## **Predictive model selection** and uncertainty Vik Shirvaikar (joint work with Stephen Walker and Chris Holmes)

**June 2025** 

## How do we view uncertainty?

### Frequentists

The parameter is fixed and the data is random

"If I ran this experiment several times, 95% of the confidence intervals I build using this procedure would contain the true mean."

### Bayesians

The data is fixed and the parameter is random

"Given the observed data and my prior, there is a 95% posterior probability that the true mean lies within this credible interval."

## Can we get the best of both worlds?

**Uncertainty directly on a parameter**, without having to interpret results through the lens of repeated sampling...

...but in a completely data-dependent framework that doesn't require the specification of a prior distribution

### Outline

1. Prequential forecasting and proper scoring

- 2. Martingale posterior distributions (Fong et al., 2023)
- 3. Model uncertainty via predictive resampling
  - A. Illustration: density estimation
  - B. Illustration: hypothesis testing

### Prior works have emphasized the value of prediction

All this is not to deny the usefulness and convenience of parametric models; clearly, the extraordinarily appealing notion of exchangeability when applied to statistical paradigms implies the mathematical existence of the lurking parameter as evidenced by the de Finetti representation theorem. Further, in some situations parameters may be important hypothetical constructs in providing an appropriate lubricant or focus for a model and in planning an experimental program whether or not they can be regarded as real physical entities. However, the time has come for statistical analyses to emphasize the observable and the finite in contrast to the potentially artificial and infinite.

Geisser (1982), "Aspects of the predictive and estimative approaches in the determination of probabilities"



### **Prior works have emphasized the value of prediction**

### SUMMARY

The prequential approach is founded on the premiss that the purpose of statistical inference is to make sequential probability forecasts for future observations, rather than to express information about parameters. Many traditional parametric concepts, such as consistency and efficiency, prove to have natural counterparts in this formulation, which sheds new light on these and suggests fruitful extensions.

probability forecast distribution  $P_{n+1}$  for the next observation  $X_{n+1}$ .

**3. PREQUENTIAL FORECASTING SYSTEMS** So now let  $X = (X_1, X_2, ...)$  be a sequence of uncertain quantities. At any time n, the prequential forecaster, with the values  $\mathbf{x}^{(n)}$  of  $\mathbf{X}^{(n)} = (X_1, X_2, \ldots, X_n)$  to hand, must issue a

Dawid (1984), "Statistical theory: the prequential approach"



## A scoring rule is a summary measure

- Suppose I have a model  $\mathcal{M}$  and start to observe some data  $x_1, x_2, \ldots$
- How do I evaluate this model's forecasts? I use a scoring rule
- For example, I could use a very simple rule where I score one point if my mostlikely guess is correct, and zero points otherwise

$$S(x, \mathcal{M}) = \begin{cases} 1, \\ 0, \end{cases}$$

Matheson and Winkler (1976), "Scoring Rules for Continuous Probability Distributions"

if 
$$x = \arg \max_{y} \mathcal{M}(y)$$
,  
y  
otherwise





## The marginal likelihood is a scoring rule

product of one-step-ahead predictive densities

$$p(x_{1:n} \mid \mathcal{M}) =$$

• By logging each side, we get a scoring rule!

$$S(x_{1:n}, \mathcal{M}) = \log p(x_{1:n} \mid \mathcal{M}) = \sum_{i=1}^{n} \log p_{\mathcal{M}}(x_i \mid x_{1:i-1})$$

• The Bayesian marginal likelihood of any model *M* can be factorized as a

$$= \prod_{i=1}^{n} p_{\mathcal{M}}(x_i \mid x_{1:i-1})$$

## The marginal likelihood is a scoring rule

- It turns out that this particular scoring rule has some nice properties:
  - It's prequential, treating the data as if they arrive in sequence, like in reality
  - It's proper, because its expectation is maximized by using the true distribution for forecasting — there's no way to "game the system"
  - It's unique in guaranteeing coherent model evaluation

$$S(x_{1:n}, \mathcal{M}) = \log p(x_{1:n} \mid \mathcal{M}) = \sum_{i=1}^{n} \log p_{\mathcal{M}}(x_i \mid x_{1:i-1})$$

Gneiting and Raftery (2007), "Strictly Proper Scoring Rules, Prediction, and Estimation" Fong and Holmes (2020), "On the Marginal Likelihood and Cross-Validation"

## The Bayes factor is a marginal likelihood (ratio)

- - $\log BF = S(x_1)$

We know that the **Bayes factor** is the ratio of evidence between two models

 $BF = \frac{p(x_{1:n} \mid \mathcal{M}_1)}{p(x_{1:n} \mid \mathcal{M}_2)}$ 

This means the Bayes factor has an interpretation based on prequential scoring

$$:_n, \mathcal{M}_1) - S(x_{1:n}, \mathcal{M}_2)$$

I.J. Good (1952), "Rational Decisions"

Dawid (1984), "Statistical theory: the prequential approach"





## The BIC is (approximately) the Bayes factor

- Our favorite model selection rule is the **BIC** (because it has "Bayesian" in it)  $BIC(\mathcal{M}_i) = -2\log i$
- Optimizing the BIC is asymptotically equivalent to optimizing the Bayes factor  $\log BF \approx -\frac{\operatorname{BIC}(\mathcal{M}_1) - \operatorname{BIC}(\mathcal{M}_2)}{2}$
- This allows the BIC to inherit key properties of the Bayes factor, specifically **consistency** — the probability of selecting the correct model converges to one

$$p(x_{1:n} \mid \mathcal{M}_j, \hat{\theta}_j) + d_j \log(n)$$

Schwarz (1978), "Estimating the dimension of a model" Kass and Raftery (1995), "Bayes factors"



### The BIC is an optimal predictive summary measure!

How do we evaluate model predictions? 

### → We use a scoring rule

- such as the marginal likelihood
  - which is captured by the Bayes factor
    - which is approximated by the BIC

the model-selection method which proceeds by maximizing the adjusted prequential likelihood, or equivalently minimizing the "adjusted stochastic complexity" of  $x^n$ ,  $-\log p_j(x^n) - \log \alpha_j$ , will be (almost everywhere) consistent.

Dawid (1992), "Prequential Analysis, Stochastic Complexity, and Bayesian Inference"



## Outline

- 1. Prequential forecasting and proper scoring
- 2. Martingale posterior distributions (Fong et al., 2023)
- 3. Model uncertainty via predictive resampling
  - A. Illustration: density estimation
  - B. Illustration: hypothesis testing

## Uncertainty comes from missing data

Clarity is even more easily attainable, perhaps, if it is recognized that an inference about finite population parameters can be interpreted as a predictive inference about a sample, namely, the unsampled part of a finite population. This sample is already in existence but has not been observed by the statis-

Roberts (1965), "Probabilistic prediction"



## Uncertainty comes from missing data

### JOURNAL ARTICLE

### Martingale posterior distributions 3

Edwin Fong, Chris Holmes 🖾, Stephen G Walker 🛛 Author Notes

Journal of the Royal Statistical Society Series B: Statistical Methodology, Volume 85, Issue 5, November 2023, Pages

### Abstract

The prior distribution is the usual starting point for Bayesian uncertainty. In this paper, we present a different perspective that focuses on missing observations as the source of statistical uncertainty, with the parameter of interest being known precisely given the entire population. We argue that the foundation of Bayesian inference is to assign a distribution on missing observations conditional on what has been observed. In the i.i.d. setting with an observed sample of size *n*, the Bayesian would thus assign a predictive distribution on the missing  $Y_{n+1:\infty}$  conditional on  $Y_{1:n}$ , which then induces a distribution on the parameter.



## The predictive resampling pipeline

# Simulate $Y_{i+1} \sim p(\cdot \mid Y_{1:i})$ Update $p(\cdot | Y_{1:i}) \rightarrow p(\cdot | Y_{1:i+1})$



### Calculate final parameter $\theta_{\infty} = \theta(Y_{1:\infty})$

Repeat over several Monte Carlo iterations (stopping in practice at some large N >> n)



## **Different perspective, same uncertainty**

likelihood prior Let  $f_{\theta}(y) = \mathcal{N}(y \mid \theta, 1)$ , with  $\pi(\theta) = \mathcal{N}(\theta \mid 0, 1)$ . Given an observed dataset  $y_{1:n}$ , the tractable posterior density takes on the form  $\pi(\theta \mid y_{1:n}) = \mathcal{N}(\theta \mid \overline{\theta}_n, \overline{\sigma}_n^2)$  where

$$\bar{\theta}_n = \frac{\sum_{i=1}^n y_i}{n+1}, \quad \bar{\sigma}_n^2 = \frac{1}{n+1}.$$

The posterior predictive density then takes on the form  $p(y | y_{1:n}) = \mathcal{N}(y | \overline{\theta}_n, 1 + \overline{\sigma}_n^2)$ . For observed data, we generated  $y_{1:n} \stackrel{\text{iid}}{\sim} f_{\theta}(y)$  for n = 10 with  $\theta = 2$ , giving  $\bar{\theta}_n = 1.84$ .

Fong et al. (2023), "Martingale posterior distributions"



## **Different perspective, same uncertainty**



The **predictive** uncertainty from

... is equivalent to the usual posterior the imputation of additional data... uncertainty from Bayesian methods

Fong et al. (2023), "Martingale posterior distributions"



## What type of uncertainty is this?



The total variance can be decomposed into aleatoric and epistemic

Fong et al. (2023), "Martingale posterior distributions"



### Outline

- 1. Prequential forecasting and proper scoring
- 2. Martingale posterior distributions (Fong et al., 2023)
- 3. Model uncertainty via predictive resampling
  - A. Illustration: density estimation
  - B. Illustration: hypothesis testing

## Acknowledging model structure

of S. It is common in statistical theory and practice to acknowledge parametric uncertainty about  $\theta$  given a particular assumed structure S; it is less common to acknowledge structural uncertainty about S itself. A widely used approach involves enlisting the aid of x to specify a plausible single 'best' choice  $S^*$  for S, and then proceeding as if  $S^*$  were known to be correct. In general this approach fails to assess and propagate structural uncertainty fully and may lead to miscalibrated uncertainty assessments about y given x. When miscalibration



Let our model  $\mathcal{M} = \{S, \theta\}$  consist of structural assumptions S, such as the link function in a GLM, and parameters  $\theta$  whose meaning depends on S

Draper (1995), "Assessment and propagation of model uncertainty"



## **Bayesian model averaging**



$$P(\xi \mid \mathscr{D}) = \sum_{k=1}^{K} P(\xi)$$

### model selection with AIC, BIC, etc. $P(\xi \mid \mathcal{D}) = \frac{P(\mathcal{D} \mid \xi) P(\xi)}{P(\mathcal{D})} \quad \longleftrightarrow \quad P(\xi \mid \mathcal{D}, \mathcal{M}^*) = \frac{P(\mathcal{D} \mid \xi, \mathcal{M}^*) P(\xi \mid \mathcal{M}^*) P(\mathcal{M}^*)}{P(\mathcal{D} \mid \mathcal{M}^*)}$

 $P(\xi \mid \mathcal{M}_k, \mathcal{D}) P(\mathcal{M}_k \mid \mathcal{D})$ 

model weighting



## Where do we get the weights?

Existing approaches share some common concerns

- Specification of priors
- Complex integration over parameter spaces
- Inexact or approximate results

So we return to our guiding principle from earlier...





## The predictive resampling pipeline



### Calculate final model



**Repeat over several Monte Carlo iterations** (stopping in practice at some large N >> n)



## The predictive resampling pipeline

Key requirement: specification of a **consistent** one-step model selection criterion, which is then converted directly into a measure of uncertainty

Convergence again comes from martingales

Calculate final model



 $\mathscr{M}_{\infty} = \mathscr{M}_{\hat{k}(\infty)}$ 

**Repeat over several Monte Carlo iterations** (stopping in practice at some large N >> n)



### Outline

- 1. Prequential forecasting and proper scoring
- 2. Martingale posterior distributions (Fong et al., 2023)
- 3. Model uncertainty via predictive resampling
  - A. Illustration: density estimation
  - B. Illustration: hypothesis testing



# **Problem:** Find the number of components in a Gaussian mixture model (GMM)



## **Problem:** Find the in a Gaussian mixture



•		
	2	
•		
•		
•		
•		
•		

## **Problem:** Find the in a Gaussian mixture



## Outline

- 1. Prequential forecasting and proper scoring
- 2. Martingale posterior distributions (Fong et al., 2023)
- 3. Model uncertainty via predictive resampling
  - A. Illustration: density estimation
  - B. Illustration: hypothesis testing

## How do we test a hypothesis?

### Frequentists

The parameter is fixed and the data is random

"If I ran this experiment several times, I'd expect to reject the null hypothesis only 5% of the time when it is actually true." Bayesians

The data is fixed and the parameter is random

"Given the observed data and my prior, there is a 95% posterior probability that the alternative hypothesis is true over the null."

## Simple point hypothesis demonstration

 $\max_{\theta_k} \mathscr{L} = \sum_{i=1}^n \log p(x_i \mid x_i)$ 

- Unknown mean parameter for  $\mathcal{N}(\theta, 1)$ 
  - Observed data  $x_1, \ldots, x_n$ 
    - $H_0: \theta = 0$
    - $H_1: \theta = 0.1$

$$\theta_k) \longrightarrow x_{n+1} \sim \mathcal{N}(\theta_{\hat{k}(n)}, 1)$$

## Simple point hypothesis demonstration





- Unknown binomial proportion parameter
  - Observed data  $x_1, \ldots, x_n$ 

    - $H_0: \theta = 0.5$  $H_1: \theta \neq 0.5$





### Alternative hypothesis $(H_1)$





### Frequentist

If the null hypothesis is true, how extreme would my result be if I repeated this experiment several times?



### Alternative hypothesis $(H_1)$





### Null hypothesis ( $H_{\Omega}$ )



### Frequentist

If the null hypothesis is true, how extreme would my result be if I repeated this experiment several times?



### Alternative hypothesis $(H_1)$



### The chance of **168 or more (5)** is 2.15%

If it's less than 2.5%/5%, reject the null hypothesis





### Alternative hypothesis $(H_1)$



### Bayesian

What is the evidence for each hypothesis given the observed data?







The ratio of these values is 61.6%, so we find that there is greater evidence for the null hypothesis



### Alternative hypothesis $(H_1)$



### What is the evidence for each hypothesis given the observed data?



The "paradox" comes from the **prior** — we've implicitly assumed that a single value ( $\theta = 0.5$ ) is just as likely as any value between 0 and 1



### Alternative hypothesis $(H_1)$



### What is the evidence for each hypothesis given the observed data?

 $\min \text{BIC} = -2\widehat{\mathscr{S}}$  $H_k$ 

- Unknown binomial proportion parameter
  - Observed data  $x_1, \ldots, x_n$ 
    - $H_0: \theta = 0.5$
    - $H_1: \theta \neq 0.5$

$$\hat{e}(x_{1:n} \mid H_k) + d_k \log(n)$$

 $\min BIC = -2\widehat{\mathscr{L}}$  $H_k$ 

 $\mathscr{L}$  at  $\theta = 0$  and  $d_k =$ 

- Unknown binomial proportion parameter
  - Observed data  $x_1, \ldots, x_n$ 
    - $H_0: \theta = 0.5$  $H_1: \theta \neq 0.5$

$$\hat{e}(x_{1:n} \mid H_k) + d_k \log(n)$$

$$= 0 \longrightarrow x_{n+1} \sim \text{Bernoulli}(0.5)$$

 $\min BIC = -2\widehat{\mathscr{L}}$  $H_k$ 

 $\mathscr{L}$  at  $\theta = \bar{x}_n$  and  $d_k =$ 

- Unknown binomial proportion parameter
  - Observed data  $x_1, \ldots, x_n$ 

    - $H_0: \theta = 0.5$  $H_1: \theta \neq 0.5$

$$\hat{e}(x_{1:n} \mid H_k) + d_k \log(n)$$

= 1 
$$\longrightarrow x_{n+1} \sim \text{Bernoulli}(\bar{x}_n)$$



Results have a prior-free interpretation in terms of posterior uncertainty

Evidence can be accumulated in favor of the null

### The final probabilities are 76.9% for the null and 23.1% for the alternative

## Thank you!

"A general framework for probabilistic model uncertainty" (Shirvaikar, Walker, and Holmes)

https://arxiv.org/abs/2410.17108

