# Outline

1) Conjugate Priors

2) Where does the prior come from?

3) Bayesian pros and cons

# Conjugate Priors

If the posterior is from the same family as the prior, we say the prior is conjugate to the likelihood.

Most common in exp. fam.s: Suppose

$$X_i \mid \eta \overset{iid}{\sim} p_\eta(x) = e^{\eta' T(x) - A(\eta)} h(x) \qquad \eta \in \Xi \subseteq \mathbb{R}^s \quad i = 1, \ldots, n$$

For carrier $\lambda_0(\eta)$, define $s+1$-dim family:

$$\lambda_{k\mu, k}(\eta) = e^{k\mu' \eta - kA(\eta) - B(k\mu, k)} \lambda_0(\eta)$$

Suff. stat $\begin{pmatrix} \eta \\ -A(\eta) \end{pmatrix} \in \mathbb{R}^{s+1}$

Nat. param. $\begin{pmatrix} k\mu \\ k \end{pmatrix}$

Then

$$\lambda(\eta \mid x_1, \ldots, x_n) \underset{\eta}{\propto} \left( \prod_{i=1}^{n} e^{\eta' T(x_i) - A(\eta)} h(x_i) \right)$$

$$\cdot e^{k\mu' \eta - kA(\eta) - B(k\mu, k)} \lambda_0(\eta)$$

$$\underset{\eta}{\propto} e^{(k\mu + \Sigma T(x_i))' \eta - (k+n) A(\eta)} \lambda_0(\eta)$$

$$= \lambda_{k\mu + n\bar{T}, \, k+n}(\eta)$$

where $\bar{T}(x) = \frac{1}{n} \sum_{i=1}^{n} T(x_i)$

<u>Interp.</u>: If we:
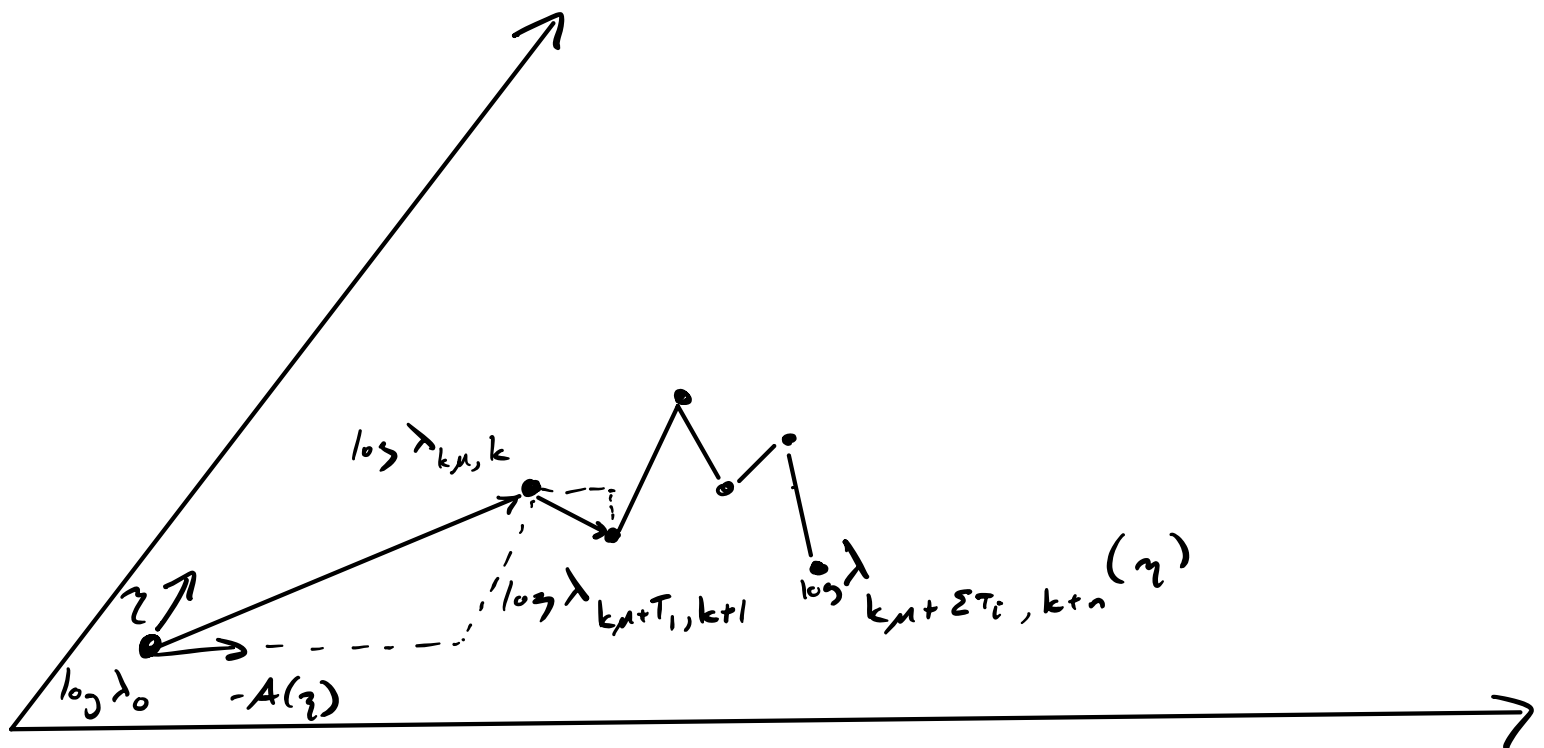
1) Take prior $\lambda_{k\mu,k}$, observe avg. suff. stat.
$\bar{T}$ on sample size $n$

OR 2) Take prior $\lambda_0$, observe avg. suff. stat.
   a) $\mu$ on sample size $k$    (pseudo-data)
   b) $\bar{T}$ on sample size $n$

We get posterior $\lambda_{k\mu+n\bar{T}, k+n}$



Often $\dfrac{k\mu + n\bar{T}}{k+n}$ is Bayes est. for $\mathbb{E}_\eta T$,

then $\hat{\mu}_{post} = \bar{T} \cdot \dfrac{n}{k+n} + \mu \cdot \dfrac{k}{k+n}$

$\uparrow$ UMVUE from data

$\uparrow$ UMVUE from pseudo data

(If so then $\lambda_0$ not a proper prior. Why?)

# Conjugate Prior Examples

| Likelihood | Prior |
|---|---|
| $X_i \mid \theta \sim \text{Binom}(n, \theta)$ $= \theta^x (1-\theta)^{n-x} \binom{n}{x}$ | $\theta \sim \text{Beta}(\alpha, \beta)$ $= \theta^{\alpha-1} (1-\theta)^{\beta-1} \dfrac{\Gamma(\alpha)\, \Gamma(\beta)}{\Gamma(\alpha+\beta)}$ |
| $X_i \mid \theta \sim N(\theta, \sigma^2)$ $\quad$ (σ² known) $= \dfrac{1}{\sqrt{2\pi\sigma^2}} \, e^{-(\theta-x)^2 / 2\sigma^2}$ | $\theta \sim N(\mu, \tau^2)$ $= \dfrac{1}{\sqrt{2\pi\tau^2}} \, e^{(\theta-\mu)^2 / 2\tau^2}$ |
| $X_i \mid \theta \sim \text{Pois}(\theta) \qquad x = 0, 1, \dots$ $= \dfrac{\theta^x e^{-\theta}}{x!}$ | $\theta \sim \text{Gamma}(\upsilon, s) \qquad \theta > 0$ $= \dfrac{1}{\Gamma(\upsilon)\, s^2} \, \theta^{\upsilon-1} \, e^{-\theta/s}$ |

## Gamma / Poisson :

$$\lambda(\theta \mid x) \propto_\theta \theta^{\upsilon - 1 + \Sigma x_i} \, e^{-(s' + n)\theta}$$

$$= \text{Gamma}\left(\upsilon + \Sigma x_i, \ (s' + n)^{-1}\right)$$

$$\Rightarrow k = s^{-1}, \quad \mu = \upsilon s$$

$$\lambda_0(\theta) = \theta^{-1} \qquad (\text{not normalizable})$$

# Where does prior come from?

Biggest issue with Bayes in practice is how to choose prior

In general, can't check goodness of fit: 1 draw of $\theta \sim \lambda$, not even directly observed.

Various ideas of how to do it:

## 1) Prior experience:

Ex. A/B testing in tech. co.s, estimating eff. of 400,000 SNPs on trait, estimating "true" 3PT% for basketball players

Prior is (relatively) non-controversial
- Can fit from data (leads to hierarchical / Empirical Bayes)
- Can test validity of prior b/c we have many draws from it

Works when encountering similar problems repeatedly (but are they really similar?)

## 2) Subjective beliefs:

Prior reflects <u>epistemic</u> uncertainty

Posterior = rational updating of beliefs.

<u>Pros</u>:

- Can't be wrong about your own opinion!
- Can bring to bear hard-to-formalize knowledge from outside the data

<u>Issues</u>:

- Philosphical conundrum: is the mass of a particle really "random"?

- Scientists find subjectivity offputting (reporting posterior is just reporting an opinion)

- Generally impossible to write down your beliefs about joint dist. of $\theta \in \mathbb{R}^{10}$

- What if people are systematically overconfident?

<u>But</u> This is the most philosophically coherent account of statistics. (Coin flip demo)

# 3) Convenience Prior

Bayes computations can be very hard:
generically, it's very hard to compute
the normalizing constant $\int_{\Omega} \lambda(\theta) p_{\theta}(x) d\theta$

If $\dim(\Omega)$ large, posterior $\approx 0$ for most of $\Omega$

Lots of Bayesian research is computational
(e.g., MCMC) and great progress has been
made.

Helps to use conjugate priors where possible

But then the subjective account basically
falls apart

# 4) "Objective" Priors

Suppose $X_i | \theta \overset{iid}{\sim} N(\theta, 1)$   $i = 1, \dots, n$

"Natural" choice is   $\theta \sim$ "flat prior"

$$\lambda(\theta) \propto_\theta 1$$

This prior is improper but it's ok:

$$\lambda(\theta | x) \propto_\theta e^{\theta \Sigma X_i - n\theta^2/2}$$

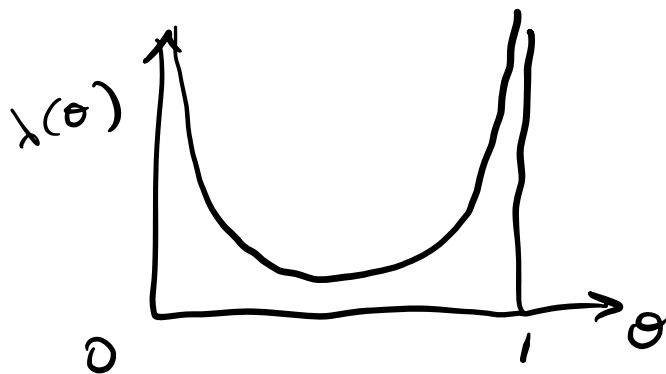$$\propto_\theta N(\bar{X}, n^{-1})$$

More generally, could always use flat prior
Arises naturally as limit of $\theta \sim N(0, \tau^2)$, $\tau^2 \to \infty$

Problem: Flat prior is not flat anymore if
we reparameterize $\theta$!

Jeffereys proposed using $\lambda(\theta) \propto_\theta |J(\theta)|^{1/2}$

This is also the Jeffereys prior after any
reparameterization

Binomial: Jeffereys prior looks like



$\lambda(\theta)$

0          1   $\theta$

"Objective" ??

# Gaussian Sequence Model

$X \sim N_d(\mu, I_d)$     $\mu \in \mathbb{R}^d$

Jeffereys prior is flat: $\lambda(\mu) \equiv 1$

$$\lambda(\mu | x) = N_d(X, I_d)$$

$$\Rightarrow \mathbb{E}[\mu | x] = X$$

Reasonable $\widehat{\text{estimator}}$: coincides with UMVU, MLE,
<span style="color:red">(but inadmissible)</span>                               minimax, ...

What about $\rho^2 = \|\mu\|^2$ ?

$$\mu \sim N_d(X, I_d) \Rightarrow \mathbb{E}[\|\mu\|^2 | X] = \|X\|^2 + d$$

Recall UMVUE was $\|X\|^2 - d$
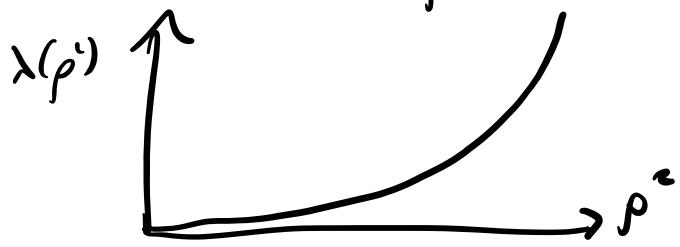
So    Bayes est. = UMVUE + $2d$

$$MSE(\theta; \delta_\Lambda) = Var_\theta(\delta_{UMVU}) + 4d^2$$

What went wrong?

$$\mu \sim N_d(0, \tau^2 I_d) \Rightarrow \rho^2 \sim \tau^2 \chi_d^2$$

Jeffereys prior takes $\tau^2 \to \infty$

$$\lambda(\rho^2) \underset{\rho^2}{\propto} (\rho^2)^{(d-1)/2}$$

$\lambda(\rho^2)$                                        "Agnostic" ?

# Advantages of Bayes

Despite difficulties above, Bayes has some major advantages over other approaches:

1) Estimator is defined straightforwardly:

$$\delta_\triangle (x) = \arg\min_d \int L(\theta, d) \, \lambda(\theta | x) d\theta$$

Problem is reduced entirely to computation

May be difficult to compute but in principle we can find it for generic $L, \lambda, \mathcal{P}, g(\theta)$

Don't need to rely on special structure like complete suff. stat, U-estimable $g$, exp. fam., simple $L$

<u>Gives us freedom to</u>:
- Use highly expressive & complex models
- Use the $L$ we actually care about
- Incorporate background subj-matter knowledge

Unparalleled expressive power for systems we know a lot abt. already

2) Appealing optimality property:
  Even if we don't "believe" prior, Bayes
    est. has best avg. case risk
  Bayes estimators are usually admissible

5) Detailed output: entire posterior, joint distr.
  over all parameters
  One computation (e.g. large sample from posterior)
    leads to estimates of any $g(\theta)$ we
    can think of

... many more

# Cons

1) Difficulty of choosing $\Lambda$, esp. in high dim.
   Avg.-case performance doesn't ensure good
      performance for the real (distribution of) $\theta$

   If we choose $\Lambda$ poorly, we might get no mass

   near    true  $\theta$

2) Flipside of ability to specify model in full
     detail  is  <u>requirement</u>  that  we  must do so.

   e.g. nonparametric estimation of $g(P) = \mathbb{E}_P X_i$
      $X_i \overset{iid}{\sim} P$.    $\bar{X}$ is UMVUE, natural choice.
      To get started on Bayes, must define
        prior  over all  distr. on $\mathbb{R}$.
   Frequentist approaches let us stay more
       parsimonious  with  our  assumptions.


... many  more