

A high-level overview of Bayesian Statistics

Bayesian Idea: Treat a parameter as having an unknown Distribution to be estimated, rather than as an unknown fixed number to be estimated.

Ex 1 Predict end-of-season win percentage from mid-season Wins and Losses,

Beta-Binomial model

$$\begin{cases} W \sim \text{Binomial}(n, P) \\ P \sim \text{Beta}(\alpha, \beta) \end{cases}$$

When we model the latent team's win probability P using a Beta Prior, we encode the prior information that P is more likely to be near $1/2$ (say, in $.3, .7$) than to be very near 0 or 1 .

Then, we found using Bayes Rule that the posterior distribution $P|W, L$ is

$$P|W, L \sim \text{Binomial}\left(n+d+\beta-2, \frac{w+\alpha-1}{W+L+d+\beta-2}\right)$$

and the Bayes estimate is the posterior mean

$$\hat{P}^{(\text{Bayes})} = \mathbb{E}[P|W, L] = \frac{w + (\alpha - 1)}{W + L + (\alpha + \beta - 2)}$$

Ex 2 Predict end-of-season batting average from mid-season batting average and number of at-bats.

Normal-Normal Model

$$\left\{ \begin{array}{l} i = \text{player } i \\ H_i = \# \text{ hits}, N_i = \# \text{ at-bats}, X_i = \frac{H_i}{N_i} \\ X_i \sim N(\mu_i, \sigma_i^2) \\ \sigma_i^2 = C/N_i, C \text{ known} \\ \mu_i \sim N(\bar{\mu}, \tau^2) \end{array} \right.$$

When we model player i 's latent quality μ_i using a Normal Prior, we encode the prior information that player i is also a baseball player, allowing us to share strength across players.

Then, we found using Bayes Rule that the posterior distribution $\mu_i | X$ is

$$\mu_i | X \sim N\left(\frac{\frac{X_i + \bar{\mu}}{\sigma_i^2 + \tau^2}}{\frac{1}{\sigma_i^2} + \frac{1}{\tau^2}}, \frac{1}{\frac{1}{\sigma_i^2} + \frac{1}{\tau^2}}\right)$$

and the Bayes estimate is the posterior mean

$$\hat{\mu}_i^{(\text{Bayes})} = \frac{\frac{X_i}{\sigma_i^2} + \frac{\bar{\mu}}{\tau^2}}{\frac{1}{\sigma_i^2} + \frac{1}{\tau^2}}$$

using $\bar{\mu}$ and τ^2 since μ, τ are unknown.

Ex 3

Bayesian Regression

Suppose we have a regression model

$$y = X\beta + \varepsilon \text{ with mean zero noise } E[\varepsilon] = 0.$$

For instance, our Park effects model (so we want to estimate the park effects β), or a Power Score model like Bradley Terry.

A Bayesian puts a *prior distribution* on the parameters β (and, a distribution on the error term ε) to end up with a *posterior distribution* $\beta|X$ or a *posterior mean* $E[\beta|X]$, quantifying our best guess of β after seeing the data X , while also incorporating *our prior information*.

* "Standard" Regression Setup:

$$y = X\beta + \varepsilon$$

$$\begin{aligned} \text{Homoscedasticity} \rightarrow \varepsilon_i &\stackrel{\text{iid}}{\sim} N(0, \sigma^2) \\ &\Rightarrow \varepsilon \sim N(0, \sigma^2 I) \end{aligned} \quad \left. \right\} \Rightarrow y \sim N(X\beta, \sigma^2 I)$$

* Bayesian Regression ex1: *uninformative Uniform Prior* $\beta \sim I$

* Bayesian Regression ex2: *Normal Prior* $\beta \sim N(0, \lambda)$

Bayesian Regression v1

$y \sim N(X\beta, \sigma^2 I)$
 $\beta \propto 1$ uniform prior
 Find posterior dist $\beta | X$

$$P(\beta | y) = \frac{P(y | \beta) P(\beta)}{P(y)} \quad \text{by Bayes Rule}$$

$$\propto P(y | \beta) \cdot P(\beta) \quad \text{since } P(y) \text{ doesn't depend on } \beta$$

$$\propto P(y | \beta) \quad \text{since } P(\beta) \propto 1$$

$$\propto \prod_{i=1}^n \exp\left(-\frac{1}{2\sigma^2} (y_i - X_i^\top \beta)^2\right)$$

$$= \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - X_i^\top \beta)^2\right).$$

Bayes estimate of β is to choose the β which maximizes the posterior probability

Maximum A-posteriori (MAP) :

$$\hat{\beta}^{(\text{Bayes})} = \operatorname{argmax}_{\beta} P(\beta | y)$$

$$= \operatorname{argmax}_{\beta} \sum_{i=1}^n (y_i - X_i^\top \beta)^2 = \hat{\beta}^{(\text{OLS})} = (X^\top X)^{-1} X^\top y.$$

* With an uninformative uniform prior, Bayesian Regression is the same as Ordinary least squares!

Bayesian Regression v2

$$P(\beta | y) = \frac{P(y | \beta) P(\beta)}{P(y)}$$

by Bayes
Rule

$$\left\{ \begin{array}{l} y \sim N(X\beta, \sigma^2 I) \\ \beta \sim N(0, \frac{1}{\lambda} I) \end{array} \right.$$

Find posterior dist $\beta | X$

$\propto P(y | \beta) \cdot P(\beta)$ since $P(y)$ doesn't depend on β

$$= P(N(X\beta, \sigma^2 I) = y) \cdot P(\beta = N(0, \lambda))$$

$$\propto \prod_{i=1}^n \exp\left(-\frac{1}{2} \frac{(x_i^\top \beta - y_i)^2}{\sigma^2}\right) \cdot \prod_{j=1}^p \exp\left(-\frac{1}{2} \frac{\beta_j^2}{1/\lambda}\right)$$

$$= \exp\left(-\frac{1}{2} \left[\lambda \sum_{j=1}^p \beta_j^2 + \sum_{i=1}^n \frac{(x_i^\top \beta - y_i)^2}{\sigma^2} \right]\right)$$

$$= \exp\left(-\frac{1}{2\sigma^2} \left[\lambda \sum_{j=1}^p \beta_j^2 + \sum_{i=1}^n (y_i - x_i^\top \beta)^2 \right]\right).$$

Bayes estimate of β is to choose the β which maximizes the posterior probability (Maximum A-posteriori (MAP)):

$$\hat{\beta}^{(\text{Bayes})} = \underset{\beta}{\operatorname{argmax}} P(\beta | y) = \underset{\beta}{\operatorname{argmax}} \lambda \sum_{j=1}^p \beta_j^2 + \sum_{i=1}^n (y_i - x_i^\top \beta)^2$$

$$\left(\text{letting } \lambda < \frac{1}{\sigma^2}\right) \quad = \hat{\beta}^{(\text{Ridge})} = (X^\top X + \lambda I)^{-1} X^\top y.$$

* With a Normal prior, Bayesian Regression is Ridge Regression. This makes perfect sense: $\beta \sim N(0, \lambda I)$ encourages β to be closer to 0. λ controls by how much, and that is the same as penalizing β !

- * So far we have mainly focused on finding the **posterior mean**, i.e. our best estimate of the parameters after seeing the data.
- * To **make decisions** in sports (e.g. player valuation or play selection), we need not only know our best estimate of the value of the player/play/decision, but also **uncertainty quantification** (e.g., error bars or prediction intervals) to describe how confident/certain we are about a value.

Ex Suppose we have a prediction of end-of-season win percentage to be 74%.
A 95% prediction interval of $[72\%, 76\%]$ is much different from $[60\%, 88\%]$ and from $[0\%, 100\%]$.

Bayesian Idea Estimate the FULL Posterior Distribution of an unknown parameter. This gets us error bars (uncertainty quantification), with which we can create more complete decisions.

* In a general Bayesian statistical model, the parameters don't all have Gaussian priors or likelihoods, and we can't write on paper a closed-form analytical solution for the posterior distribution.

* Typically need to approximate the posterior distribution using MCMC sampling methods like

{ Gibbs sampling
Hamiltonian Monte Carlo
No U-Turn Sampling,

which we won't cover here (take Shane's class), using a probabilistic programming language

[Stan
NumPyro]

* Next Time : a full example of Bayesian modeling in sports.

Takeaway Create a fully Bayesian probability model, which allows quantify both the effect size and uncertainty of a phenomenon of interest by approximating the full posterior distribution.