

# LECTURE 4: BAYESIAN ANALYSIS

CHRIS CONLON

NYU STERN

FEBRUARY 27, 2019

## QUICK REFRESH: BAYES RULE

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Given a **positive test result** what is the probability a patient actually has cancer?

**Table:** Test Accuracy

	Cancer (1%)	No Cancer (99%)
Positive Test	80%	9.6%
Negative Test	20%	90.4%

## QUICK REFRESH: BAYES RULE

Calculate  $Pr(\text{Cancer} \& \text{PositiveTest})$  and  $Pr(\text{NoCancer} \& \text{PositiveTest})$

**Table:** Joint Probabilities

	Cancer (1%)	No Cancer (99%)
Positive Test	$(0.8)(0.01)=0.008$	$(0.9)(0.096)=0.09504$
Negative Test	$(0.2)(0.01)=0.002$	$(0.9)(0.904)=0.89496$

$$Pr(\text{Cancer}|\text{PositiveTest}) = \frac{Pr(\text{Cancer}, \text{PosTest})}{Pr(\text{Cancer}, \text{PosTest}) + Pr(\text{NoCancer}, \text{PosTest})} = .008 / .10304 = 0.0776$$

- Suppose that we toss a coin several times with  $x_i \in \{H, T\} = \{1, 0\}$
- $\mathbf{X} = \{H, T, H, H, \dots\}$ .
- Suppose that the probability of heads  $Pr(x_i = H) = p$ .
- What is the likelihood of an observed sequence of  $\mathbf{X}$ ? where  $x_i$  are I.I.D.

$$Pr(x_i|p) = p^{x_i}(1-p)^{1-x_i}$$

$$Pr(\mathbf{X}|p) = p^{\sum_i x_i}(1-p)^{\sum_i (1-x_i)}$$

# INTRODUCTION: MLE FOR COIN TOSS

Can construct the **log likelihood** and find the MLE.

$$\ell(\mathbf{X}|p) = \left(\sum_i x_i\right) \ln p + \left(N - \sum_i x_i\right) \ln(1 - p)$$

$$\frac{\partial \ell(p)}{\partial p} = \left(\sum_i x_i\right) \frac{1}{p} - \left(N - \sum_i x_i\right) \frac{1}{1 - p} = 0$$

$$\frac{1 - p}{p} = \frac{\left(\frac{1}{N} \cdot N - \frac{1}{N} \cdot \sum_i x_i\right)}{\frac{1}{N} \cdot \sum_i x_i} \rightarrow \hat{p} = \frac{1}{N} \cdot \sum_i x_i$$

# INTRODUCTION: MLE FOR COIN TOSS

Can also construct the properties of  $\hat{p}$ .

$$\mathbb{E}[\hat{p}] = \mathbb{E}\left[\frac{1}{N} \cdot \sum_i x_i\right] = \left[\frac{1}{N} \cdot \sum_i \mathbb{E}x_i\right] = \mu_x = p_0$$

$$\mathbb{V}[\hat{p}|\mathbf{X}] = \mathbb{V}\left[\frac{1}{N} \cdot \sum_i x_i\right] = \frac{1}{N^2} \cdot \sum_i \mathbb{V}(x_i) = \frac{N}{N^2} p(1-p)$$

Which gives us a CI of:  $\left(\bar{x} \pm 1.96 \cdot \sqrt{\frac{1}{N}\bar{x}(1-\bar{x})}\right)$

A different idea:

- Start with a (diffuse) initial guess for the distribution of  $p$ :  $f_P(p)$ .
- Incorporate information from likelihood:  $f(x_i|p)$
- Construct **posterior density** estimate  $f(p|x_i)$ .
  - ▶ This doesn't characterize a best estimate  $\hat{p}$  but a full distribution.
  - ▶ We can calculate  $\mathbb{E}[p|x_i]$  or  $\mathbb{V}[p|x_i]$  or any other functions of the posterior density.
- Challenge: How to choose initial  $f_P(p)$ .

# BAYESIAN STATISTICS: BRIEF INTRODUCTION

One possible guess is the uniform distribution  $f(x) = 1$  on  $0 \leq x \leq 1$ .

- **Marginal/Prior Distribution:**  $f_P(p) = 1$  for  $0 \leq p \leq 1$ .

- **Conditional Distribution/Likelihood:**  $f_{X|P}(x|p) = p^x(1-p)^{1-p}$

- **Joint Distribution :**

$$f_{X,P}(x, p) = f_{X|P}(x|p) \cdot f_P(p) = p^x(1-p)^{1-p} \cdot 1 = p^x \cdot (1-p)^{1-p}$$

- ▶ This is only defined for  $p \in [0, 1]$  and  $x \in \{0, 1\}$ . It is zero elsewhere.

- What about **Marginal Distribution** for  $x$ ?

$$\begin{aligned} \int_0^1 f_{PX}(x, p) dp &= x \cdot \int_0^1 p dp + (1-x) \cdot \int_0^1 (1-p) dp \\ &= x \cdot \frac{1}{2} + (1-x) \cdot \frac{1}{2} = \frac{1}{2} \propto 1 \end{aligned}$$



The object we are usually interested in is the **Posterior Distribution**

$$f_{P|X}(p|x) = \frac{f_{X|P}(x|p) \cdot f_P(p)}{\int_0^1 f_{X|P}(x|p) \cdot f_P(p) dp} = 2p^x(1-p)^{1-x} \propto p^x(1-p)^{1-x}$$

- We are back at the p.m.f. of the **Bernoulli** which is maybe comforting.
- This is true because  $f_X(x) \propto 1$  and  $f_P(p) \propto 1$ .
- $f_{P|X}(p|x=0) = (1-p)$  and  $f_{P|X}(p|x=1) = p$ .

- Let's try a different **prior distribution** than the uniform we used last time. This time we will use a  $Beta(\alpha, \beta)$  distribution:

$$f_P(p|\alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha) \cdot \Gamma(\beta)} p^{\alpha-1} (1-p)^{\beta-1}$$

$$\mathbb{E}[p|\alpha, \beta] = \frac{\alpha}{\alpha + \beta}$$

$$\mathbb{V}[p|\alpha, \beta] = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}$$

- This has the advantage that it places nicely with the Binomial.
- Consider  $\alpha = 16, \beta = 8$ . This gives  $\mathbb{E}[p] = \frac{2}{3}$  and  $\text{SE}[p] = 0.094$ .

Consider the case where  $x = 1$  (we get one piece of new data).

$$f_P(p) \cdot f_{X|P}(x|p) = \frac{\Gamma(\alpha + \beta)}{\underbrace{\Gamma(\alpha) \cdot \Gamma(\beta)}_{C(\alpha, \beta)}} p^{\alpha-1} (1-p)^{\beta-1} \cdot p \propto p^{\alpha} (1-p)^{\beta-1}$$

- The resulting distribution is now  $(p|x = 1) \sim \text{Beta}(\alpha + 1, \beta)$ .
- Our posterior has mean = 0.68 and SE = 0.091.
- Estimate of mean increases and SE decreases.
- Likewise if  $x = 0$  we get  $(p|x = 0) \sim \text{Beta}(\alpha, \beta + 1)$
- There is a **conjugacy** relationship between the Beta and the Binomial.

# GENERAL CASE

$$\overbrace{f_{\theta|X}(\theta|X)}^{\text{posterior}} = \frac{\overbrace{f_{X,\theta}(X,\theta)}^{\text{joint}}}{\underbrace{f_X(X)}_{\text{marginal of } X}} = \frac{\overbrace{f_{X|\theta}(X|\theta)}^{\text{likelihood}} \cdot \overbrace{f_\theta(\theta)}^{\text{prior}}}{\int f_{X|\theta}(X|\theta) \cdot f_\theta(\theta) d\theta}$$

There is a shortcut because the denominator doesn't depend on  $\theta$

$$f_{\theta|X}(\theta|X) \propto f_{X|\theta}(X|\theta) \cdot f_\theta(\theta) = \mathcal{L}(\theta|X) \cdot f_\theta(\theta)$$

We can cheat because there exists a constant  $c$  so that  $c \int \mathcal{L}(\theta|X) \cdot f_\theta(\theta) d\theta = 1$ .

## A NORMAL EXAMPLE

Assume  $X \sim N(\mu, 1)$  and  $\mu \sim N(0, 100)$ . What is  $f_{\mu|X}(\mu|X = x)$ ?

$$\begin{aligned}f_{\mu|X}(\mu|X) &\propto \exp\left(-\frac{1}{2}(x - \mu)^2\right) \cdot \exp\left(-\frac{1}{2 \cdot 100}\mu^2\right) \\&= \exp\left(-\frac{1}{2}(x^2 - 2x\mu + \mu^2 + \mu^2/100)\right) \\&\propto \exp\left(-\frac{1}{2(100/101)}(\mu - (100/101)x)^2\right)\end{aligned}$$

It happens that  $(\mu|X) \sim N(100x/101, 100/101)$ .

In general the posterior will not be well defined.

# KALMAN UPDATE: A MORE COMPLICATED NORMAL

Assume  $X \sim N(\mu, \sigma^2)$  with  $\sigma^2$  known.  $\mu \sim N(\mu_0, \tau^2)$ . What is  $f_{\mu|X}(\mu|X=x)$ ?

$$\begin{aligned} f_{\mu|X}(\mu|X) &\propto \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right) \cdot \exp\left(-\frac{1}{2\tau^2}(\mu-\mu_0)^2\right) \\ &\propto \exp\left[-\frac{1}{2}\left(\frac{x^2}{\sigma^2} - \frac{2x\mu}{\sigma^2} + \frac{\mu^2}{\sigma^2} + \frac{\mu^2}{\tau^2} - \frac{2\mu\mu_0}{\tau^2} + \frac{\mu_0^2}{\tau^2}\right)\right] \\ &\propto \exp\left[-\frac{1}{2}\left(\mu^2\frac{\sigma^2 + \tau^2}{\tau^2\sigma^2} - \mu\frac{2x\tau^2 + 2\mu_0\sigma^2}{\tau^2\sigma^2}\right)\right] \\ &\propto \exp\left[-\frac{1}{2(1/(\tau^2 + 1/\sigma^2))}\left((\mu - (x/\sigma^2 + \mu_0/\tau^2)) / (1/\sigma^2 + 1/\tau^2)\right)^2\right] \end{aligned}$$

The resulting distribution is Normal with mean and variance

$$\mathbb{E}[\mu|X=x] = \frac{\frac{x}{\sigma^2} + \frac{\mu_0}{\tau^2}}{\frac{1}{\sigma^2} + \frac{1}{\tau^2}}, \quad \mathbb{V}(\mu|X) = \frac{1}{\frac{1}{\sigma^2} + \frac{1}{\tau^2}}$$

# KALMAN UPDATE: A MORE COMPLICATED NORMAL

Despite being a giant mess this makes sense:

$$\mathbb{E}[\mu|X = x] = \frac{\frac{x}{\sigma^2} + \frac{\mu_0}{\tau^2}}{\frac{1}{\sigma^2} + \frac{1}{\tau^2}}, \quad \mathbb{V}(\mu|X) = \frac{1}{\frac{1}{\sigma^2} + \frac{1}{\tau^2}}$$

- Posterior mean is a weighted average of **prior mean** and **sample mean**.
- Weights depend on **precision** of two samples.
- Posterior **Precision** is sum of precision of each sample  $\frac{1}{\mathbb{V}(\cdot)}$
- Probably we want to choose a relatively **uninformative** prior with large  $\tau^2$ .
- $\tau^2 \rightarrow \infty$  implies an **improper prior distribution** because it no longer integrates to one. But because of  $\propto$  still mostly ok.

This is straightforward:

$$p(\theta|X_1, \dots, X_N) \propto \mathcal{L}(\theta|X_1, \dots, X_N) \cdot p(\theta)$$

- Still depends on: **prior**, **likelihood** to construct **posterior**.
- Can update one observation at a time or all at once.



## Bernstein von-Mises Theorem

*A posterior distribution converges as you get more and more data to a multivariate normal distribution centred at the maximum likelihood estimator with covariance matrix given by  $n^{-1}I(\theta_0)^{-1}$ , where  $\theta_0$  is the true population parameter (Edit: here  $I(\theta_0)$  is the Fisher information matrix at the true population parameter value).*

Under these conditions (and some more):

1. MLE is consistent
2. Fixed number of parameters
3.  $\theta_0$  in interior of  $\Theta$  (true value of SD can't = 0).
4. The prior density must be non-zero in a neighborhood of  $\theta_0$ .
5. log-likelihood needs to be smooth (two derivatives at the true value and more)

# CONJUGATE PRIORS

# WHAT IS A CONJUGATE PRIOR?

For a given **likelihood**  $f_{X|\theta}(x|\theta)$  we can choose a **prior**  $f_{\theta}(\theta)$  so that the **posterior** is proportional to a known parametric distribution.

- This makes life easy because now the posterior has a known parametric distribution (normal, beta, gamma, etc.)
- Other than convenience, this alone doesn't tell us that our choice of  $f_{\theta}(\theta)$  is the **best** prior by any metric.
- Using a non-conjugate prior is entirely defensible, just less convenient.

Prior has **hyper parameters**  $(\alpha, \beta)$ :

$$Pr(p|\alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} p^{\alpha-1} (1-p)^{\beta-1}, \quad E[p] = \frac{\alpha}{\alpha + \beta}$$

Likelihood

$$Pr(y = k|n, p) = \binom{n}{k} p^k (1-p)^{n-k}$$

Posterior

$$\begin{aligned} f(p|n, k, \alpha, \beta) &\propto p^{k+\alpha} (1-p)^{n-k+\beta} \\ &\sim \text{Beta}(\alpha + k, \beta + n - k) \end{aligned}$$

# DIRICHLET-MULTINOMIAL

Prior defined on **unit simplex** when  $\sum_{i=1}^k p_i = n$

$$Pr(p_1, \dots, p_K; \alpha_1, \dots, \alpha_K) = \frac{1}{B(\alpha)} \prod_{i=1}^K p_i^{\alpha_i - 1}$$

Likelihood

$$L(x_1, \dots, x_k | n, p_1, \dots, p_k) = \frac{n!}{x_1! \dots x_k!} p_1^{x_1} \times \dots \times p_k^{x_k}$$

Posterior

$$\begin{aligned} f(p_1, \dots, p_k | x_1, \dots, x_k, \alpha_1, \dots, \alpha_k) &\propto p_1^{x_1 + \alpha_1} \times \dots \times p_k^{x_k + \alpha_k} \\ &\sim \text{Dirichlet}(\alpha_k + x_k) \end{aligned}$$

- Dirichlet is the Multinomial generalization of a beta
- In the Beta-Binomial we can write things as if  $m = \alpha + \beta$  is the number of **pseudo observations** and  $E[p] = \frac{\alpha}{\alpha + \beta}$ .
- In the Dirichlet-Multinomial we can write things as if  $m = \sum_k \alpha_k$  and 
$$E[p_1, \dots, p_k] = \left( \frac{\alpha_1}{\sum_{k=1}^K \alpha_k} \cdots \frac{\alpha_K}{\sum_{k=1}^K \alpha_k} \right)$$
- In both cases it is as if we see  $m$  observations from before the data and  $n$  observations from the data.

Prior (Gamma) has  $\alpha$  occurrences in  $\beta$  intervals

$$Pr(\theta|\alpha, \beta) = \frac{\theta^{\alpha-1} e^{-\theta/\beta}}{\beta^\alpha \Gamma(\alpha)} \text{ for } \theta > 0, \alpha > 0, \beta > 0$$

Likelihood (Poisson) we observe  $k$  events in each of the  $n$  periods:

$$Pr(x_i = k|\theta) = \frac{\theta^k e^{-\theta}}{k!}$$

Posterior is **Gamma**

$$\theta \sim \text{Gamma}(\alpha + \sum_{i=1}^n x_i, \beta + n)$$

# MULTIVARIATE NORMAL DISTRIBUTION

- Multivariate Normal Likelihood
- Various Prior Options: Inverse Wishart for Variance.
- This is useful later on but a mess of matrix algebra for now.
- I will skip it.



# EMPIRICAL BAYES

# WHAT IS EMPIRICAL BAYES?

- Priors can be an important modeling choice
- But what makes a good prior?
  - Sufficiently diffuse
  - As non-informative as possible
  - Don't tip the scales
  - Don't rule out the truth
- Idea: can we use the data itself to construct a prior?
  - If everything is a function of data, are we back in frequentist paradigm?
  - Can we get benefits of Bayes estimation without unpalatable assumptions?

# A (FAMOUS) BASEBALL EXAMPLE

Suppose we want to estimate batting averages (AVG) for some baseball players

- $AVG = \frac{\#hits}{\#AtBats}$
- Use data on the first  $n = 45$  at bats and hits  $x_i$  for the 1970 season.
- Predict the batting average  $\mu_i$  for the end of the season ( $n = 400 - 500$  at bats).
- Obvious estimate is batting average after 45 at bats:  $\hat{\mu}_i^{MLE} = x_i/45$ .
- Is there a better estimate?

# A BASEBALL EXAMPLE

**Table 1.1:** Batting averages  $z_i = \hat{\mu}_i^{(\text{MLE})}$  for 18 major league players early in the 1970 season;  $\mu_i$  values are averages over the remainder of the season. The James–Stein estimates  $\hat{\mu}_i^{(\text{JS})}$  (1.35) based on the  $z_i$  values provide much more accurate overall predictions for the  $\mu_i$  values. (By coincidence,  $\hat{\mu}_i$  and  $\mu_i$  both average 0.265; the average of  $\hat{\mu}_i^{(\text{JS})}$  must equal that of  $\hat{\mu}_i^{(\text{MLE})}$ .)

Name	hits/AB	$\hat{\mu}_i^{(\text{MLE})}$	$\mu_i$	$\hat{\mu}_i^{(\text{JS})}$
Clemente	18/45	.400	<b>.346</b>	.294
F Robinson	17/45	.378	<b>.298</b>	.289
F Howard	16/45	.356	<b>.276</b>	.285
Johnstone	15/45	.333	<b>.222</b>	.280
Berry	14/45	.311	<b>.273</b>	.275
Spencer	14/45	.311	<b>.270</b>	.275
Kessinger	13/45	.289	<b>.263</b>	.270
L Alvarado	12/45	.267	<b>.210</b>	.266
Santo	11/45	.244	<b>.269</b>	.261
Swoboda	11/45	.244	<b>.230</b>	.261
Unser	10/45	.222	<b>.264</b>	.256
Williams	10/45	.222	<b>.256</b>	.256
Scott	10/45	.222	<b>.303</b>	.256
Petrocelli	10/45	.222	<b>.264</b>	.256
E Rodriguez	10/45	.222	<b>.226</b>	.256
Campaneris	9/45	.200	<b>.286</b>	.252
Munson	8/45	.178	<b>.316</b>	.247
Alvis	7/45	.156	<b>.200</b>	.242
Grand Average		.265	<b>.265</b>	.265

## A (FAMOUS) BASEBALL EXAMPLE

Probably we can do better than the MLE here:

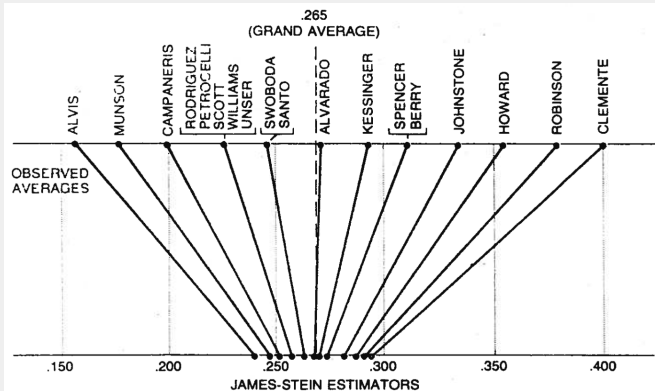
- Thurman Munson wins Rookie of the Year and ends up batting  $\mu_i = .316$ . If he batted .178 all year, his career would not have lasted long.
- Clemente's .400 seems unlikely to hold up. Last player to hit  $> .400$  was Ted Williams .406 in 1941.
- But how?

Idea is to take an average between the observed average  $y_i$  and the overall mean  $\bar{y}$ :

$$\hat{\mu}_i^{JS} = (1 - \lambda) \cdot \bar{y} + \lambda \cdot y_i, \quad \lambda = 1 - \frac{(m - 3)\sigma^2}{\sum_i (y_i - \bar{y})^2}$$

- This has the effect of **shrinking**  $y_i$  towards the **prior mean**  $\bar{y}$ .
- In this case the **prior mean** is just  $\bar{y}$  the grand-mean of all players
- How can information about unrelated players inform us about  $\mu_i$ ?
- Also consider proportion of foreign cars in Chicago as an additional  $y_i$ , can this help too?
- The **shrinkage factor**  $\lambda$  depends on sample size and variance, but how is it chosen?

# A BASEBALL EXAMPLE



**JAMES-STEIN ESTIMATORS** for the 18 baseball players were calculated by “shrinking” the individual batting averages toward the overall “average of the averages.” In this case the grand average is .265 and each of the averages is shrunk about 80 percent of the distance to this value. Thus the theorem on which Stein’s method is based asserts that the true batting abilities are more tightly clustered than the preliminary batting averages would seem to suggest they are.

## **EXAMPLE: DIVERSION RATIOS**



# THE DIVERSION RATIO: CONLON AND MORTIMER (2018)

Raise price of good  $j$ . People leave. What fraction of leavers switch to  $k$ ?

$$D_{jk} = \frac{\frac{\partial q_k}{\partial p_j}}{\left| \frac{\partial q_j}{\partial p_j} \right|}$$

It's one of the best ways economists have to characterize competition among sellers.

- High Diversion: Close Substitutes → Mergers more likely to increase prices.
- Very low diversion → products may not be in the same market.  
(ie: Katz & Shapiro)
- Demand Derivatives NOT elasticities.
- No equilibrium responses.

Remove a product  $j$  measure  $\Delta q_j$  and  $\Delta q_k$ .

$$\overline{D}_{jk} = \frac{\widehat{\Delta q_k}}{|\widehat{\Delta q_j}|} = \frac{E[q_k|Z=1] - E[q_k|Z=0]}{|E[q_j|Z=1] - E[q_j|Z=0]|} = \frac{E[q_k|Z=1] - E[q_k|Z=0]}{E[q_j|Z=0]}$$

## USING A BETA-BINOMIAL PRIOR

How to restrict  $D_{jk} \in [0, 1]$ ?

$$\Delta q_k | \Delta q_j, D_{jk} \sim \text{Bin}(n = \Delta q_j, p = D_{jk})$$

# USING A BETA-BINOMIAL PRIOR

How to restrict  $D_{jk} \in [0, 1]$ ?

$$\Delta q_k | \Delta q_j, D_{jk} \sim \text{Bin}(n = \Delta q_j, p = D_{jk})$$

$$D_{jk} | \beta_1, \beta_2 \sim \text{Beta}(\beta_1, \beta_2)$$

$$E[D_{jk} | \beta_1, \beta_2, \Delta q_j, \Delta q_k] = \frac{\beta_1 + \Delta q_k}{\beta_1 + \beta_2 + \Delta q_j}$$

$$\mu_{jk} = \underbrace{\frac{\beta_1}{\beta_1 + \beta_2}}_{m_{jk}}, \quad \lambda = \frac{m_{jk}}{m_{jk} + \Delta q_j}$$

$$\widehat{D_{jk}} = \lambda \cdot \mu_{jk} + (1 - \lambda) \frac{\widehat{\Delta q_k}}{\widehat{\Delta q_j}}$$

$\mu_{jk}$  is prior mean;  $m_{jk}$  is no. pseudo-obs;  $\lambda$  weights our prior mean.  
When we have a lot of experimental obs, prior receives little weight.

# USING A DIRICHLET PRIOR

How do restrict  $D_j \in \Delta$ ?

- Same idea as before, but use **Dirichlet Prior**.
- Acts like pseudo-observations from the **multinomial** distribution.
- If we had same number of treated observations for each substitute we would have conjugacy/closed form (We don't).
- Likelihood is  $\Delta q_k \sim \text{Binomial}(\Delta q_j, D_{jk})$  not **multinomial**.
- We use  $m = 3.05$  pseudo-observations for the Dirichlet prior.
- Estimator is still technically **non-parametric**. Why?

# SHRINKAGE ESTIMATOR, INTUITION

- We “shrink” towards the prior mean when we have experimental estimates that are imprecise.
- Idea is very simple: when we have lots of data, use the experimental measure.
- When data are scarce: put more weight on the prior/model-based measure.
  - ▶ In practice: FTC/DOJ tend to assume diversion proportional to marketshare
  - ▶ Use plain logit (could also use more complicated model)
  - ▶ Logit sets the mean of the prior as:  $\mu_{jk} = \frac{s_k}{1-s_j}$

# RESULTS: MARS

Firm	Product	# Weeks	$\Delta q_k$	$\Delta q_j$	$\frac{\Delta q_k}{\Delta q_j}$	Beta(j)	Beta(300)	Dirichlet(4.15)
Snickers Removal								
Mars	M&M Peanut	176	375.52	-954.30	39.35	37.04	30.80	18.40
Mars	Twix Caramel	134	289.60	-702.39	41.23	37.86	29.49	15.88
Pepsi	Rold Gold (Con)	174	161.37	-900.11	17.93	16.84	13.95	7.54
Nestle	Butterfinger	61	72.95	-362.82	20.11	17.07	11.19	4.45
Mars	M&M Milk Chocolate	97	71.76	-457.36	15.69	13.83	9.85	4.14
Kraft	Planters (Con)	136	78.01	-759.87	10.27	9.57	7.80	3.81
Kellogg	Zoo Animal Cracker	177	65.72	-970.22	6.77	6.48	5.68	2.92
Pepsi	Sun Chip	159	45.30	-866.09	5.23	4.98	4.33	2.07
Hershey	Choc Hershey (Con)	41	29.78	-179.57	16.58	12.17	6.30	2.01
	Outside Good	180	460.89	-970.22	47.50			23.12
M&M Peanut Removal								
Mars	Snickers	218	296.58	-1239.29	23.93	22.90	19.91	16.47
Mars	Twix Caramel	176	110.93	-1014.32	10.94	10.39	8.88	6.76
Mars	M&M Milk Chocolate	99	73.47	-529.58	13.87	12.46	9.18	6.26
Nestle	Raisinets	181	71.82	-1001.14	7.17	6.82	5.82	4.37
Kraft	Planters (Con)	190	61.42	-1046.10	5.87	5.62	4.90	3.60
Hershey	Twizzlers	62	32.98	-332.99	9.90	8.32	5.32	3.35
Kellogg	Rice Krispies Treats	46	22.37	-220.17	10.16	7.90	4.43	2.51
Pepsi	Frito	160	37.25	-902.42	4.13	3.95	3.47	2.37
	Outside Good	218	606.18	-1238.49	48.95			36.35

# RESULTS: KELLOGG'S

Firm	Product	# Weeks	$\Delta q_k$	$\Delta q_j$	$\frac{\Delta q_k}{\Delta q_j}$	Beta(j)	Beta(300)	Dirichlet(4,15)
Animal Crackers Removal								
Pepsi	Rold Gold (Con)	132	114.39	-440.80	25.95	22.90	16.21	9.89
Mars	Snickers	145	92.44	-483.63	19.11	17.26	13.04	7.58
Mars	M&M Peanut	142	77.72	-469.44	16.55	14.98	11.43	6.47
Kellogg	CC Famous Amos	144	66.18	-478.20	13.84	12.40	9.15	5.39
Pepsi	Baked Chips (Con)	134	62.55	-447.60	13.97	12.46	9.13	5.27
Mars	Twix Caramel	110	50.17	-338.97	14.80	12.75	8.74	4.58
Sherwood	Ruger Wafer (Con)	119	48.20	-368.65	13.07	11.28	7.63	4.28
Hershey	Choc Herhsey (Con)	30	33.60	-132.57	25.34	17.14	7.86	3.81
Kellogg	Rice Krispies Treats	13	23.52	-37.80	62.22	23.24	7.16	2.99
Kar's Nuts	Kar Sweet&Salty Mix	95	30.06	-334.50	8.99	7.72	5.27	2.73
Misc	Popcorn (Con)	56	25.72	-226.89	11.34	8.92	5.08	2.61
Kraft	Planters (Con)	114	28.05	-380.25	7.38	6.53	4.78	2.43
Mars	M&M Plain	73	22.67	-295.07	7.68	6.47	4.26	2.15
	Outside Good	145	240.52	-482.91	49.81			21.98
Famous Amos Removal								
Pepsi	Sun Chip	139	143.60	-355.68	40.37	34.39	22.66	15.75
Kraft	Planters (Con)	121	82.11	-332.61	24.69	20.89	13.68	8.75
Hershey	Choc Hershey (Con)	38	48.60	-66.84	72.72	36.93	13.36	7.18
Pepsi	Frito	119	49.88	-313.21	15.93	13.44	8.85	5.32
Misc	Rasbry Knotts	133	46.62	-345.38	13.50	11.45	7.49	4.81
Pepsi	Grandmas Choc Chip	95	39.99	-259.21	15.43	12.51	7.62	4.49
Pepsi	Dorito Buffalo Ranch	72	38.11	-224.24	17.00	13.28	7.53	4.43
Pepsi	Chs PB Frito Cracker	34	26.87	-83.65	32.13	18.16	7.14	3.74
Kellogg	Choc Sandwich FA	57	27.97	-122.04	22.91	15.06	6.84	3.69
Pepsi	Rold Gold (Con)	147	32.62	-392.22	8.32	7.40	5.54	3.19
Kraft	Oreo Thin Crisps	29	20.73	-43.29	47.89	19.20	6.12	3.05
Mars	Combos (Con)	98	23.56	-274.54	8.58	7.03	4.34	2.61



# COMPARISON OF ASSUMPTIONS: SNICKERS EXPERIMENT

	Total	Assn 1	Assn 2	Assn 3 ( $m = K$ )	Assn 4 ( $m = 4.15$ )
Products with $D_{jk} < 0$	51	24	26	0	0
Products with $0 \leq D_{jk} \leq 10$	51	13	15	43	48
Products with $10 \leq D_{jk} \leq 20$	51	5	5	5	2
Products with $D_{jk} > 20$	51	9	5	3	1
Sum of all positive $D_{jk}$ s	51	402.84	301.95	265.41	98.72
Sum of all negative $D_{jk}$ s	51	-238.90	-239.07	0.00	0.00

Note: Table includes only products for which there were at least 50 sales of the focal product in control weeks, on average.

# STAN CODE

```
% Main Specification: Dirichlet Prior
data {
  int<lower=1> J;           // number of products, including outside good
  int<lower=1> N[J];        // number of trials
  int<lower=0> y[J];        // number of successes for each product j
  vector[J] priors;        // mean of the distribution of alpha
}

parameters {
  simplex[J] theta;
}

model {
  theta ~ dirichlet(priors);
  for (j in 1:J) {
    y[j] ~ binomial(N[j], theta[j]);
  }
}
```

# BAYESIAN ESTIMATION

# HOW DO WE ESTIMATE THESE MODELS?

A few options

- With conjugate priors we can closed forms
- Can do it by hand if we have **sufficient statistics**
  - Clear in beta-binomial or poisson-gamma relationship.
  - Mostly not the case.
- Mostly we do what is known as **Markov Chain Monte Carlo**
- The goal is to draw from  $f(\theta|\mathbf{X})$  or to compute moments of the distribution  $E_f[g(\theta)]$ .

- Suppose we want to calculate a function

$$E_f[g(x)] = \int g(x)f(x)dx$$

- How do we do it?

1. Draw from  $\hat{x}_s \sim f(x)$
2. Calculate  $g(\hat{x}_s)$
3. Repeat for  $s = 1 \dots, S$
4. Calculate  $E[\hat{g}] = \frac{1}{S} \sum_{s=1}^S g(\hat{x}_s)$

# MONTE CARLO EXAMPLE

- Let's integrate  $g(x) = \phi(x)$  (normal pdf) over  $f(x) = \text{Unif}(0, 1)$  from  $[0, 1]$ .
- We know the answer is  $\Phi(1) - \Phi(0)$ .

```
Integral <- function(n){  
  X <- runif(n)  
  Y <- exp(-X^2/2)/sqrt(2*pi)  
  Int <- sum(Y)/n  
  Error <- Int-(pnorm(1)-pnorm(0))  
  list(Int, Error)}
```

# GIBBS SAMPLING

The first building block is known as **Gibbs Sampling**

- Suppose that  $p(x, y)$  is a p.d.f or p.m.f that is hard to sample directly from.
- But suppose that  $p(x|y)$  or  $p(y|x)$  are easy to sample from.
- Gibbs sampler says:
  1. Initialize  $(x_0, y_0)$ .
  2. Randomly draw  $y_1 \sim g(y|x_0)$ .
  3. Randomly draw  $x_1 \sim f(x|y_1)$ .
  4. Randomly draw  $y_2 \sim g(y|x_1)$ .
  5. Rinse and Repeat.
- This sequence  $(x_0, y_0), (x_1, y_1), (x_2, y_2), \dots$  is a **Markov Chain**.
- Why? Because  $(x_k, y_k) | (x_{k-1}, y_{k-1})$  doesn't depend on  $x_{k-h}$  or  $y_{k-h}$  for  $h \geq 2$ .
  - ▶ This does not mean that  $(x_3, y_3)$  and  $(x_1, y_1)$  are I.I.D!

**THANKS!**