

Bayesian Methods

Alexander. Rothkopf @ uis.wo
@ rothkopf AI

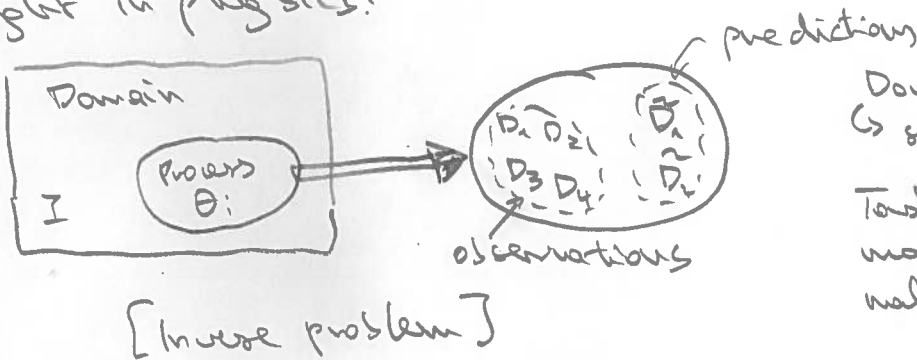
Motivation: - Combining theory and experimental results to explain and predict measurable processes (swathing, transport..)

- Extracting properties of the quark-gluon-plasma from heavy-ion collision data conditioned on a phys. model
J. Bernhard Nature Physics 15 113 (2019)
- Extract from lattice QCD real-time dynamics of quarks & gluons \rightarrow spectral functions, transport coeff., PDF's
A.R. Faust. Phys. 10 28995 (2022)
- Determine Wilson coeff. and their uncertainties from a combined theory & data analysis.
S. Wosonowski J. Phys. A. 43 074001 (2010)

Books: Bayesian Rethinking, R. McElrath

Information Theory, Inference & Learning Algorithms, D. Mackay

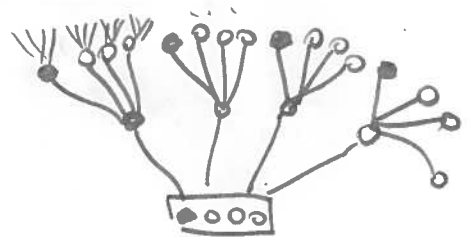
Insight in physics:



Domain knowledge
 \hookrightarrow set of potential models
Task: Infer most plausible model and parameters, then make predictions.

Probability Theory: Efficient formalism for the forward problem: known process, estimate outcomes. [frequentist: probabilities as outcome of repeated experiments]

How probable to draw ? $\frac{1}{4} \cdot \frac{3}{4} \cdot \frac{1}{4} = \frac{3}{64}$



No uncertainty on state of system

Bayesian statistics "Plausibility theory"

Efficient formalism for the inverse problem:
unknown process (both type & parameters) inferred
from observations.



[plausibility: generalized concept of (un)certainty
assigned NOT only to data but parameters & models]
{ still between {0,1} etc }

all possible
compositions
of 4 balls

What is the most plausible composition of an urn, given
observed draws? (prior knowledge: 4 balls):
uncertainty on state of system, need to survey all potential
cases.

Intuitively: The more observations the more certain

Bayesian Inference (or Bayesian learning) via skilful application
of conditional probabilities. Forces us to make explicit ALL
uncertainty by assigning plausibility distributions to data & models.

Starting point: joint probability

$$P(D, \theta, I) = P(D|\theta, I) P(\theta, I) \left. \vphantom{P(D, \theta, I)} \right\} P(\theta|D, I) = \frac{P(D|\theta, I) P(\theta, I)}{P(D, I)}$$

↑ ↑ ↑ ↑ ↑
hour 4 balls
margin

likelihood prior
normalization

likelihood: How is data generated (conditioned on model)

Prior: Plausibility of model & parameters from domain knowledge

Normalization: Plausibility relative to all potential models

Example of a simple hierarchical / multilevel model
(see also Bayesian graph networks)




$$P(\dots | \dots, 4) = \frac{P(\dots | \dots, 4) P(\dots, 4)}{\text{all possible compositions}}$$

What about $P(\dots, 4)$? A priori no reason to favor one composition
over other 1/s. (BEWARE: choosing a good representation of
ignorance may be non-trivial, see e.g. Jeffreys' prior)

$L[000]$	prior	post.
0000	0	0
•000	3/64	0.15
••00	3/64	0.4
•••0	9/64	0.45
••••	0	0

- Distribution over parameters
- simple to count L for 3 draws

What if new draws comes in 

Efficient Bayesian updating: Posterior \rightarrow new prior
likelihood \rightarrow prob. of STORAGE
new draws cond. on model

$L[0]$	prior	post.
0	0	0
1/4	0.15	0.07
3/4	0.4	0.35
3/4	0.45	0.54
1	0	0

New observation favors 10000

\Rightarrow iterated learning process

Interpretation of prior plausibility: model (parameter) uncertainty informed by domain prior knowledge.

(in the past prior often chosen for computational convenience "conjugate prior" or on generic information theory arguments)

BEWARE: Bayes is explicit in modeling choices. Non-Bayesian approaches often contain hidden prior assumptions, which are difficult to probe.

Example Machine learning: Training work data, choice of learning cost functional often includes regularizers, structure of NN.

likelihood: Basis for χ^2 -fitting \equiv maximum likelihood fit in Bayesian view constant prior $P(\theta, I) = 1$

For a process with two outcomes Binomial $P_B(n_a, n_{tot}, p) = \frac{n_{tot}!}{n_a!(n_{tot}-n_a)!} p^{n_a} (1-p)^{n_{tot}-n_a}$
(relevant for classification tasks)

For continuous variables with observations fulfilling the central limit theorem

$P(D_n | \theta, I) \sim N(\mu(D_n, I, \theta), \sigma)$
(example body height & weight)

linear regression

$\mu = \theta_1 + \theta_2 \cdot D_n + \dots$
 \leftarrow predictor variable

(Model has "=" not "~")

In Bayesian setting: expressivity of model is uncertain
 (functional dependence on predictor & # of params)
 and parameters θ comp prior distribution.

(Jargon: linear models with Gaussian prior on slope parameter
 called ridge regression.)

Evaluating the posterior

① Quick and cheap: Maximum a posteriori (MAP) + Fisher curvature

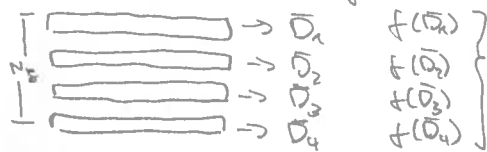
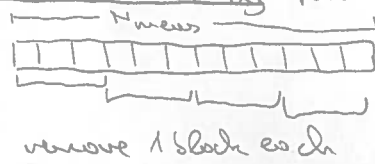
1) compute point estimate $\frac{\partial P(\theta|D, I)}{\partial \theta} \Big|_{\theta = \theta_{MAP}} = 0$ "most probable θ ,"

2) Quadratic approx to posterior $\partial^2 P / \partial \theta^2$ "how shallow is maximum?"

BEWARE: Gaussian approx of tailed distribution (or even multimodal)

\Rightarrow Explicit explicit dependence of $P(\theta|D, I)$ on data uncertainty & prior
 statistical error systematic error

Data uncertainty from (Jackknife) resampling:



$$\left. \begin{array}{l} \rightarrow \bar{D}_1 \quad f(\bar{D}_1) \\ \rightarrow \bar{D}_2 \quad f(\bar{D}_2) \\ \rightarrow \bar{D}_3 \quad f(\bar{D}_3) \\ \rightarrow \bar{D}_4 \quad f(\bar{D}_4) \end{array} \right\} \bar{f} \quad \sigma_{\bar{f}}^2 = \frac{N-1}{N} \sum_{i=1}^N (f(\bar{D}_i) - \bar{f})^2$$

If one needs $\sigma_{\bar{f}}$ for $f(\bar{D}, \sigma_{\bar{D}})$ auto correlations can be lessened by



Estimation of systematic error here "in systematic". Variation of prior
 probability from judicious choice.

② Full Rank - Chain Monte Carlo (from a prob. perspective 1701.02434)
 hands-on tutorials.

Powerful tool for inference and prediction since access to full posterior.
 This includes BOTH stat. & systematic uncertainties.

Inference $P(\theta|D, I)$ Posterior predictive distribution
 $P(\tilde{D}|D) = \int d\theta P(\tilde{D}|\theta, I) P(\theta|D, I)$

Allows cross-validation: estimate $P(\theta|D, I)$ on "training" subset of data
 and compare predictions $P(\tilde{D}|D)$ with "validation"
 data set.

How to summarize posterior? MAP = mode, mean, median?

Depends on context: If the loss incurred due to mis prediction scales with certain power \rightarrow optimal choice

loss $|P_{est} - P_{true}| \rightarrow$ median $(P_{est} - P_{true})^2 \rightarrow$ mean

Model uncertainties: von Neumann's elephant

"with 4 parameters I can fit an elephant with 5 I can make his funk wiggle."

- Intuitive:
- Occam's Razor: "models with fewer assumptions preferred"
 - ignorance of cause: potential causes that produce data in more ways are more plausible
 - adding more parameters improves fit (on training data)

Goal: model learns to generalize - predictive accuracy

(think of model fitting as data compression c.f. variational autoencoders)

Need measure of distance between distributions: via information entropy

Information: Reduction in uncertainty, derived from learning outcomes
 uncertainty $\hat{=}$ entropy $H(p) = - \langle \log(p) \rangle = - \int dx p(x) \log(p(x))$

["Maximum Entropy" principle looks for "least surprising distribution"]
 May be a good 1st guess BUT not always best choice when e.g. prior info avail.

Assume we know the correct target distribution $t(\theta)$

Concept of Divergence: Additional uncertainty induced by using a distribution p different from t .

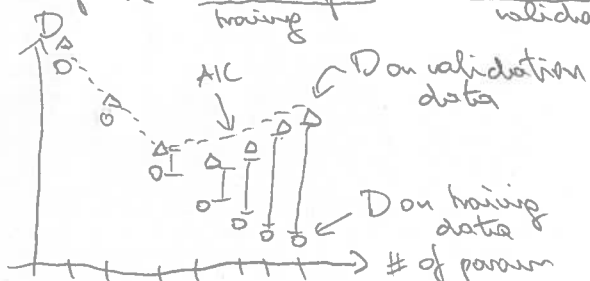
\hookrightarrow Kullback-Leibler $D_{KL}(t, p) = \int dx t(x) \left(\log(t(x)) - \log(p(x)) \right)$

$\begin{matrix} \text{cross-entropy} \\ \uparrow \quad \downarrow \\ \log(t(x)) \quad \log(p(x)) \\ \text{entropy of } t \end{matrix}$

In practice t is unknown so instead often used deviance

$D(p) = -2 \int dx \log(p(x))$ Simplest model comp. based on $p = P(D|\theta, Z)$
 D is the log likelihood.

Compare in-sample to out-sample deviance for model with diff. # of param.



Difference offers a measure for predictive accuracy.

[BEWARE: at this point prior not considered]

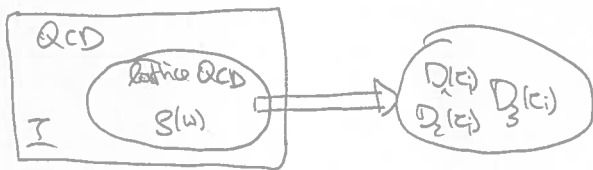
Akaike Information criterion: attempts to estimate the prediction variance for models with flat prior $AIC = D_{train} + 2 \#param. \approx D_{validate}$

For a more robust measure of pred. accuracy when using non-flat priors: WAIC (widely applicable IC) or LOO (leave-one-out cross validation) for more details see 1507.04544

How to compare models \Rightarrow average over models using info criteria assign a model weight $w_m = \frac{\exp(-\frac{1}{2} IC_m)}{\sum_m \exp(-\frac{1}{2} IC_m)}$ works for IC's based on deviation scale, exponentiation gives prob. sample from ensemble of models weighted by w_m .

Interpretation of w_m^{AIC} : "estimates the probability that the model will make the best prediction on new data conditional on set of models considered."

QCD Example: Bayesian inference of spectral functions



- wish to infer $g(w)$ from $D(z)$
- ill-posed since data sparse to approximate $C(z_i) + \epsilon_i$

① infinitely many $g(w)$'s reproduce $C(z)$ within errors. ② Naive inversion leads to exponential dependence of resulting g on error in D

$P(g|D,I) = \frac{P(D|g,I) P(g,I)}{P(D,I)}$ Model for likelihood given by QCD $D \sim \mathcal{N}(\mu, \sigma)$ $\mu(z) = \int dw K(w,z) g(w)$

No uncertainty in $K(w,z)$ since ab-initio input.

$K(w,z) = \exp(-wz)$ $T=0$ spectroscopy, NRQCD at $T=0$ & $T>0$, PDFs from hadronic tensor
 $K(w,z) = \frac{\cosh(w(z-\beta/2))}{\sinh(w\beta/2)}$ $T>0$ hadrons
 $K(x,y) = \cos(x,y)$ $T=0$ pseudo PDFs

Note that correlations not only along MC time (auto correlation) but also between the means of correlator at different τ 's. i.e. covariance matrix var. f. mean $Cov_{ij} = \frac{1}{N(N-1)} \sum_{\ell=1}^{N-1} (D_\ell(\tau_i) - \bar{D}(\tau_i))(D_\ell(\tau_j) - \bar{D}(\tau_j))$ is not diagonal assuming NO autocorr

Instead of modeling cross-correlations \rightarrow decorrelate data

$Cov = R \text{diag}(\sigma_i) R^T$ transform data in basis where Cov is diagonal \rightarrow independent Gaussians.

Prior distribution Acts as a regulator for the ill-posed inverse problem. Prior knowledge picks a most probable g among the degenerate extrema of the likelihood



In the literature: Exploit positivity of g

$$P_T(g, I) \propto \exp(-\alpha \int dw (g(w) - m(w))^2)$$

Tikhonov Gaussian prior (does not exploit positivity)

Jargon: $m(w)$ "default model" parametrizes maximum of prior. α width parameter

$$P_{TV}(g, I) \propto \exp(-\alpha \int dw | \frac{dg}{dw} |)$$

L_1 / total variance denoising (also does not use positivity)

$$P_{MEM}(g, I) \propto \exp\left(\alpha \int dw \left[g(w) - m(w) - g(w) \log \left(\frac{g(w)}{m(w)} \right) \right] \right)$$

Shannon-Jaynes entropy

see e.g.

Phys. Rep. 269 (1996) 133

justified via axioms: Do not introduce correlations where none are present in the data. Challenge: originally from 2d image reconstruction and axioms refer to that setting & assumes g is probability distribution w/o units.

$$P_{BR}(g, I) \propto N \exp\left(\int dw \alpha(w) \left[1 - \frac{g(w)}{m(w)} + \log \left(\frac{g(w)}{m(w)} \right) \right] \right)$$

see Y. Burnier, A.R. 1307.6106

justified via axioms: Scale invariance - units of g may not matter and smoothness of g is assumed. Weakest known regulator: let's the data speak but tendency to show more ringing than other priors.

In Bayesian continuous limit ($N_z \rightarrow \infty, \Delta D \rightarrow 0$) Bayes theorem guarantees correct inference of g . For finite # of datapoints N_z and finite error ΔD choice of prior determines how efficiently limit is approached and what artifacts encountered on the way (oversmoothing vs. ringing e.g.)