

























Note: model is (for now!) correct





- Model instantiated to  $Y = \sum_{j=0}^{k} \beta_j X^j + \epsilon$
- Let's experiment to see what happens if data are sampled from following "true" distribution: X<sub>i</sub> ~ Unif.[-1, 1], i.i.d.

 $Y_i = X_i + \epsilon_i, \epsilon_i \sim \text{Normal}(0, 1), \text{i.i.d.}$ 

- Note: model is (for now!) correct
- ...and Bayes works perfectly well: posterior concentrates around  $\beta=(0,1,0,\ldots,0)$  after just a few outcomes and keeps on doing so for ever
- Least Squares/ML would perform terribly here!

# Bayesian High Dimensional Regression

Two standard approaches:

- Bayesian Ridge/Lasso/Horseshoe Regression
- Bayesian Model Selection/Model Averaging

### Bayes Factor Model Selection (Hypothesis Testing)

 $\mathcal{M}_k = \{p_\theta \mid \theta \in \Theta_k\}, k = 1, 2, ...$  $\hat{k}(y^n)$  is *k* maximizing a posteriori probability

$$p(y^n \mid \mathcal{M}_k)\pi(k)$$

$$(\mathcal{M}_k \mid y^n) = \frac{p(y \mid \mathcal{M}_k)\pi(w)}{\sum_{k \in \mathcal{K}} p(y_n \mid \mathcal{M}_k)\pi(k)}$$

 $p(y^n \mid \mathcal{M}_k) = \int_{\theta \in \Theta_k} p_{\theta}(y^n) \pi(\theta \mid k) d\theta$ 

#### $\pi(k)$ is prior **on** models

 $\pi(\theta \mid 1), \pi(\theta \mid 2),$  are priors within models

# **Bayes Factor Model Selection**

- Standard Bayesian method to select a model based on the data
- Can be used to select degree of polynomial
  - Bayes has a built-in Occam's Razor: automatic regularization
  - Complex models 'penalized' automatically (even if flat priors are used within these models)
- Close Relation to information-theoretic Minimum Description Length (MDL) Model Selection



# **Bayes Factor Model Selection**

- Standard Bayesian method to select a model based on the data
- Can be used to select degree of polynomial
  Bayes has a built-in Occam's Razor:
  - automatic regularization
- If the model is correct i.e. well-specified, then this guarantees that the 'right' degree will be selected given enough time, with probability 1 (consistency)
  - In contrast standard least squares would always select a polynomial with 0 error on the data and degree equal to number of data points -1



- ...and Bayes factor model selection works perfectly
- well, selects 0-degree model after just a few outcomes and keeps on doing so for ever





















Freund, Schapire, Cesa-Bianchi, Lugosi, Hazan, Kale, ..., Koolen, Van Erven











# Explaining SafeBayes via Ridge

- Frequentist (non-Bayesian) Ridge Regression:
- For each  $\lambda$  select  $\widehat{\beta_{\lambda}}$  minimizing

$$\sum_{i=1}^{n} (y_i - \sum_{j=0}^{k} \beta_j g_j(x_i))^2 + \lambda \sum_{j=0}^{k} \beta_j$$

- Select final  $\hat{\beta} = \hat{\beta}_{\lambda_{cv}}$  by cross-validating over  $\lambda$
- $\hat{\beta}_{\lambda}$  is also the posterior mean/MAP of Bayesian ridge regression with fixed variance  $\sigma^2 = 2 \lambda$

# Explaining SafeBayes via Ridge

• For each  $\lambda$  select  $\widehat{\beta_{\lambda}}$  minimizing

$$\sum_{i=1}^{n} (y_i - \sum_{j=0}^{k} \beta_j g_j(x_i))^2 + \lambda \sum_{j=0}^{k} \beta_j^2$$

- Select final  $\hat{\beta} = \hat{\beta}_{\lambda_{cv}}$  by cross-validating over  $\lambda$
- $\hat{\beta}_{\lambda}$  is also the posterior mean/MAP of Bayesian ridge regression with fixed variance  $\sigma^2 = 2\lambda$
- under 'bad' misspecification, putting prior on  $\lambda$  leads to posterior concentrating on way smaller  $\lambda$  than  $\lambda_{cv}$





### Main Result as a Slogan

- "Generalized Bayes is great...
   ...once you know the right learning rate"
- "Safe Bayes is great... ...even if you don't know the learning rate"

#### Main Result as a Slogan

- "Generalized Bayes is great...
   ...once you know the right learning rate"
- "Safe Bayes is great... ...even if you don't know the learning rate"

Most Extensive Explanation So-Far: Inconsistency of Bayesian Inference for Misspecified Linear Models, and a Proposal for Repairing It http://arxiv.org/abs/1412.3730

# Final Remarks - I

- Bayes can be in trouble when model is wrong but useful; adding learning rate helps
- There are other issues with Bayes as well, e.g. in nonparametrics. These are *not* the classical objections 'it is subjective' or 'where do you get your prior'
  - see 'Larry and Jamie take on a Nobel Prize Winner' on Larry Wasserman's blog
- ...but also amazing successes

# **Final Remarks**

- For sequential prediction, the learning-rate approach
   is common among theorists and extremely robust
- Just a few practical applications, but these are very successful
- ...e.g. online prediction of electricity demand in Paris region
  - M Devaine, P Gaillard, Y Goude, G Stoltz, Machine Learning 90 (2), 231-260
- They actually used SafeBayes (unconsciously)







