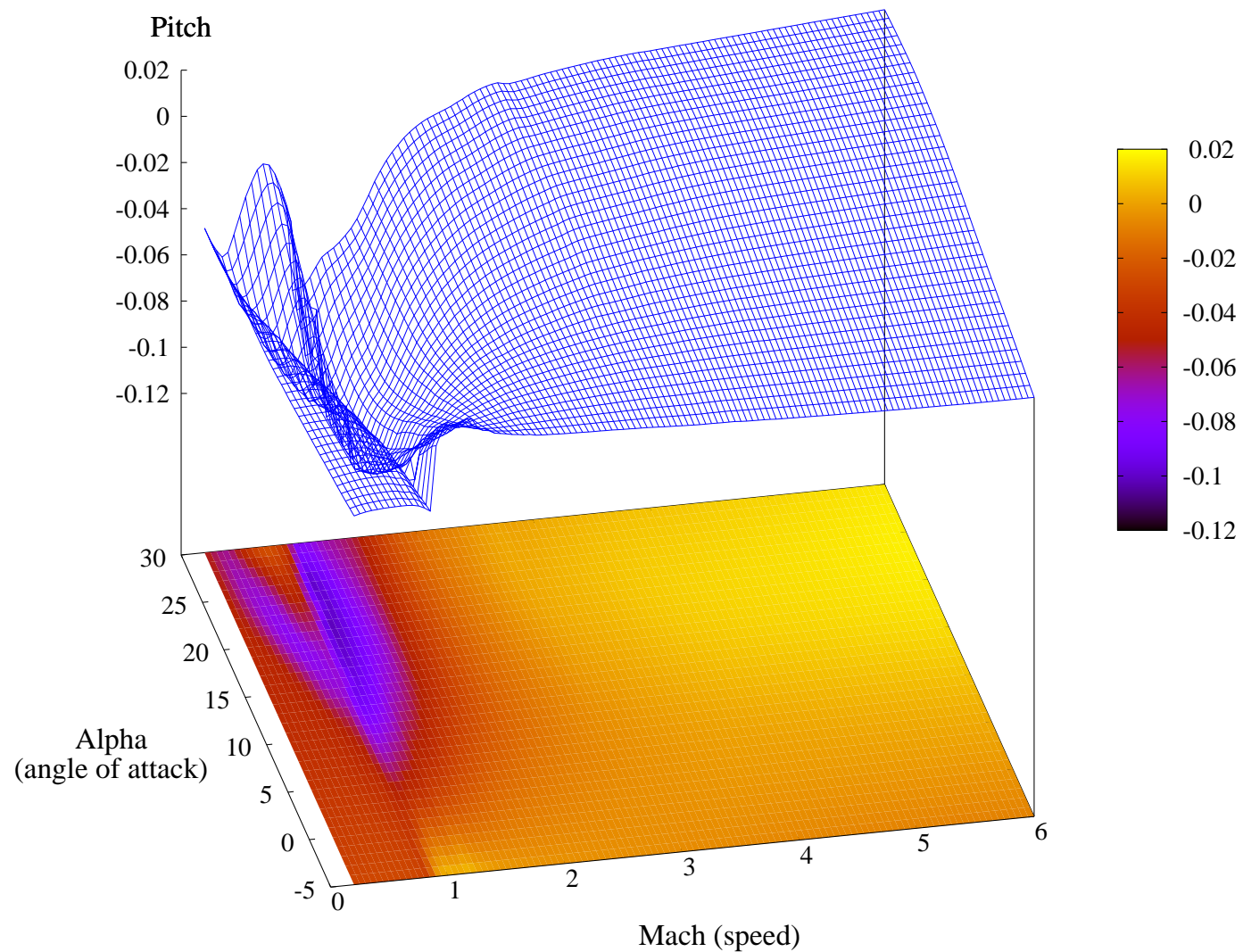
# ADAPTIVE SAMPLING ON LGBB: DRAG

Mean posterior predictive -- Drag  
fixing Beta (side slip angle) to zero



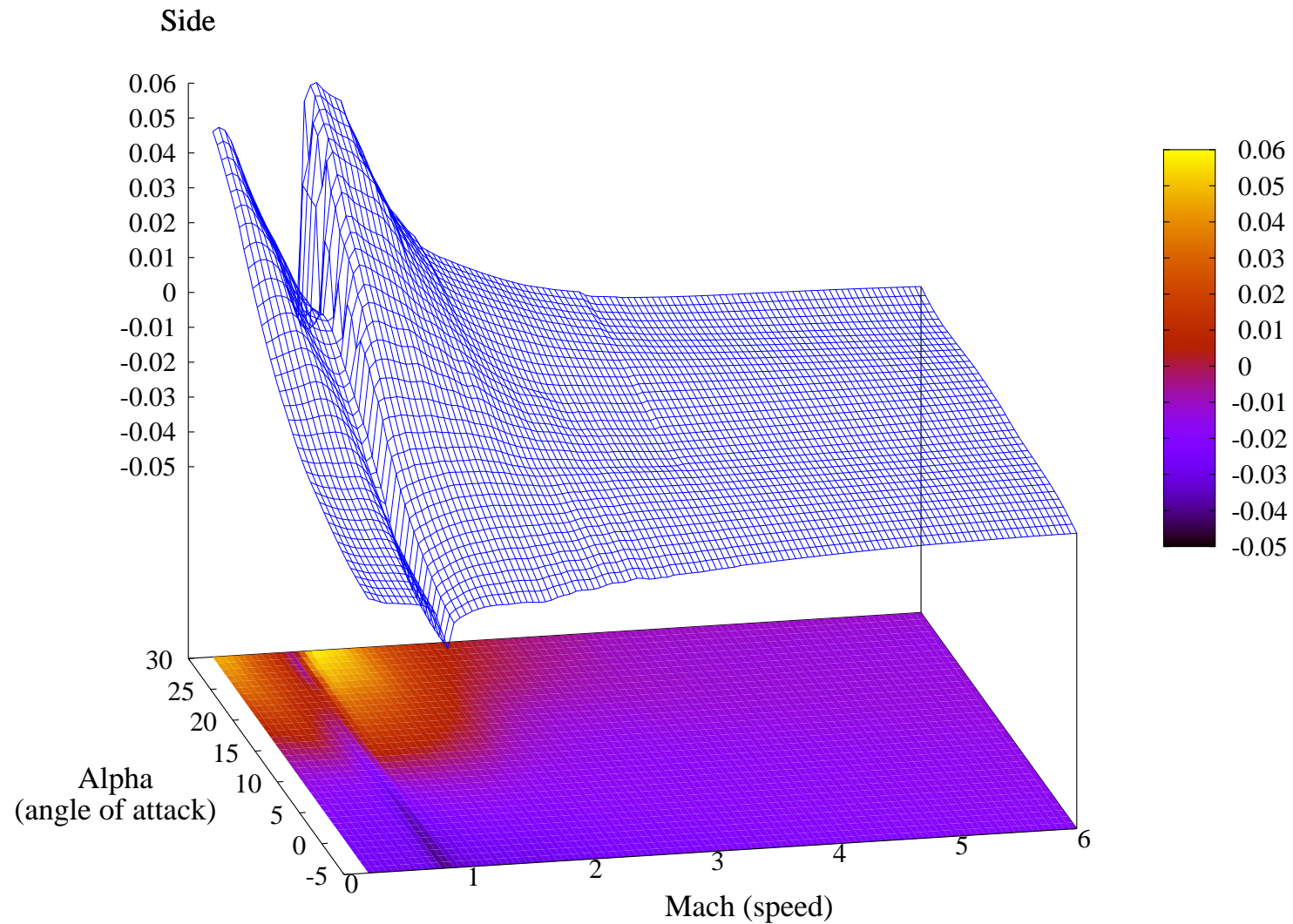
# ADAPTIVE SAMPLING ON LGBB: PITCH

Mean posterior predictive -- Pitch  
fixing Beta (side slip angle) to zero



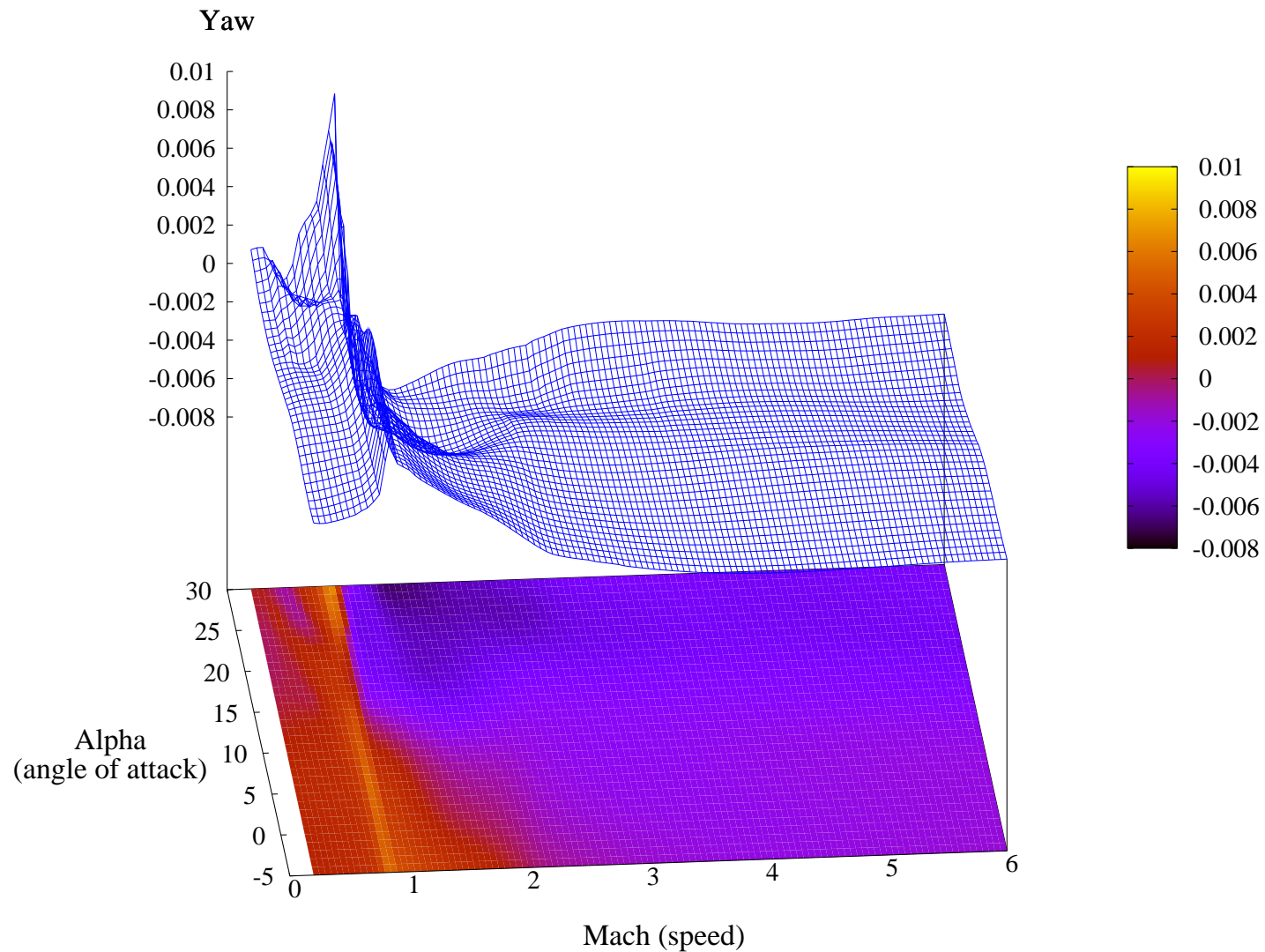
# ADAPTIVE SAMPLING ON LGBB: SIDE

Mean posterior predictive -- Side  
fixing Beta (side slip angle) to 2



# ADAPTIVE SAMPLING ON LGBB: YAW

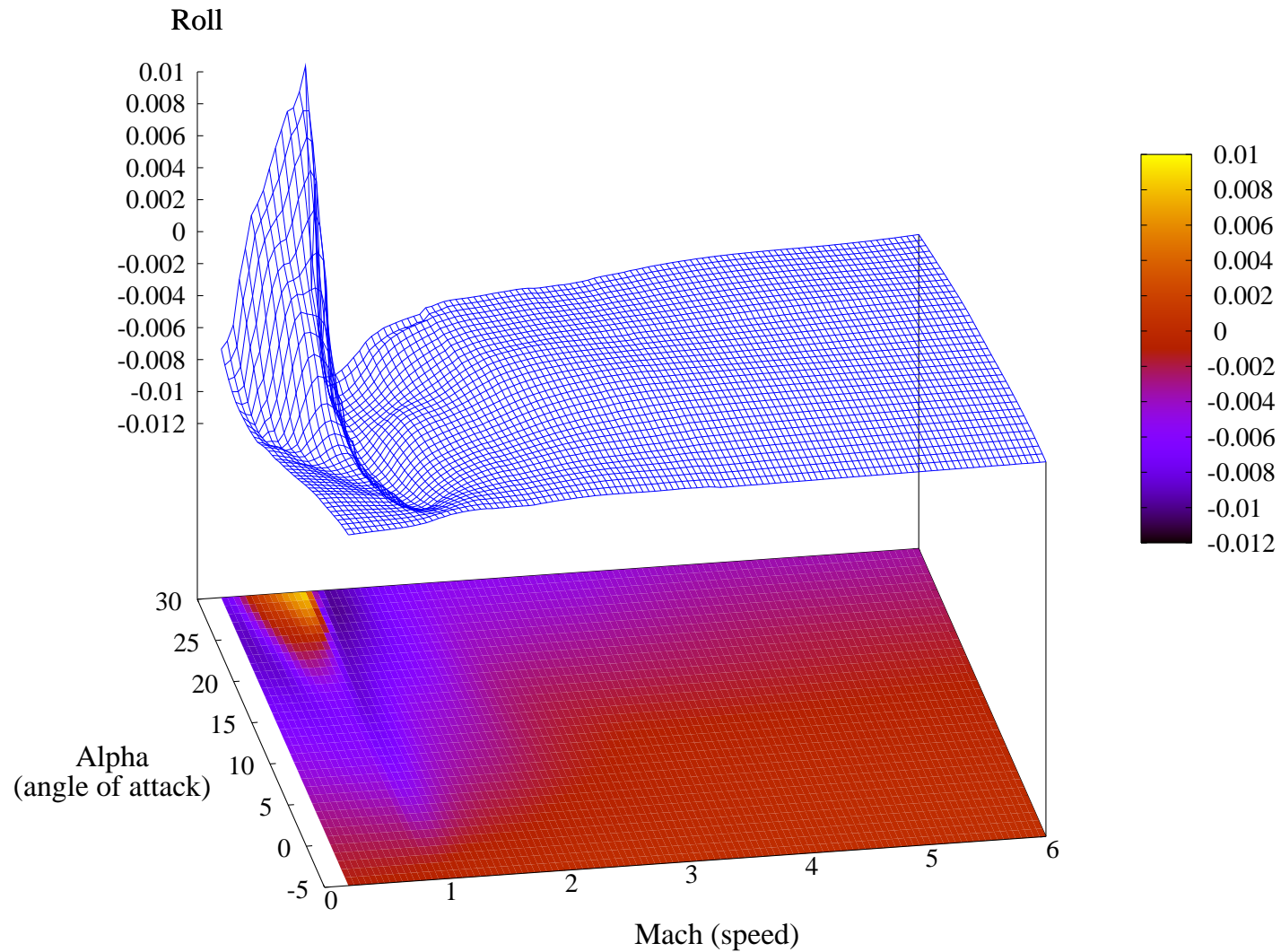
Mean posterior predictive -- Yaw  
fixing Beta (side slip angle) to 2



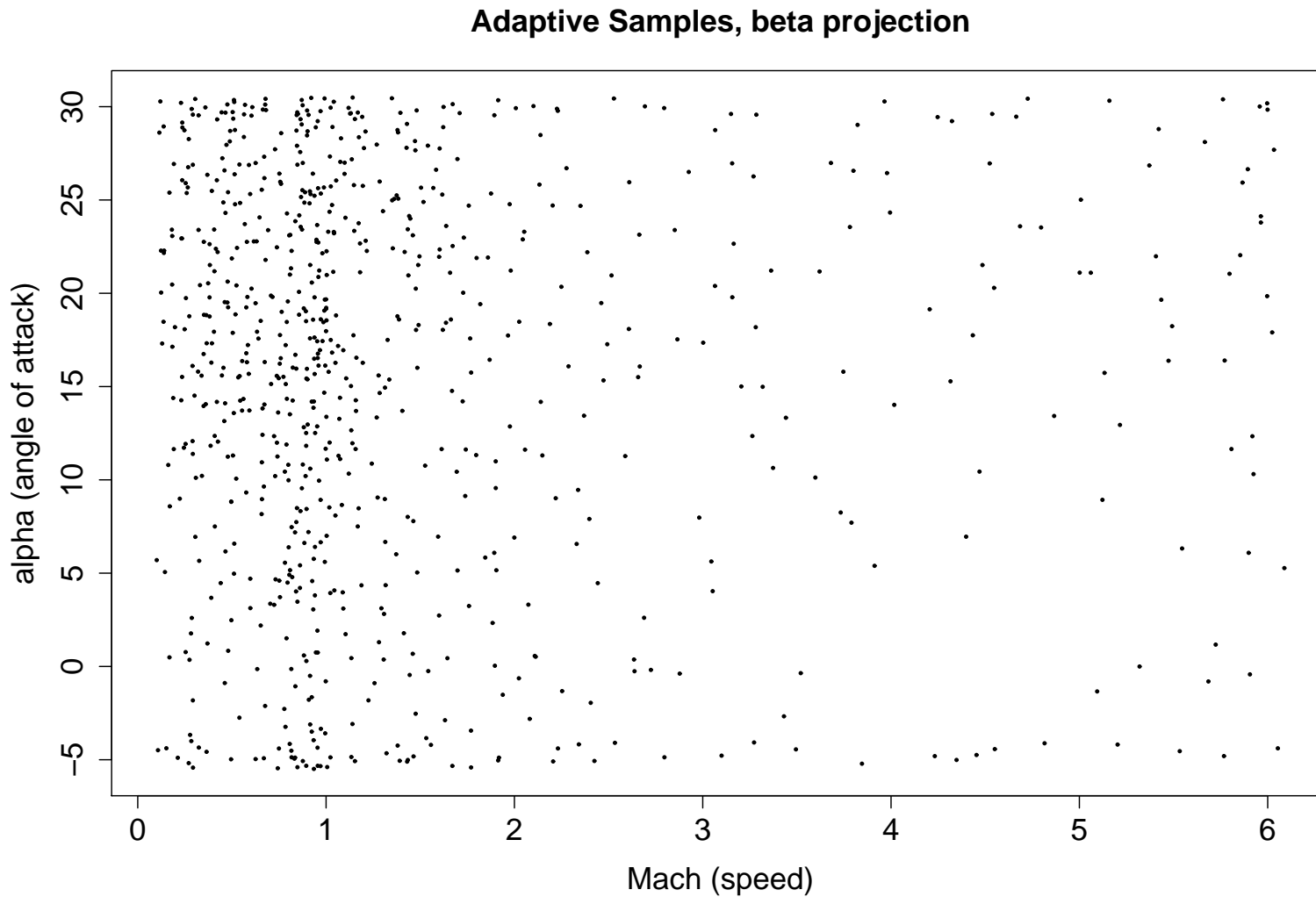


# ADAPTIVE SAMPLING ON LGBB: ROLL

Mean posterior predictive -- Roll  
fixing Beta (side slip angle) to 2



# ADAPTIVE SAMPLING ON LGBB



780 adaptive samples, compared to more than 3,250



# IMPLEMENTATION & R PACKAGE

---

## Implementation & computing details

- C++ (trees) and C (GPs) with LAPACK/BLAS
- Pthreads for (shared memory) parallelization:

R package called `tgp` on CRAN

- Implements all model combinations and special cases
  - `blm`, `bgp`, `btlm`, `btgp`, `bgpllm`, `btgpllm`
- Adaptive sampling: ALM, ALC, etc.