

MULTI-PARAMETER BAYESIAN MODELS



Learning Objectives

- Review properties of the multivariate normal distribution
- Review conjugate priors for multivariate normal data
- Explore Wishart and Inverse Wishart distributions
- Work with multinomial and Dirichlet distributions
- Apply Bayesian inference to real data examples

Definition

A k -dimensional vector \mathbf{y} follows a multivariate normal with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$ if:

$$p(\mathbf{y}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \propto |\boldsymbol{\Sigma}|^{-1/2} \exp\left(-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{y} - \boldsymbol{\mu})\right)$$

- $\boldsymbol{\mu} \in \mathbb{R}^k$ and $\boldsymbol{\Sigma}$ is a symmetric and positive definite matrix. If $\boldsymbol{\Sigma}$ is not of full rank, the distribution is degenerate and does not have a density.
- A positive definite matrix $\boldsymbol{\Sigma}$ satisfies:
 1. $\mathbf{x}^T \boldsymbol{\Sigma} \mathbf{x} > 0$ for all nonzero vector \mathbf{x} . $\mathbf{x}^T \boldsymbol{\Sigma} \mathbf{x} = 0$ if and only if $\mathbf{x} = \mathbf{0}$.
 2. Eigenvalues of $\boldsymbol{\Sigma}$ are all positive.
 3. The determinant of $\boldsymbol{\Sigma}$ is positive (since it is the product of all eigenvalues)

- Any subset of \mathbf{y} has a normal distribution.
- Any linear combination of \mathbf{y} is also normal.

$$\mathbf{A}\mathbf{y} \sim N(\mathbf{A}\boldsymbol{\mu}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)$$

- The conditional distributions are also normal. WLOG, let $\mathbf{y} = (\mathbf{y}_1, \mathbf{y}_2)$,

$$p(\mathbf{y}_1|\mathbf{y}_2) \sim N(\mathbf{m}, \mathbf{W})$$

where

$$\mathbf{m} = \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{y}_2 - \boldsymbol{\mu}_2)$$

$$\mathbf{W} = \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}$$

For n observations $\mathbf{y} = (\mathbf{y}_1, \dots, \mathbf{y}_n)^T$, notice we can rewrite the likelihood

$$p(\mathbf{y}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \propto |\boldsymbol{\Sigma}|^{-n/2} \exp\left(-\frac{1}{2} \sum_i (\mathbf{y}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{y}_i - \boldsymbol{\mu})\right)$$

into

$$p(\mathbf{y}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \propto |\boldsymbol{\Sigma}|^{-n/2} \exp\left(-\frac{1}{2} \text{tr}(\boldsymbol{\Sigma}^{-1} \tilde{\mathbf{S}})\right)$$

with $\tilde{\mathbf{S}} = \sum_i^n (\mathbf{y}_i - \boldsymbol{\mu})(\mathbf{y}_i - \boldsymbol{\mu})^T$.

Note that $\tilde{\mathbf{S}} = \sum_{i=1}^n (\mathbf{y}_i - \bar{\mathbf{y}})(\mathbf{y}_i - \bar{\mathbf{y}})^T + n(\bar{\mathbf{y}} - \boldsymbol{\mu})(\bar{\mathbf{y}} - \boldsymbol{\mu})^T$. That is, the sample mean and variance are sufficient statistics for $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$.

Multivariate normal with known Σ

- Condition: $\mathbf{y}|\boldsymbol{\mu}, \boldsymbol{\Sigma} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ with known $\boldsymbol{\Sigma}$
- Conjugate prior: $\boldsymbol{\mu} \sim N(\boldsymbol{\mu}_0, \boldsymbol{\Lambda}_0)$
- Posterior: $\boldsymbol{\mu}|\mathbf{y}, \boldsymbol{\Sigma} \sim N(\boldsymbol{\mu}_n, \boldsymbol{\Lambda}_n)$ where:

$$\boldsymbol{\mu}_n = \boldsymbol{\Lambda}_n(\boldsymbol{\Lambda}_0^{-1}\boldsymbol{\mu}_0 + n\boldsymbol{\Sigma}^{-1}\bar{\mathbf{y}}), \quad \boldsymbol{\Lambda}_n^{-1} = \boldsymbol{\Lambda}_0^{-1} + n\boldsymbol{\Sigma}^{-1}$$

- Notice for the predictive distribution of a new data point $\tilde{\mathbf{y}}$,

$$p(\tilde{\mathbf{y}}, \boldsymbol{\mu}|\mathbf{y}) = N(\tilde{\mathbf{y}}|\boldsymbol{\mu}, \boldsymbol{\Sigma})N(\boldsymbol{\mu}|\boldsymbol{\mu}_n, \boldsymbol{\Lambda}_n)$$

So $(\tilde{\mathbf{y}}, \boldsymbol{\mu})$ has a joint normal posterior distribution. That is $\tilde{\mathbf{y}}|\mathbf{y}$ is also normally distributed.

- Posterior predictive for new data point $\tilde{\mathbf{y}}$:

$$\mathbb{E}(\tilde{\mathbf{y}}|\mathbf{y}) = \boldsymbol{\mu}_n, \quad \text{var}(\tilde{\mathbf{y}}|\mathbf{y}) = \boldsymbol{\Sigma} + \boldsymbol{\Lambda}_n$$

Definition

For a $k \times k$ positive definite matrix Σ , the Inverse Wishart distribution with $\nu > k - 1$ degrees of freedom and scale matrix Λ (positive definite) has density:

$$p(\Sigma|\nu, \Lambda) \propto |\Sigma|^{-\frac{\nu+k+1}{2}} \exp\left(-\frac{1}{2} \text{tr}(\Lambda\Sigma^{-1})\right)$$

An important property is that the prior $\Sigma \sim \text{Inv-Wishart}(\nu, \Lambda)$ is equivalent to a Wishart prior on the precision matrix $\Omega = \Sigma^{-1}$:

$$\Omega \sim \text{Wishart}(\nu, \Lambda^{-1})$$

- The Wishart distribution is defined through the construction:
If $z_1, \dots, z_n \sim N(0, \Lambda^{-1})$ independently, then $\sum_i z_i z_i^T \sim \text{Wishart}(n, \Lambda^{-1})$
- An Inverse Wishart prior on Σ with known μ leads to:

$$\Sigma | \mathbf{y} \sim \text{Inv-Wishart}(n + \nu, \Lambda + \tilde{S})$$

where the scale matrix is updated by the sample sum of squares

Wishart and Inverse Wishart: Properties I

- The prior mean of Σ under $\text{Inv-Wishart}(\nu, \Lambda)$ is $\frac{1}{\nu-k-1} \Lambda$
- To understand the partition structure, let $\mathbf{y} = (\mathbf{y}_1, \mathbf{y}_2)$ and

$$\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$$

- Define the **Schur complement**: $\Sigma_{22,1} = \Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12}$
- This satisfies: $|\Sigma| = |\Sigma_{11}| |\Sigma_{22,1}|$

Wishart and Inverse Wishart: Properties II

- The Inverse Wishart distribution has a special partition property:

$$\Sigma_{11} \sim \text{Inv-Wishart}(\nu - k_2, \Lambda_{11})$$

$$\Sigma_{22,1} \sim \text{Inv-Wishart}(\nu, \Lambda_{22,1})$$

and Σ_{11} and $\Sigma_{22,1}$ are independent

- For $k_1 = 1$, the Inverse Wishart prior implies Inverse Gamma on marginal variances:

$$\sigma_1^2 \sim \text{Inv-Gamma}\left(\frac{\nu - k + 1}{2}, \frac{\Lambda_{11}}{2}\right)$$

- Note that the inverse Wishart has only a single degrees of freedom parameter governing all dimensions.

Normal-Inverse-Wishart Prior

Denote $(\mu, \Sigma) \sim \text{NIW}(\mu_0, \kappa_0, \nu_0, \Lambda_0)$ if:

$$\Sigma \sim \text{Inv-Wishart}(\nu_0, \Lambda_0), \quad \mu | \Sigma \sim N(\mu_0, \Sigma / \kappa_0)$$

The posterior is also NIW with parameters:

$$\mu_n = \frac{\kappa_0 \mu_0 + n \bar{\mathbf{y}}}{\kappa_0 + n}$$

$$\kappa_n = \kappa_0 + n$$

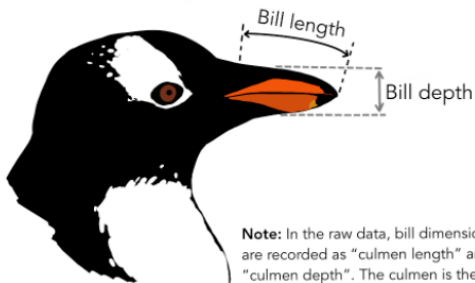
$$\nu_n = \nu_0 + n$$

$$\Lambda_n = \Lambda_0 + \mathbf{S} + \frac{\kappa_0 n}{\kappa_0 + n} (\bar{\mathbf{y}} - \mu_0)(\bar{\mathbf{y}} - \mu_0)^T$$

The posterior predictive distribution of $\tilde{\mathbf{y}}$ is multivariate t .

Example: Palmer penguins

- As an illustration, let us look at the Palmer penguins dataset. It collects various measurements for 344 penguins in Palmer station in Antarctica.
- Let us build a linear model for the bill length.



Example: Palmer penguins

We will use only one of the three penguin species in this example.

```
library(palmerpenguins)
data <- subset(penguins, !is.na(flipper_length_mm) &
               !is.na(bill_length_mm) &
               !is.na(species) &
               !is.na(sex) &
               species == "Gentoo")
head(penguins)
```

	species	island	bill_length_mm	bill_depth_mm	flipper_length_mm	body_mass_g	sex	year
1	Gentoo	Biscoe	46.10	13.20	211	4500	female	2007
2	Gentoo	Biscoe	50.00	16.30	230	5700	male	2007
3	Gentoo	Biscoe	48.70	14.10	210	4450	female	2007
4	Gentoo	Biscoe	50.00	15.20	218	5700	male	2007
5	Gentoo	Biscoe	47.60	14.50	215	5400	male	2007
6	Gentoo	Biscoe	46.50	13.50	210	4550	female	2007

Example: Palmer penguins

Consider estimating the joint distribution of flipper length and bill length using a two-dimensional Gaussian model.

```
# 2-d data matrix y
y <- as.matrix(data[, c("flipper_length_mm", "bill_length_mm")])
n <- dim(y)[1]
K <- dim(y)[2]

# Priors
mu0 <- rep(0, K)
nu0 = K + 1
Lambda0 = diag(K)
kappa0 = 1

# Posterior distribution parameters
y_bar <- apply(y, 2, mean)
Lambda_n <- Lambda0 +
  (t(y)-y_bar)%*%t(t(y)-y_bar) +
  kappa0*n/(kappa0+n)*(mu0-y_bar)%*%t(mu0-y_bar)

# to avoid numerical issue with matrix not symmetric
Lambda_n <- (Lambda_n+t(Lambda_n))/2
```

Example: Palmer penguins

Draw samples from the prior predictive distribution

$$\Sigma \sim \text{Inv-Wishart}(\nu_0, \Lambda_0)$$

$$\mu | \Sigma \sim N(\mu_0, \Sigma / \kappa_0)$$

$$\mathbf{y} | \mu, \Sigma \sim N(\mu, \Sigma)$$

```
library(LaplacesDemon)

pred_prior <- matrix(0, 1e4, K)

for(i in 1:1e4){
  Sigma_draw <- rinvwishart(nu0, Lambda0)
  mu_draw <- t(rmvn(1, mu0, 1/(kappa0)*Sigma_draw))
  pred_prior[i, ] <- rmvn(1, t(mu_draw), Sigma_draw)
}
```

Example: Palmer penguins

Draw samples from the posterior predictive distribution

$$\Sigma | \mathbf{y} \sim \text{Inv-Wishart}(\nu_n, \Lambda_n)$$

$$\mu | \Sigma, \mathbf{y} \sim N(\mu_n, \Sigma / \kappa_n)$$

$$\tilde{\mathbf{y}} | \mu, \Sigma, \mathbf{y} \sim N(\mu, \Sigma)$$

```
pred_post <- matrix(0, 1e4, K)

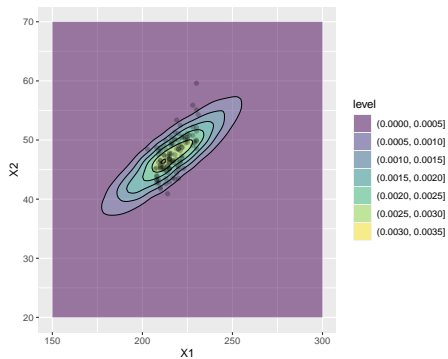
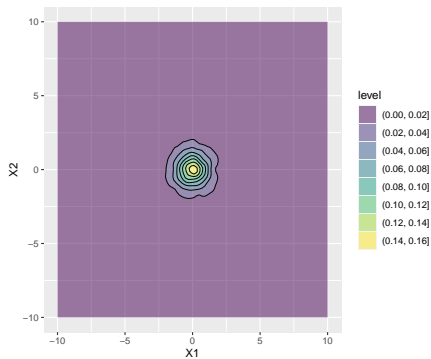
for(i in 1:1e4){
  Sigma_draw <- rinvwishart(nu0+n, Lambda_n)
  mu_draw <- t(rmvn(1, (kappa0*mu0+n*y_bar)/(kappa0+n),
                  1/(kappa0+n)*Sigma_draw))
  pred_post[i, ] <- rmvn(1, t(mu_draw), Sigma_draw)
}
```

Example: Palmer penguins

Plot the prior and posterior predictive distribution of y .

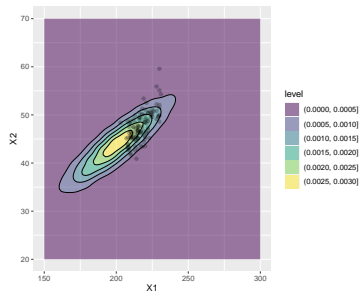
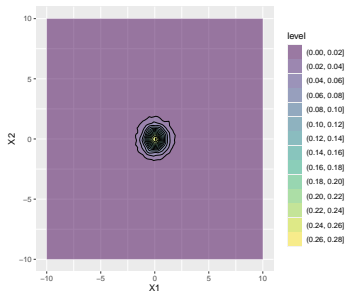
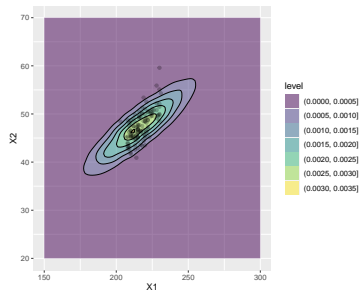
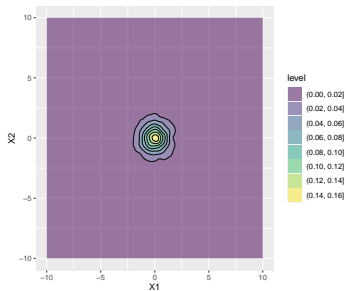
```
# Draw density plot
g1 <- ggplot(data = data.frame(pred_prior), aes(x = X1, y = X2)) +
  geom_density_2d_filled(alpha = 0.5)+
  geom_density_2d(linewidth = 0.25, colour = "black") +
  xlim(c(-10, 10)) + ylim(c(-10, 10))
g2 <- ggplot(data = data.frame(pred_post), aes(x = X1, y = X2)) +
  geom_density_2d_filled(alpha = 0.5)+
  geom_density_2d(linewidth = 0.25, colour = "black") +
  xlim(c(150, 300)) + ylim(c(20, 70)) +
  geom_point(data = data.frame(y),
            aes(x = flipper_length_mm, y = bill_length_mm),
            alpha = 0.25)
```

Example: Palmer penguins



Example: Palmer penguins

Comparing with $\kappa_0 = 10$, where the prior for μ is more concentrated at 0.



Noninformative priors

- Setting $\Sigma \sim \text{Inv-Wishart}(k + 1, \mathbf{I})$ is a commonly used noninformative prior. It has a nice property that the marginal distribution of each off-diagonal elements in the correlation matrix is uniform (see Barnard, McCulloch, Meng, 2000).
- Another noninformative prior is the Jefferys prior,

$$p(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \propto |\boldsymbol{\Sigma}|^{-(k+1)/2}$$

which is the limit of NIW as $\kappa_0 \rightarrow 0$, $\nu_0 \rightarrow -1$ and $|\Lambda_0| \rightarrow 0$

- Note again, that in many situations (especially later in hierarchical models), it makes more sense to use a weakly informative prior rather than seeking the noninformativeness.

Alternatives to Inverse Wishart

- A major issue with Inverse Wishart prior for the covariance matrix is that there is only one df parameter for all dimensions. This can be seen from the marginal prior distribution of the variance components.
- Another issue is that the Inverse Wishart imposes a prior dependence between correlation and marginal variances.
- When we start to model covariance matrix in hierarchical models, we could also run into issues where the posterior distribution of Σ is heavily concentrated on a degenerate matrix (i.e., the boundary), as the density increases as $|\Sigma| \rightarrow 0$.
- Many alternatives have been proposed based on a decomposition

$$\Sigma = \text{diag}(\mathbf{D})\mathbf{R}\text{diag}(\mathbf{D})$$

including the scaled inverse-Wishart, LKJ prior for the correlation matrix, etc.

Generating normal and inverse Wishart random variables

- To generate $x \sim N(\mu, \Sigma)$, we can start with the Cholesky decomposition of $\Sigma = LL^T$ where L is a lower triangular matrix.
- Then we generate $z \sim N(0, I)$ and compute $x = Lz + \mu$.
- To generate $\Sigma \sim \text{Inv-Wishart}(\nu, \Lambda)$, first generate $\Sigma^{-1} = \Omega \sim \text{Wishart}(\nu, \Lambda^{-1})$ and invert the matrix.
- The Wishart distribution can be sampled by generating ν independent samples $\alpha_1, \dots, \alpha_\nu \sim N(0, \Lambda^{-1})$ and let $\Omega = \sum_{i=1}^{\nu} \alpha_i \alpha_i^T$.
- A faster alternative is to use Bartlett's decomposition:
 1. Generate a lower triangular matrix A with $a_{ii} \sim \sqrt{\chi_{\nu-i+1}^2}$ and $a_{ij} \sim N(0, 1)$ for $j < i$.
 2. Compute the Cholesky decomposition $\Lambda = LL^T$.
 3. Compute $\Omega = LAA^T L^T$.

This method requires only $k(k+1)/2$ random variable generations and automatically produces the Cholesky decomposition for Ω .

- There are R packages with highly optimized implementations, e.g., mvtnorm, MCMCpack, LaplacesDemon, ...

Definition

Consider a categorical random variable with K possible outcomes. Count occurrences of each type for n trials. The count vector \mathbf{y} has density:

$$p(\mathbf{y}|\boldsymbol{\theta}) \propto \prod_{j=1}^K \theta_j^{y_j}$$

with $\sum_{j=1}^K \theta_j = 1$.

- Generalization of the binomial distribution to K categories
- Conjugate prior: the Dirichlet distribution

Definition

A distribution on the K -dimensional simplex:

$$p(\boldsymbol{\theta}|\boldsymbol{\alpha}) \propto \prod_{j=1}^K \theta_j^{\alpha_j-1}$$

with $\sum_{j=1}^K \theta_j = 1$ and $\alpha_j > 0$ for all j .

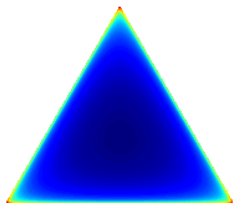
The Dirichlet is the conjugate prior for the multinomial. Under a multinomial likelihood, the posterior is:

$$p(\boldsymbol{\theta}|\mathbf{y}) \propto \prod_{j=1}^K \theta_j^{y_j+\alpha_j-1}$$

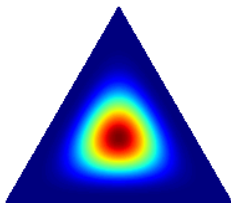
The posterior is also Dirichlet with updated parameters.

The Dirichlet distribution

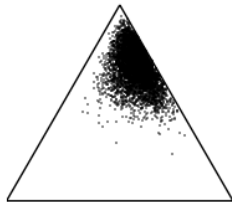
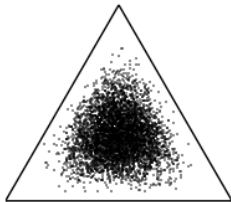
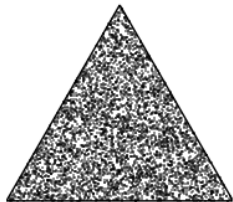
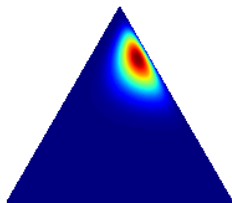
$\alpha = (0.999, 0.999, 0.999)$



$\alpha = (5.000, 5.000, 5.000)$



$\alpha = (2.000, 5.000, 15.000)$



- There is an important relationship between Dirichlet and gamma distributions. Let

$$Z_j \sim \text{Gamma}(\alpha_j, \beta)$$

for $j = 1, \dots, K$ independently, then the vector θ with

$$\theta_j = \frac{Z_j}{\sum_j Z_j}, j = 1, \dots, K$$

follow a Dirichlet distribution with parameter α .

- The conditional distribution of a sub-vector of θ given the remaining elements of θ is also a Dirichlet (after rescaling to sum to 1).

Slovenia opinion poll example

- In 1990, a plebiscite was held in Slovenia at which the adult citizens voted on the question of independence.
- A Slovenian public opinion survey had been conducted that included several questions concerning likely plebiscite attendance and voting. In that survey, 2074 Slovenians were asked:
 1. Are you in favor of independence?
 2. Will you attend the plebiscite?
- The rules of the plebiscite were such that only those attending and voting 'yes' would be counted as being in favor of independence.

<i>Table of counts</i>			
<i>Independence</i>			
<i>Attendance</i>	<i>Yes</i>	<i>No</i>	<i>DK</i>
Yes	1,439	78	159
No	16	16	32
DK	144	54	136

Dirichlet model for contingency tables

$n_{11} = 1439$	$n_{10} = 78$	$n_{1\star} = 159$
$n_{01} = 16$	$n_{00} = 16$	$n_{0\star} = 32$
$n_{\star 1} = 144$	$n_{\star 0} = 54$	$n_{\star\star} = 136$

- If we ignore any counts of DK's, a naive estimate of the fraction of people attending voting yes is $n_{11}/(n_{11} + n_{10} + n_{01} + n_{00}) = 0.93$.
- We can try to understand how the DK counts should be distributed into one of the four cells. Let the unobserved true counts $(x_{11}, x_{10}, x_{01}, x_{00}) \sim \text{Mult}(n, \theta)$.
- Given $n_{\star\star}$ people who we know nothing about, it may be reasonable to assume that

$$(n_{\star\star}^{(11)}, n_{\star\star}^{(10)}, n_{\star\star}^{(01)}, n_{\star\star}^{(00)}) \sim \text{Mult}(n_{\star\star}, \theta)$$

where $n_{\star\star}^{(11)}$ is the number of these people who will attend and vote yes, etc.

- Similarly for the other cells, we can use similar models given the partial information, e.g.,

$$(n_{1\star}^{(11)}, n_{1\star}^{(10)}) \sim \text{Mult}(n_{1\star}, (\frac{\theta_{11}}{\theta_{11} + \theta_{10}}, \frac{\theta_{10}}{\theta_{11} + \theta_{10}}))$$

Dirichlet model for contingency tables

- Thus we can complete a Bayesian model with prior on θ

$$\theta \sim \text{Dir}(\alpha)$$

$$(x_{11}, x_{10}, x_{01}, x_{00}) \sim \text{Mult}(n, \theta)$$

where $x_{ij} = n_{ij} + n_{i\star}^{(ij)} + n_{\star j}^{(ij)} + n_{\star\star}^{(ij)}$ and the corresponding latent variables follow

$$(n_{\star\star}^{(11)}, n_{\star\star}^{(10)}, n_{\star\star}^{(01)}, n_{\star\star}^{(00)}) \sim \text{Mult}(n_{\star\star}, \theta)$$

$$(n_{1\star}^{(11)}, n_{1\star}^{(10)}) \sim \text{Mult}(n_{1\star}, (\frac{\theta_{11}}{\theta_{11} + \theta_{10}}, \frac{\theta_{10}}{\theta_{11} + \theta_{10}}))$$

$$(n_{0\star}^{(01)}, n_{0\star}^{(00)}) \sim \text{Mult}(n_{0\star}, (\frac{\theta_{01}}{\theta_{01} + \theta_{00}}, \frac{\theta_{00}}{\theta_{01} + \theta_{00}}))$$

$$(n_{\star 1}^{(11)}, n_{\star 1}^{(01)}) \sim \text{Mult}(n_{\star 1}, (\frac{\theta_{11}}{\theta_{11} + \theta_{01}}, \frac{\theta_{01}}{\theta_{11} + \theta_{01}}))$$

$$(n_{\star 0}^{(10)}, n_{\star 0}^{(00)}) \sim \text{Mult}(n_{\star 0}, (\frac{\theta_{10}}{\theta_{10} + \theta_{00}}, \frac{\theta_{00}}{\theta_{10} + \theta_{00}}))$$

- The model for latent counts here is implicitly assuming data are missing at random. We will define this assumption more explicitly later in the course.

Dirichlet model for contingency tables

- Red line shows the posterior mean. Blue line shows the naive estimates ignoring the DK's.

