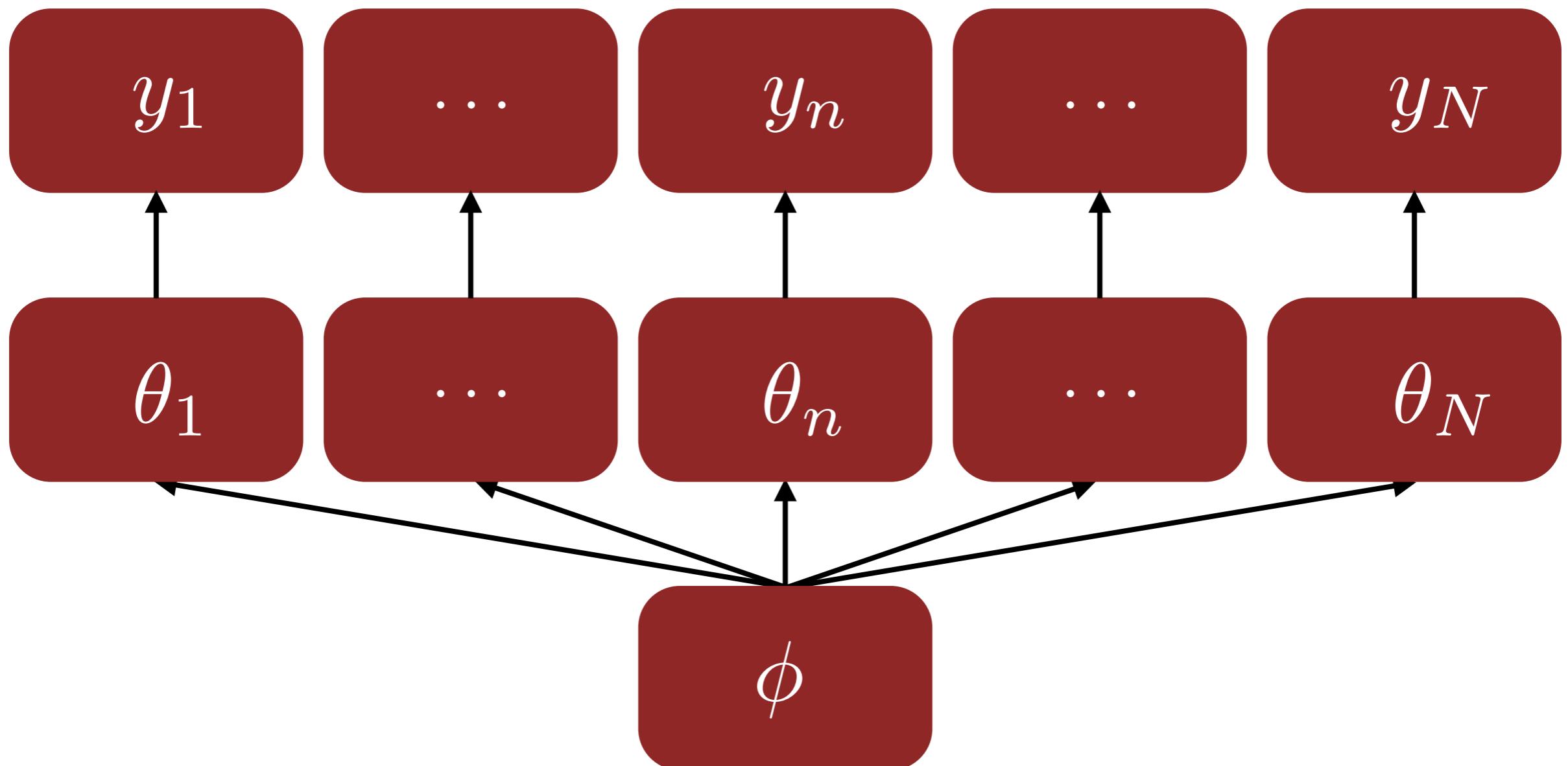


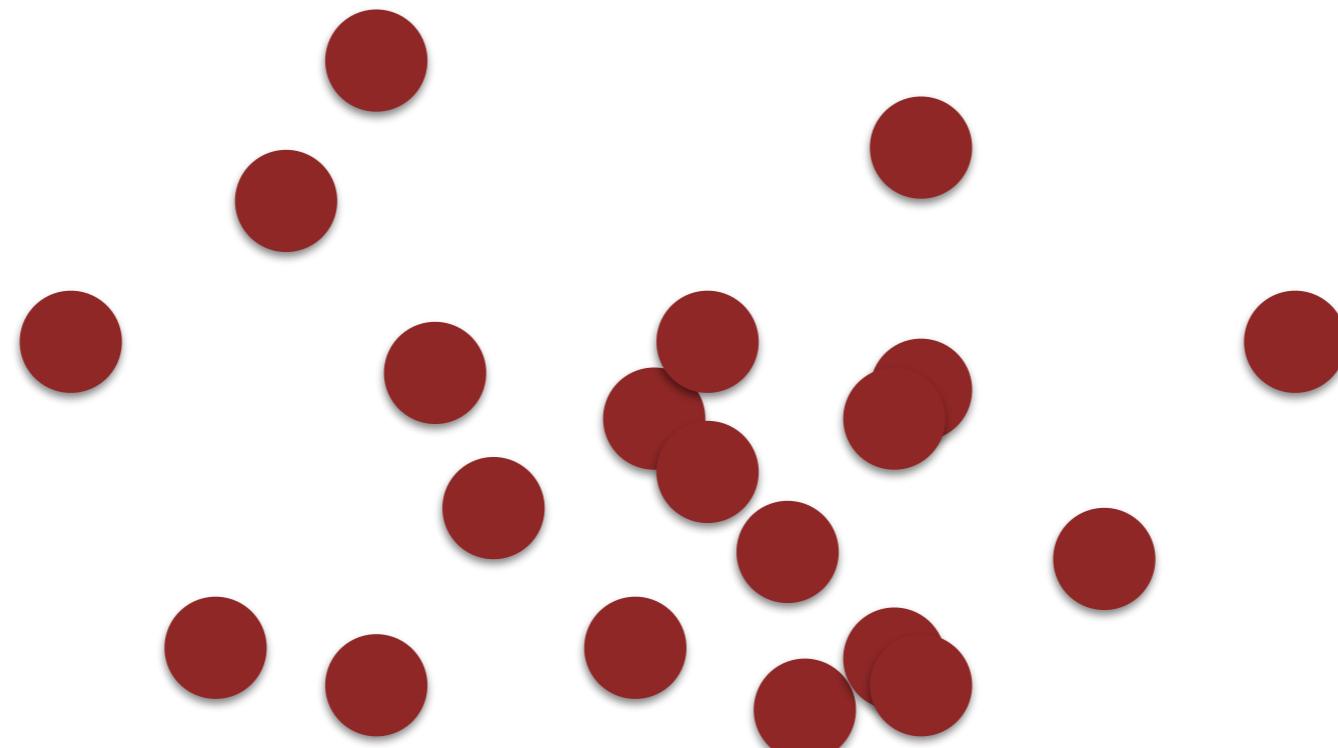
# Hierarchical Models



A common problem in applied statistics  
is modeling individuals of a *population*.

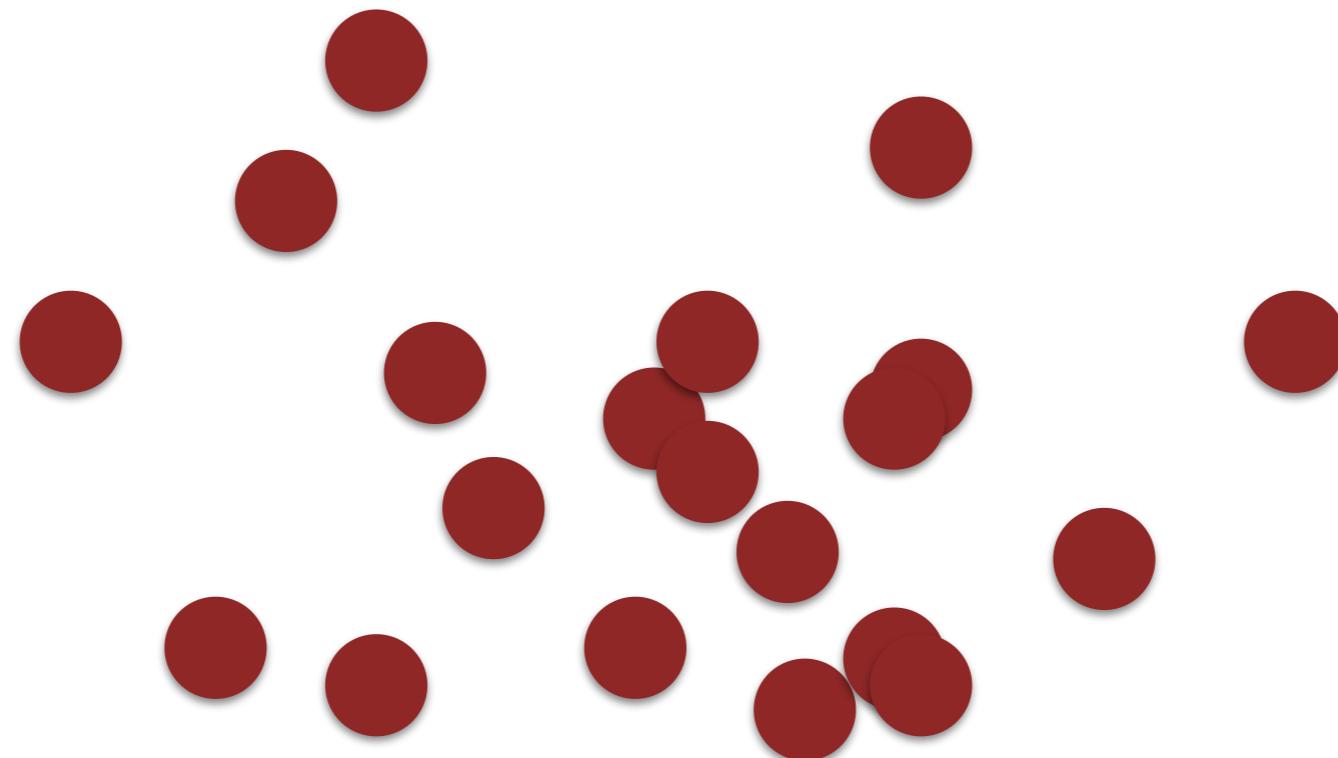
$$y \sim \mathcal{N}(\mathbf{X}^T \boldsymbol{\beta} + \alpha, \sigma)$$

A common problem in applied statistics  
is modeling individuals of a *population*.



$$y \sim \mathcal{N}(\mathbf{X}^T \boldsymbol{\beta} + \alpha, \sigma)$$

If we assume that every individual is equivalent then we can pool the data, but only at the expense of bias.

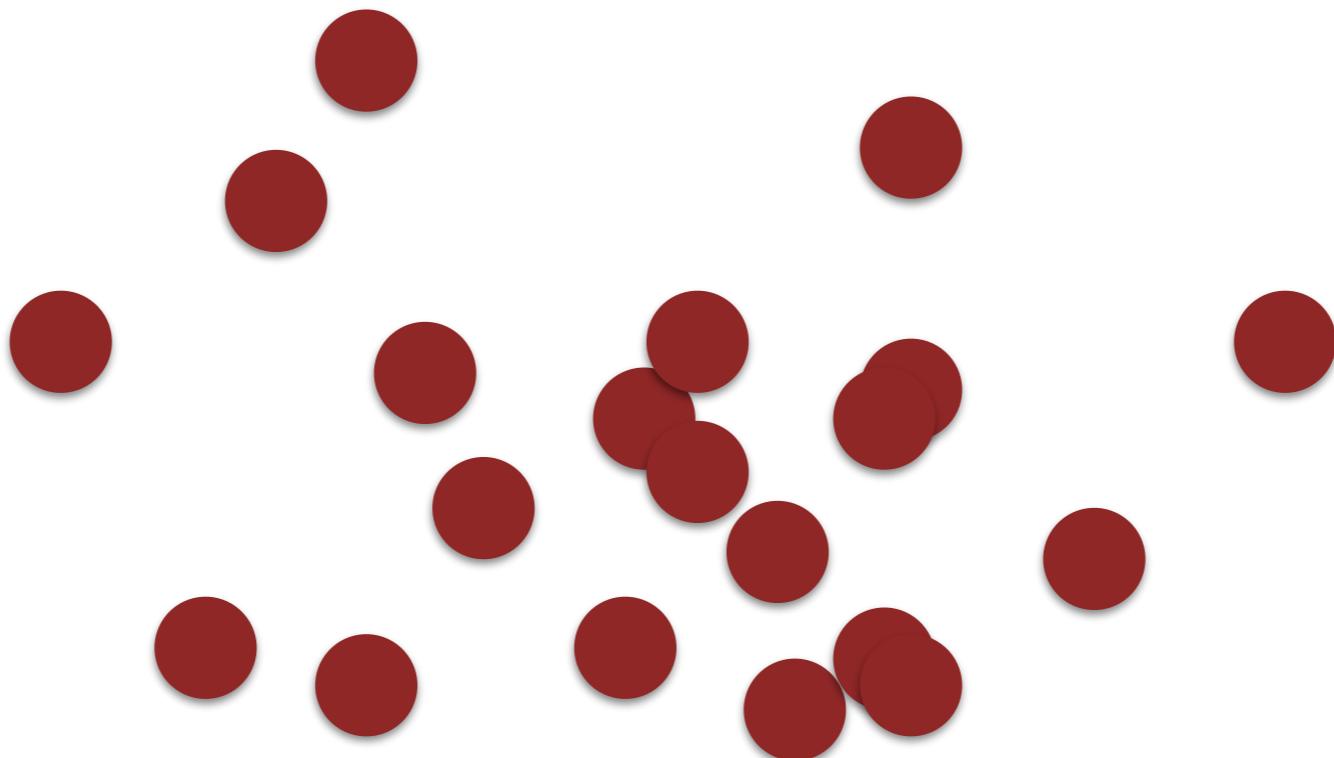


If we assume that every individual is equivalent then we can pool the data, but only at the expense of bias.

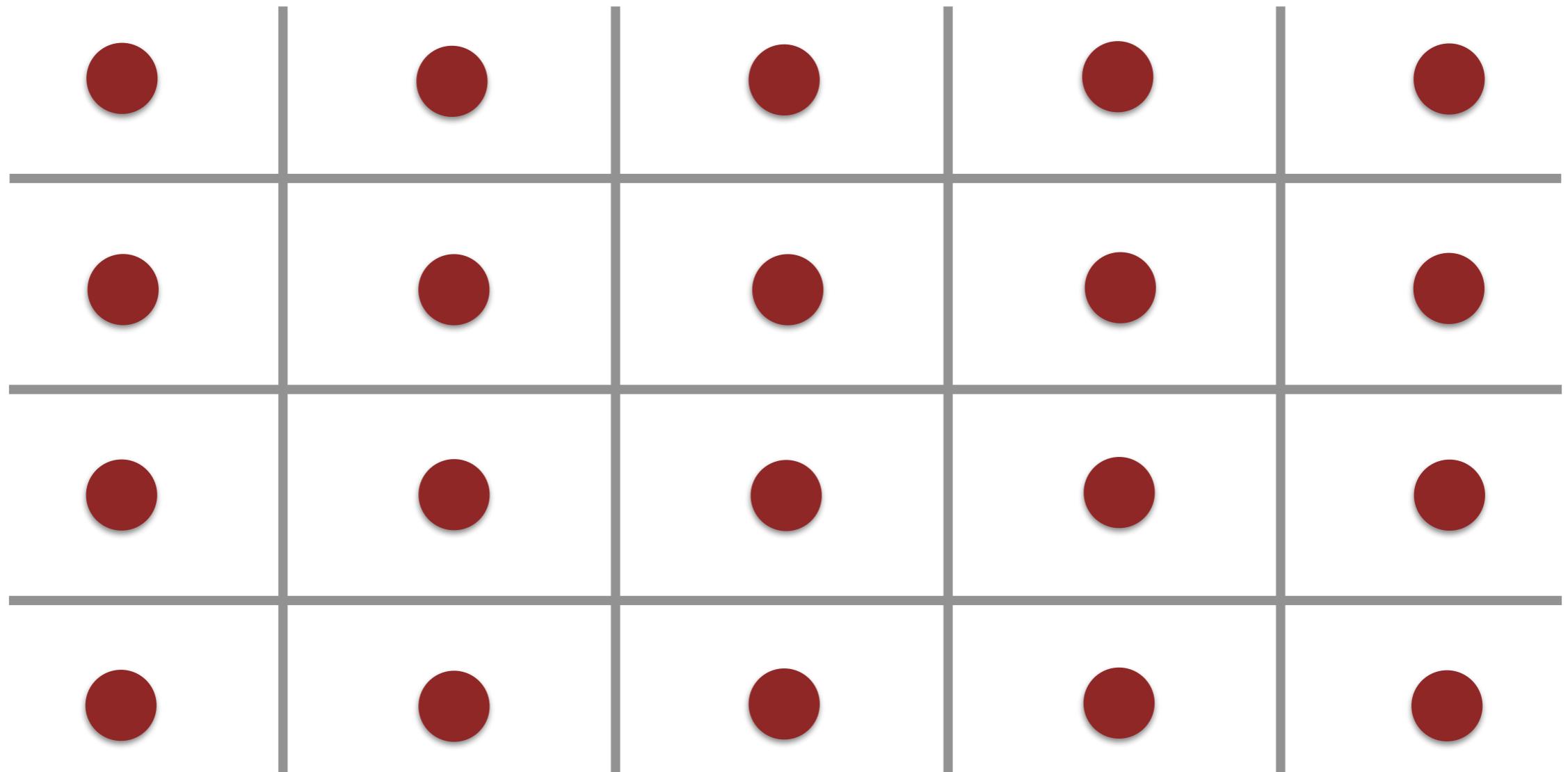


$$y_n \sim \mathcal{N}(\mathbf{X}^T \boldsymbol{\beta} + \alpha, \sigma)$$

Modeling every individual separately avoids any bias, but then the data becomes very sparse and inferences weak.

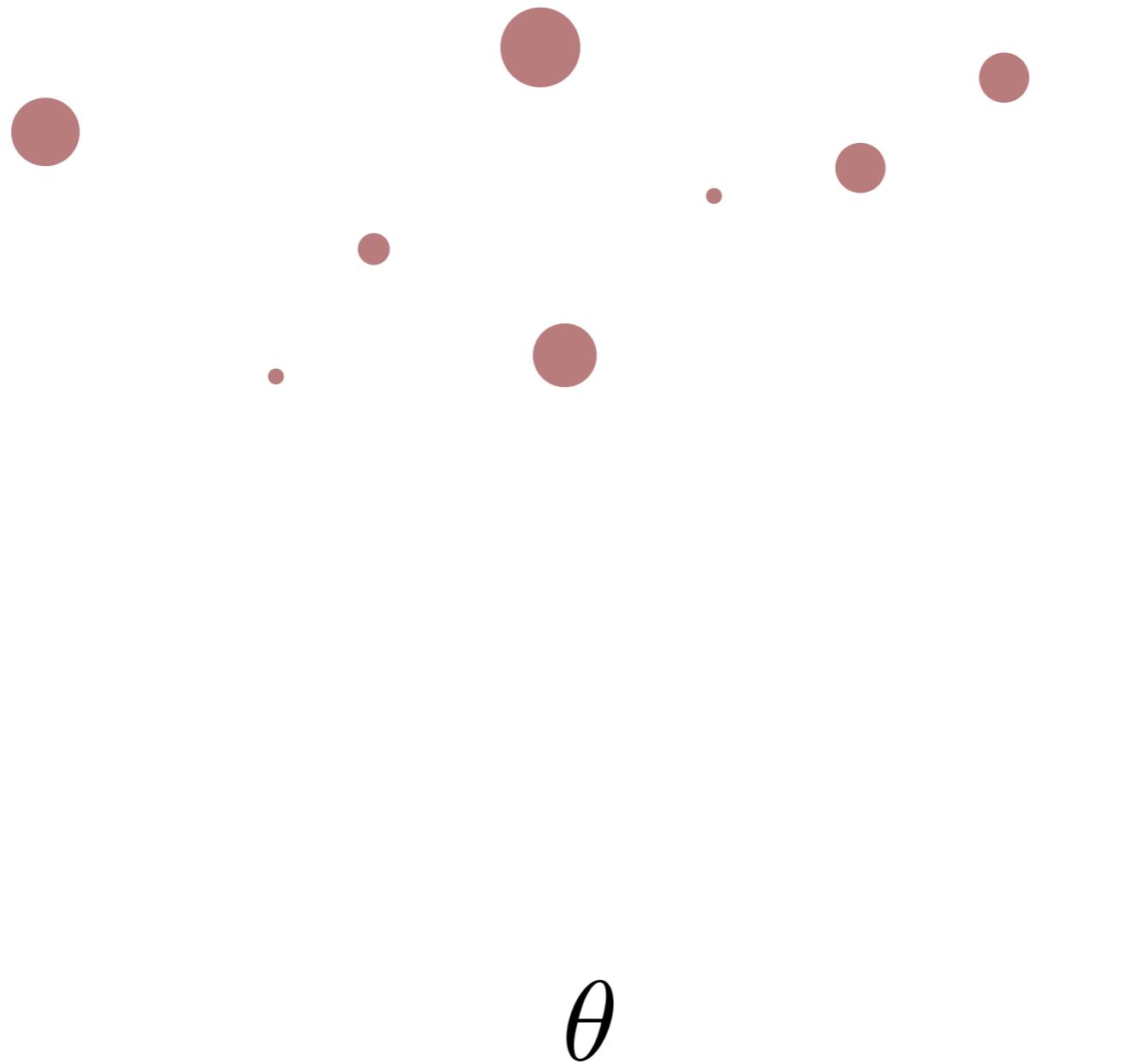


Modeling every individual separately avoids any bias, but then the data becomes very sparse and inferences weak.

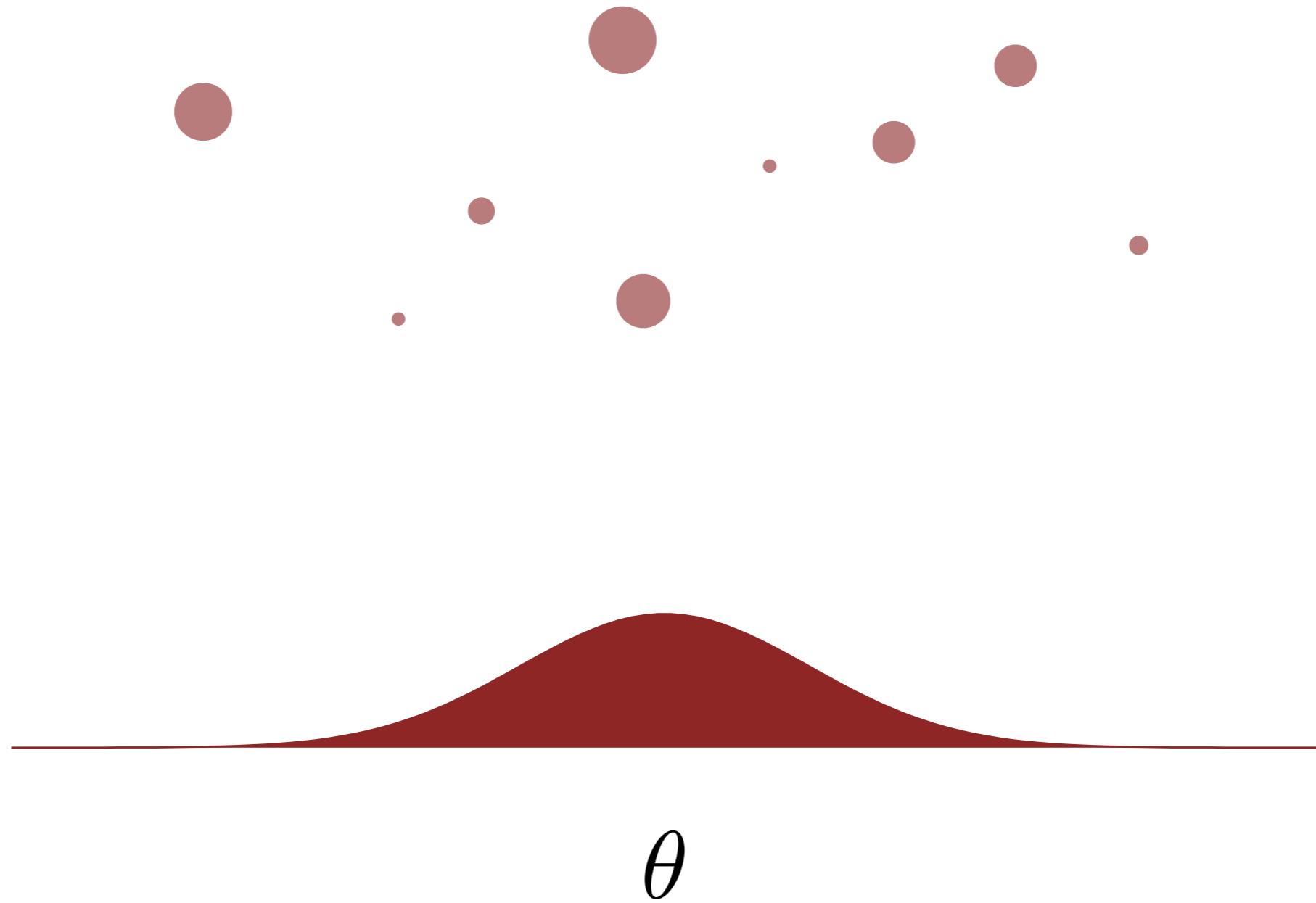


$$y_n \sim \mathcal{N}(\mathbf{X}^T \boldsymbol{\beta}_n + \alpha_n, \sigma)$$

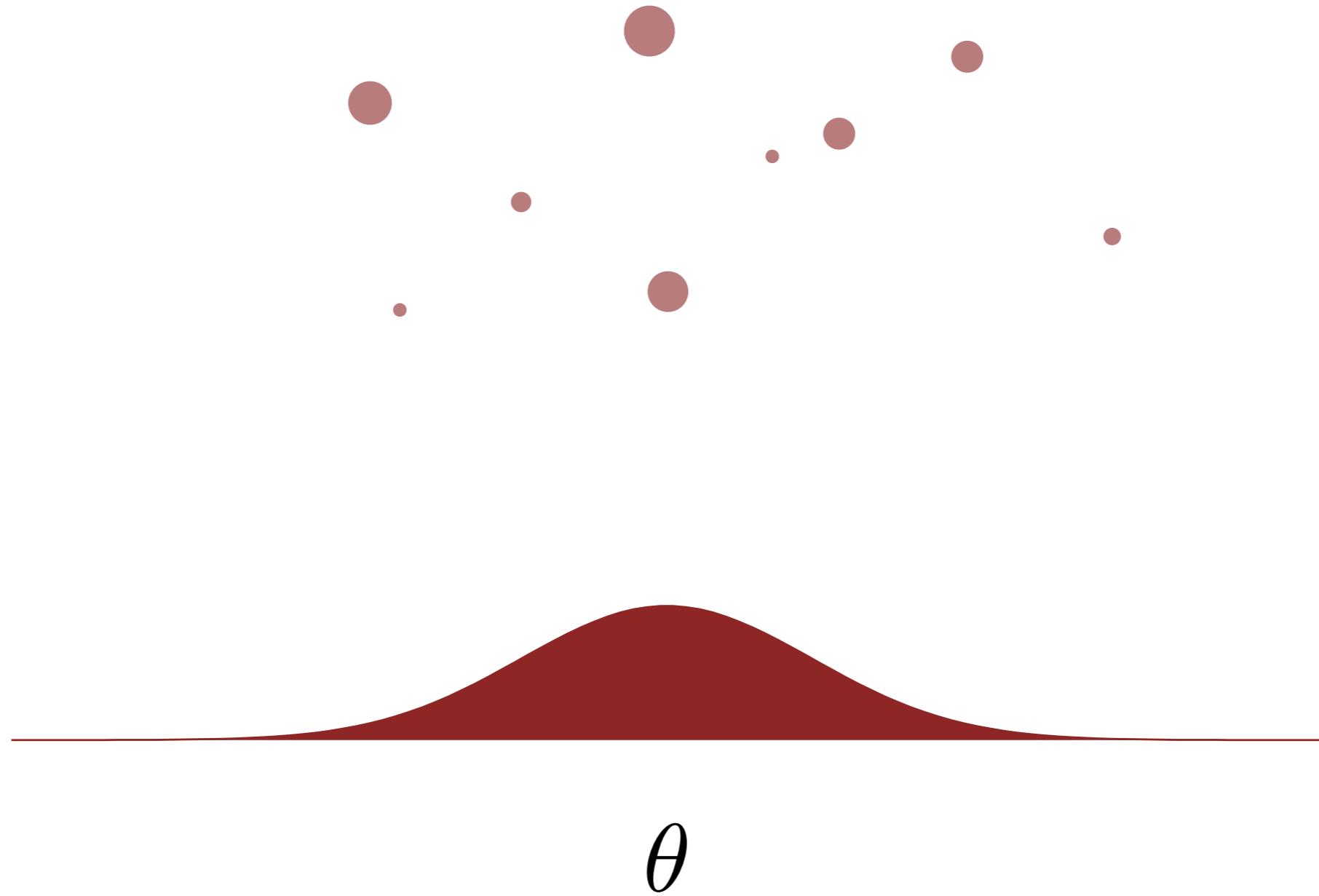
A compromise between complete pooling and no pooling that could balance bias and variance would be ideal.



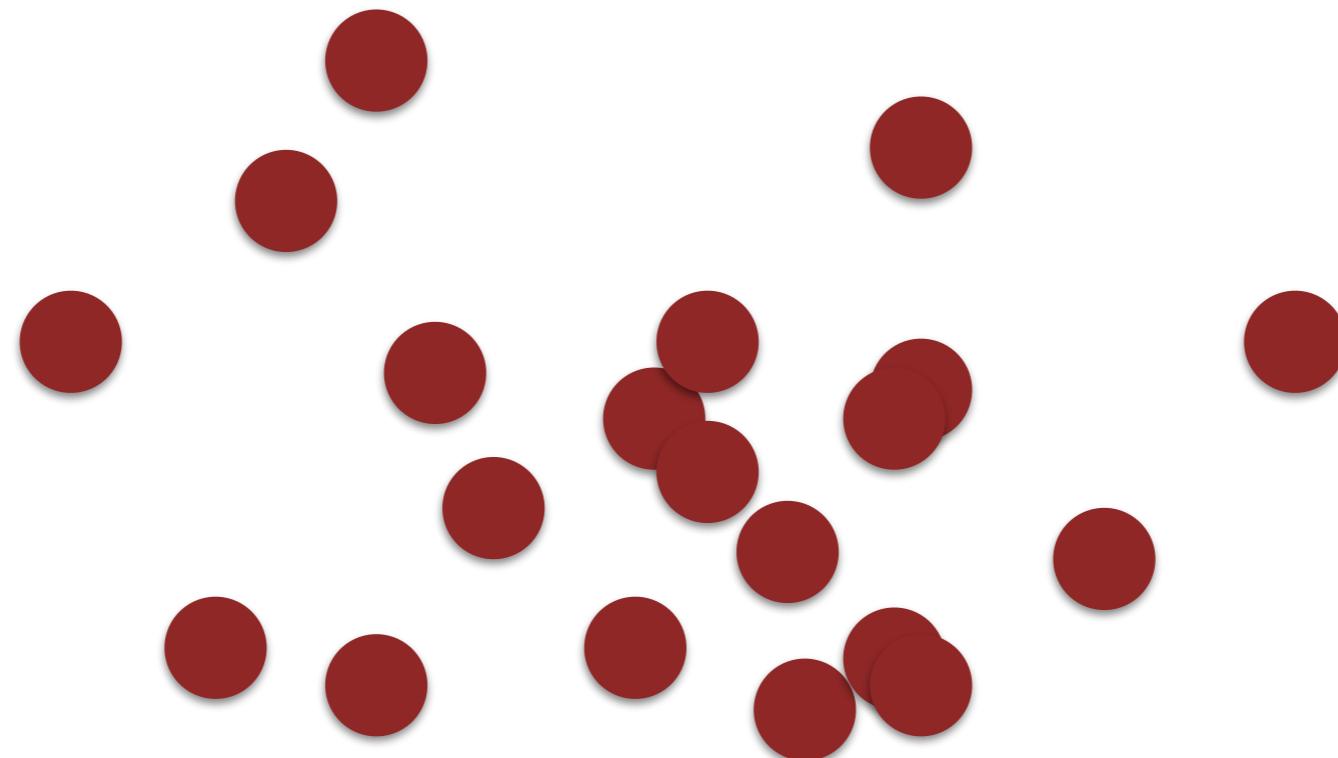
A compromise between complete pooling and no pooling  
that could balance bias and variance would be ideal.



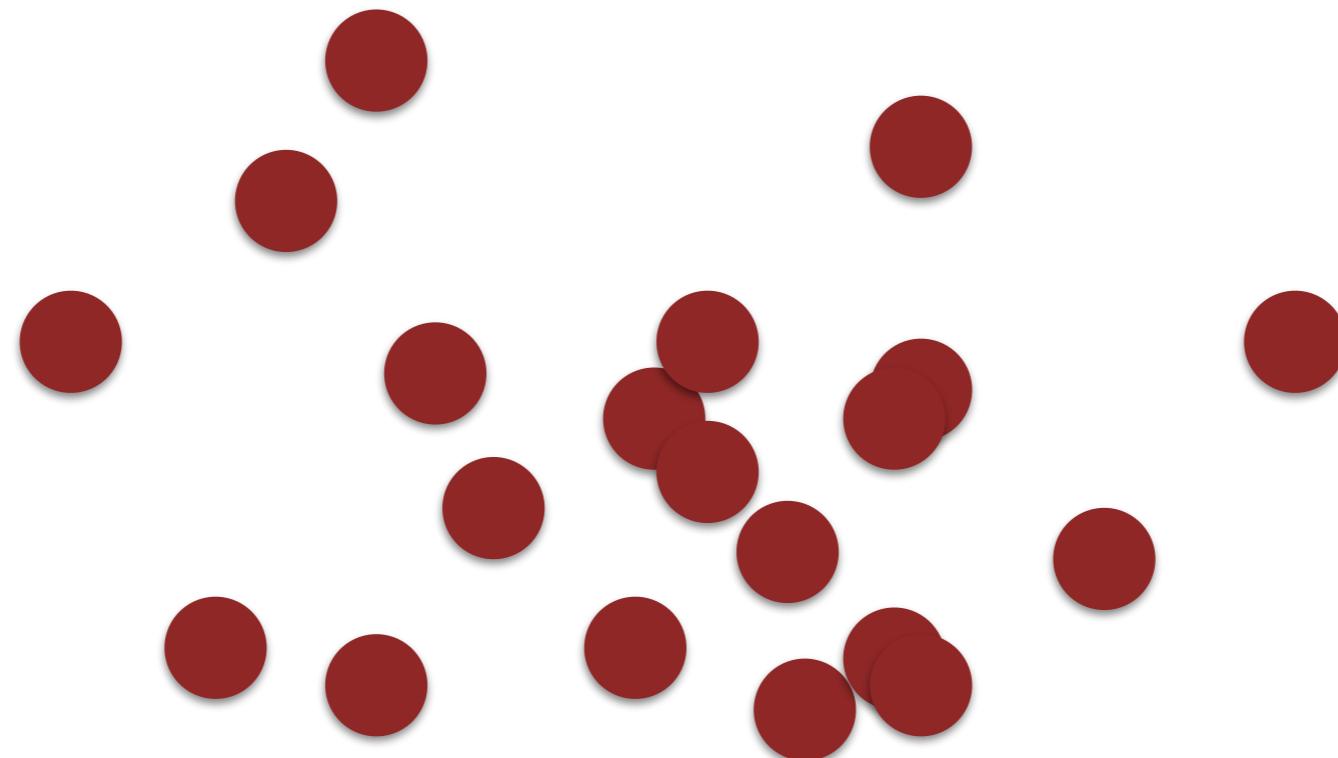
A compromise between complete pooling and no pooling  
that could balance bias and variance would be ideal.



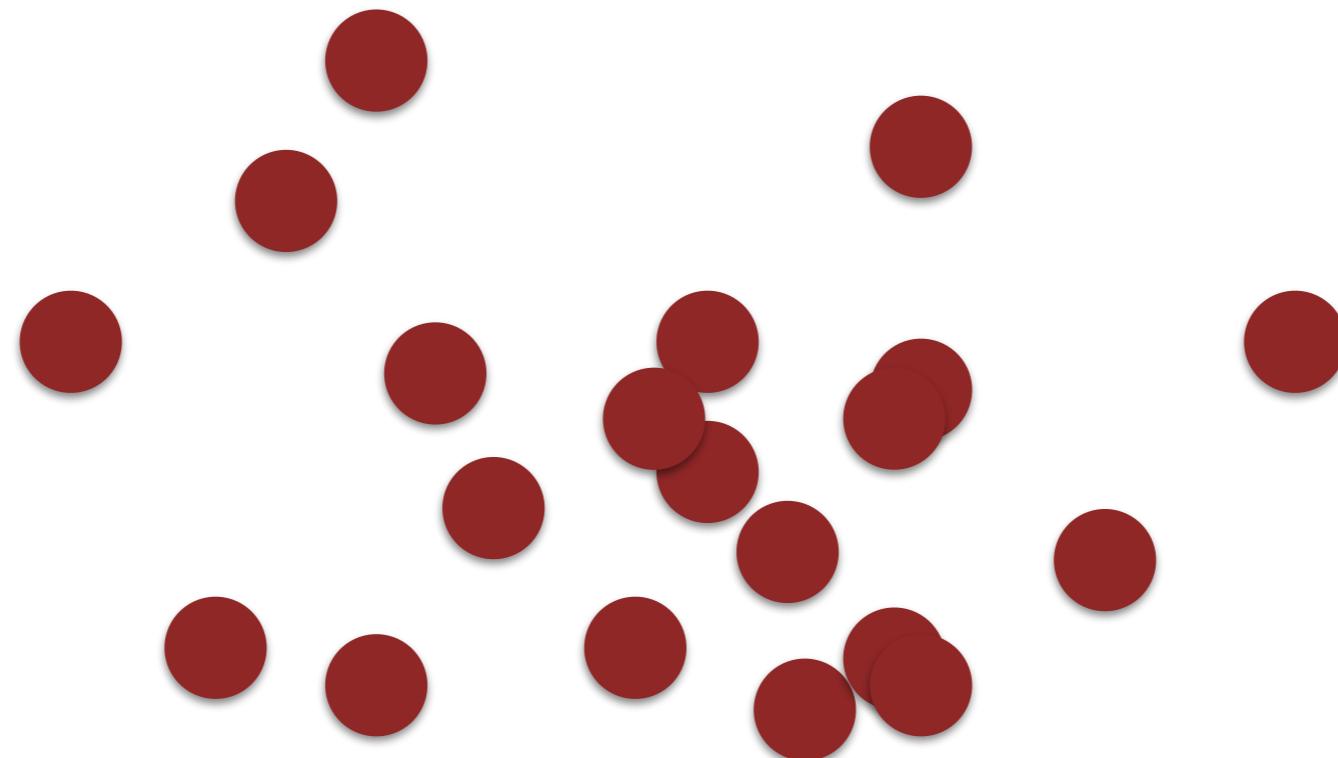
In order to formalize this approach we  
need to consider *exchangeability*.



In order to formalize this approach we  
need to consider *exchangeability*.



In order to formalize this approach we  
need to consider *exchangeability*.



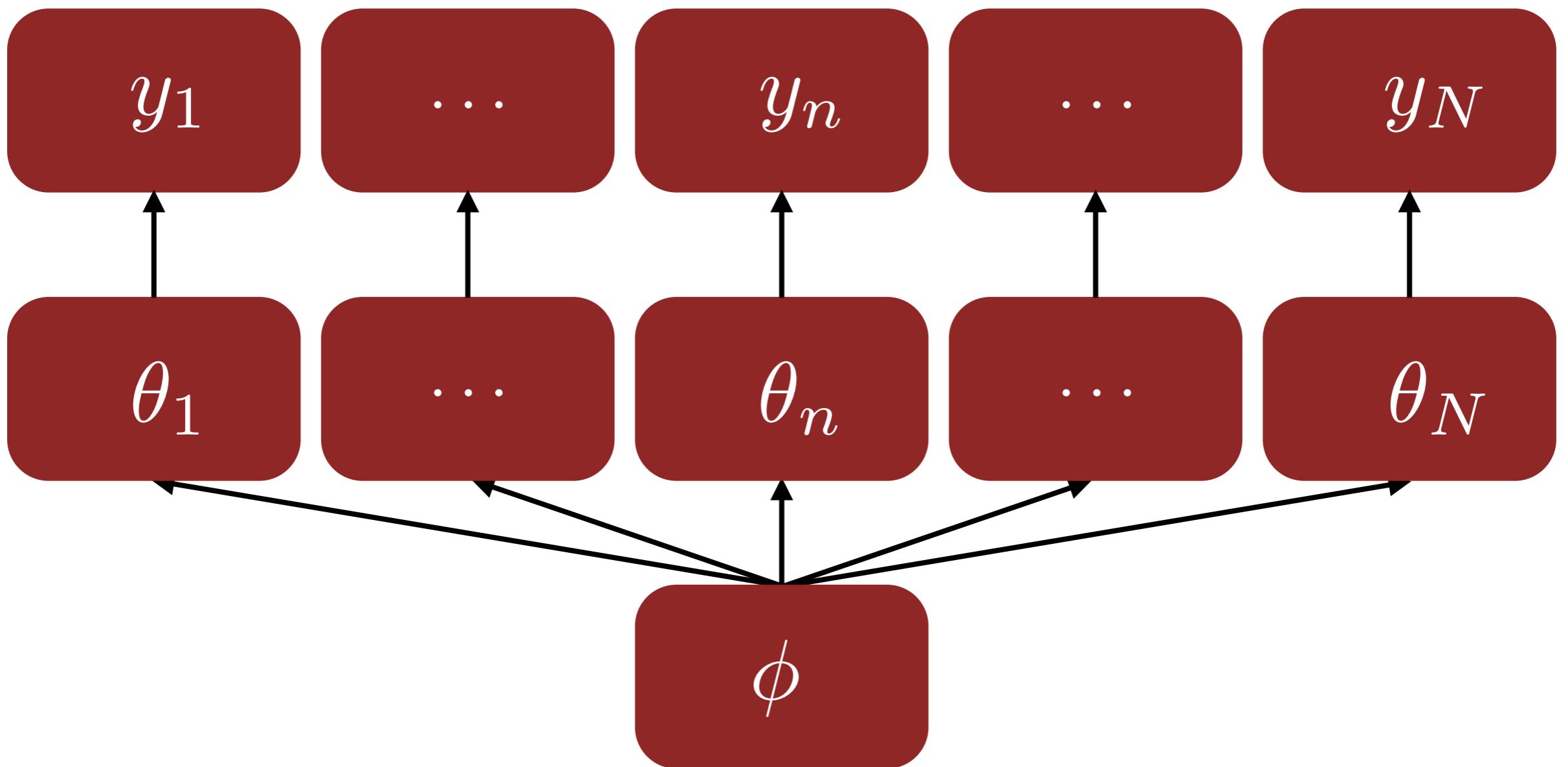
As the population grows, the only prior distribution that respects exchangeability is a *hierarchical* distribution.

$$\pi(\theta) = \int d\phi \prod_{n=1}^N \pi(\theta_n | \phi) \pi(\phi)$$

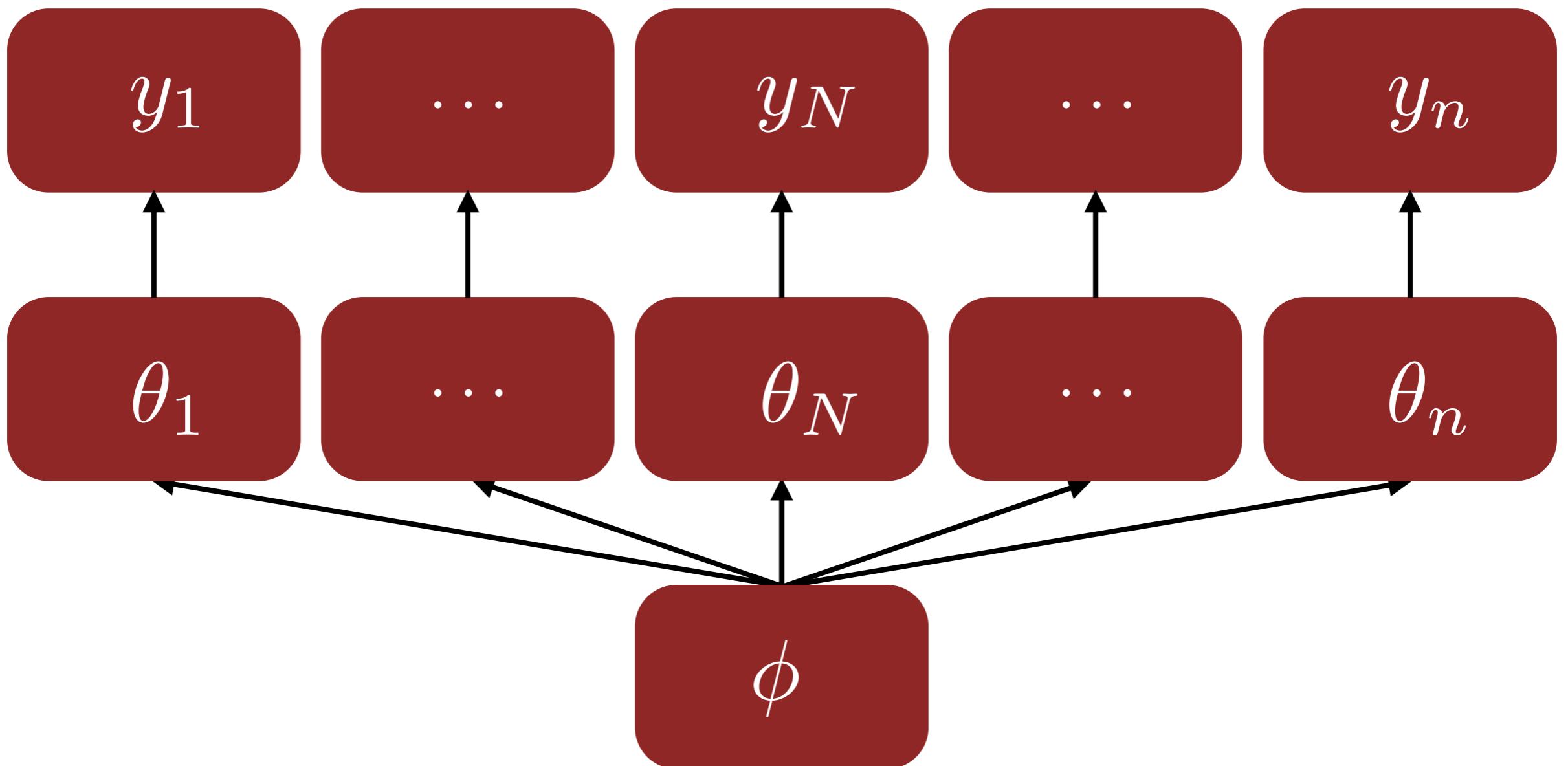
As the population grows, the only prior distribution that respects exchangeability is a *hierarchical* distribution.

$$\pi(\mathbf{y}, \boldsymbol{\theta}) = \int d\phi \prod_{n=1}^N \pi(y_n | \theta_n) \pi(\theta_n | \phi) \pi(\phi)$$

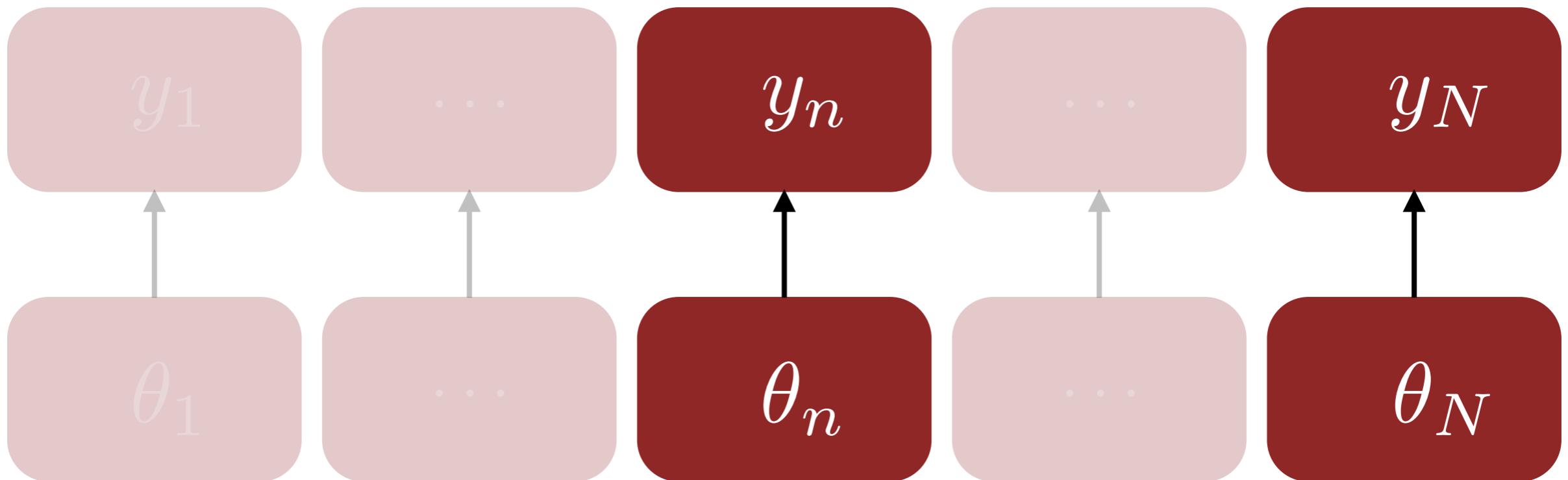
As the population grows, the only distribution that respects exchangeability is a *hierarchical* distribution.



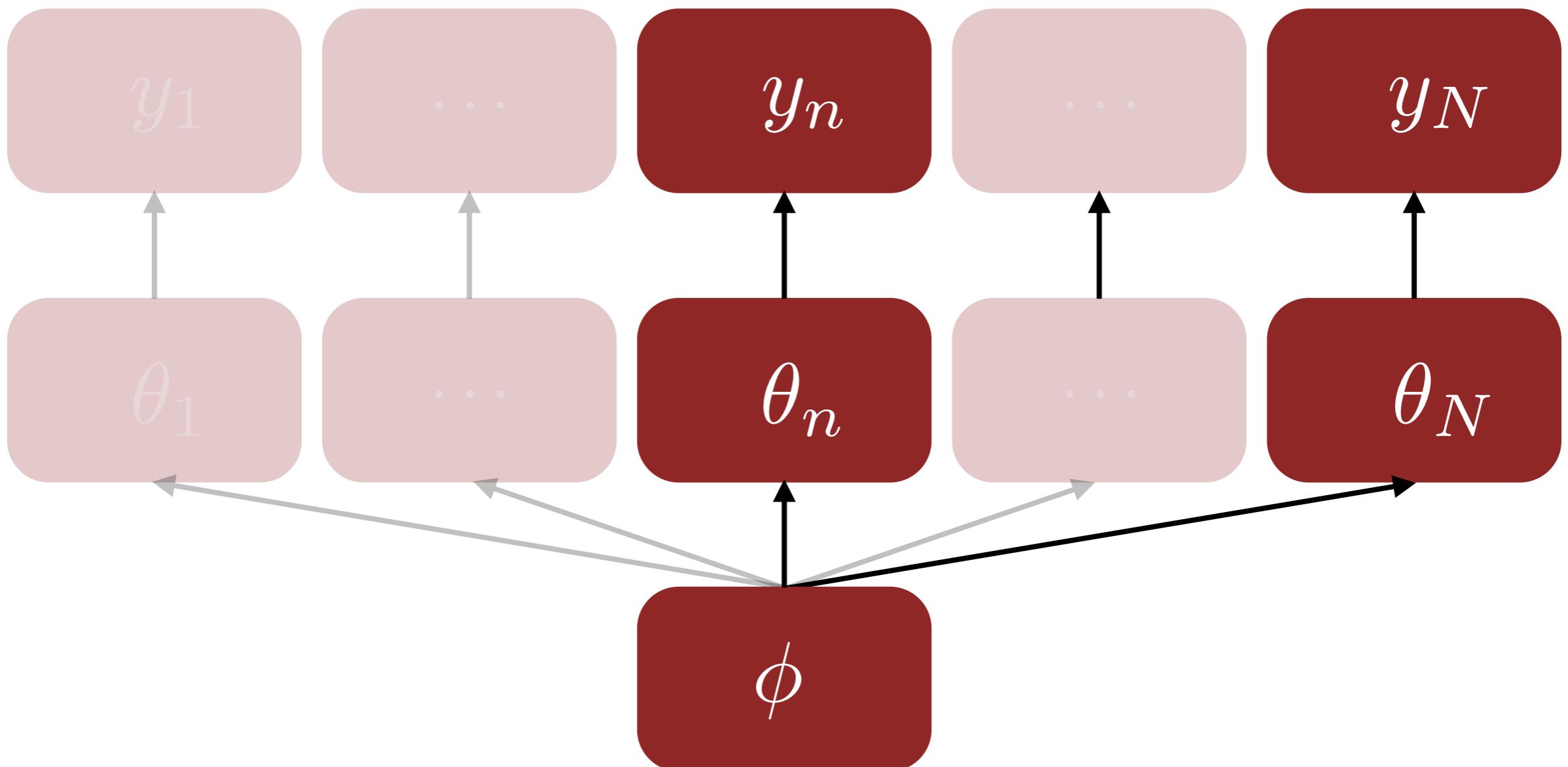
As the population grows, the only distribution that respects exchangeability is a *hierarchical* distribution.



The hyperparameters couple all of the groups together, partially pooling the data and balancing bias and variance.



The hyperparameters couple all of the groups together, partially pooling the data and balancing bias and variance.



The most common population model is a Gaussian, where we model the population mean and variance.

$$\prod_{n=1}^N \pi(y_n | \theta_n) \pi(\theta_n | \phi) \pi(\phi)$$

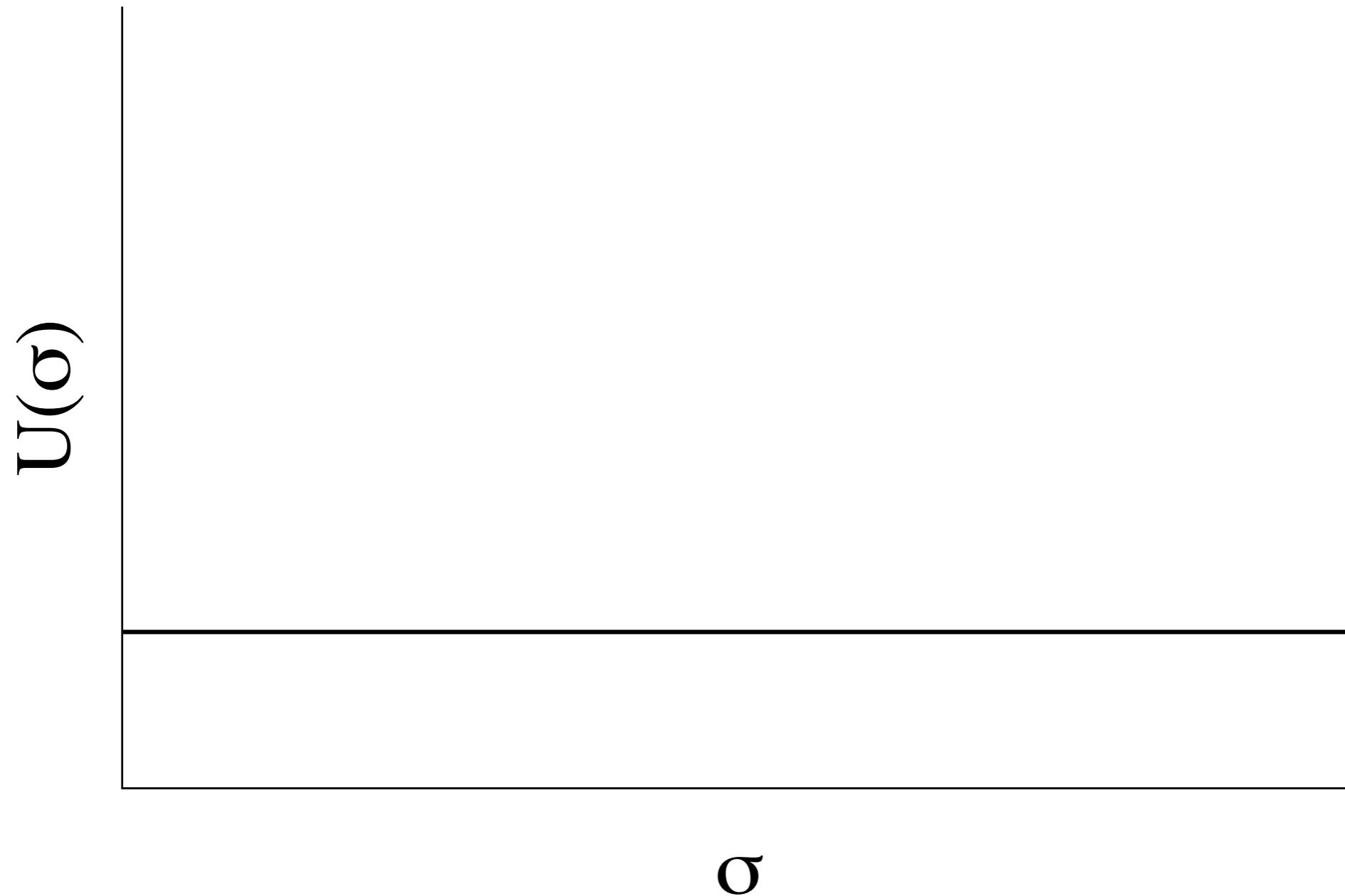
The most common population model is a Gaussian, where we model the population mean and variance.

$$\prod_{n=1}^N \pi(y_n | \theta_n) \mathcal{N}(\theta_n | \mu, \sigma) \pi(\mu) \pi(\sigma)$$

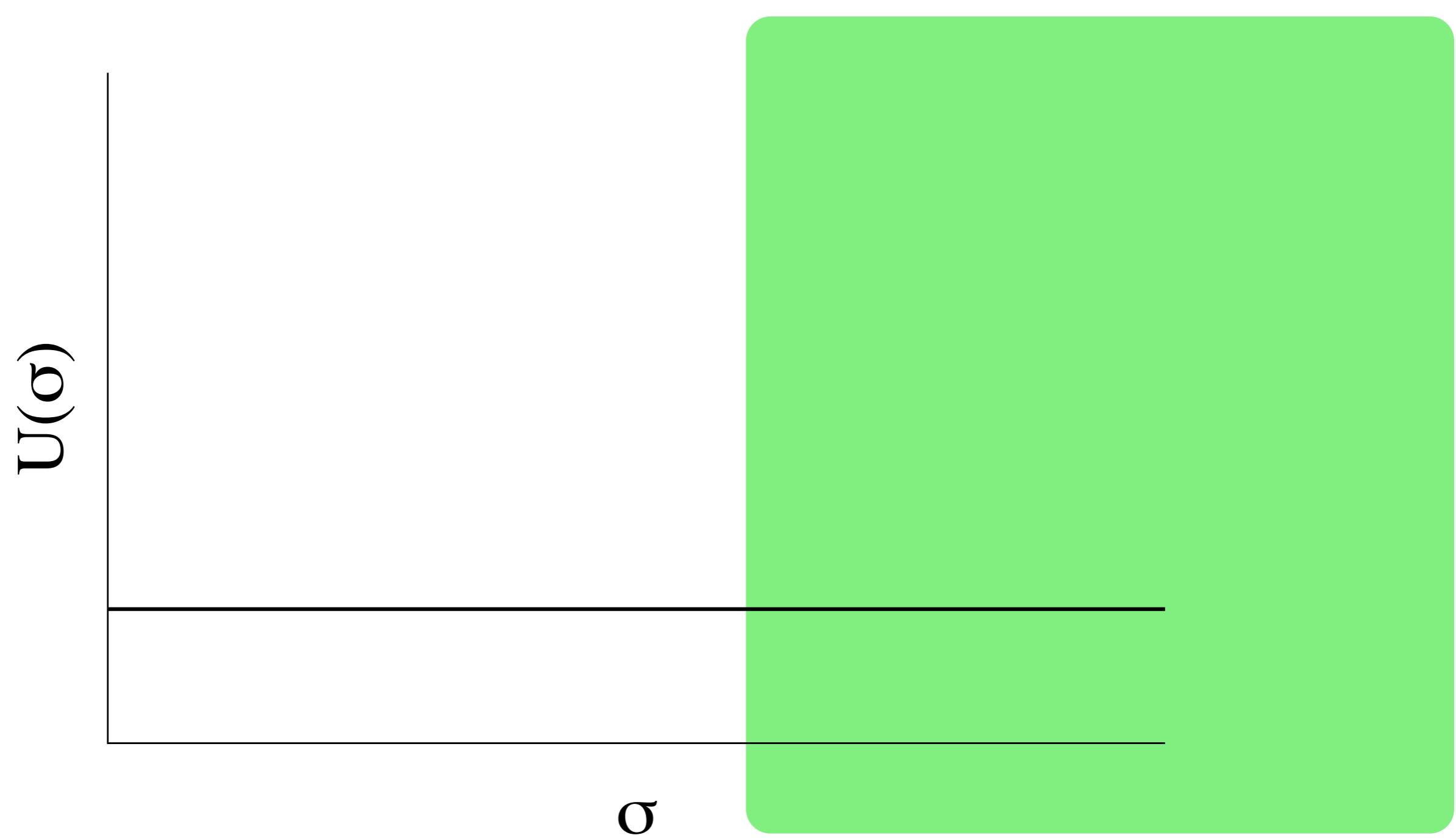
Informative prior distributions on the population hyperparameters are incredibly important.

$$\pi(\mu) \pi(\sigma)$$

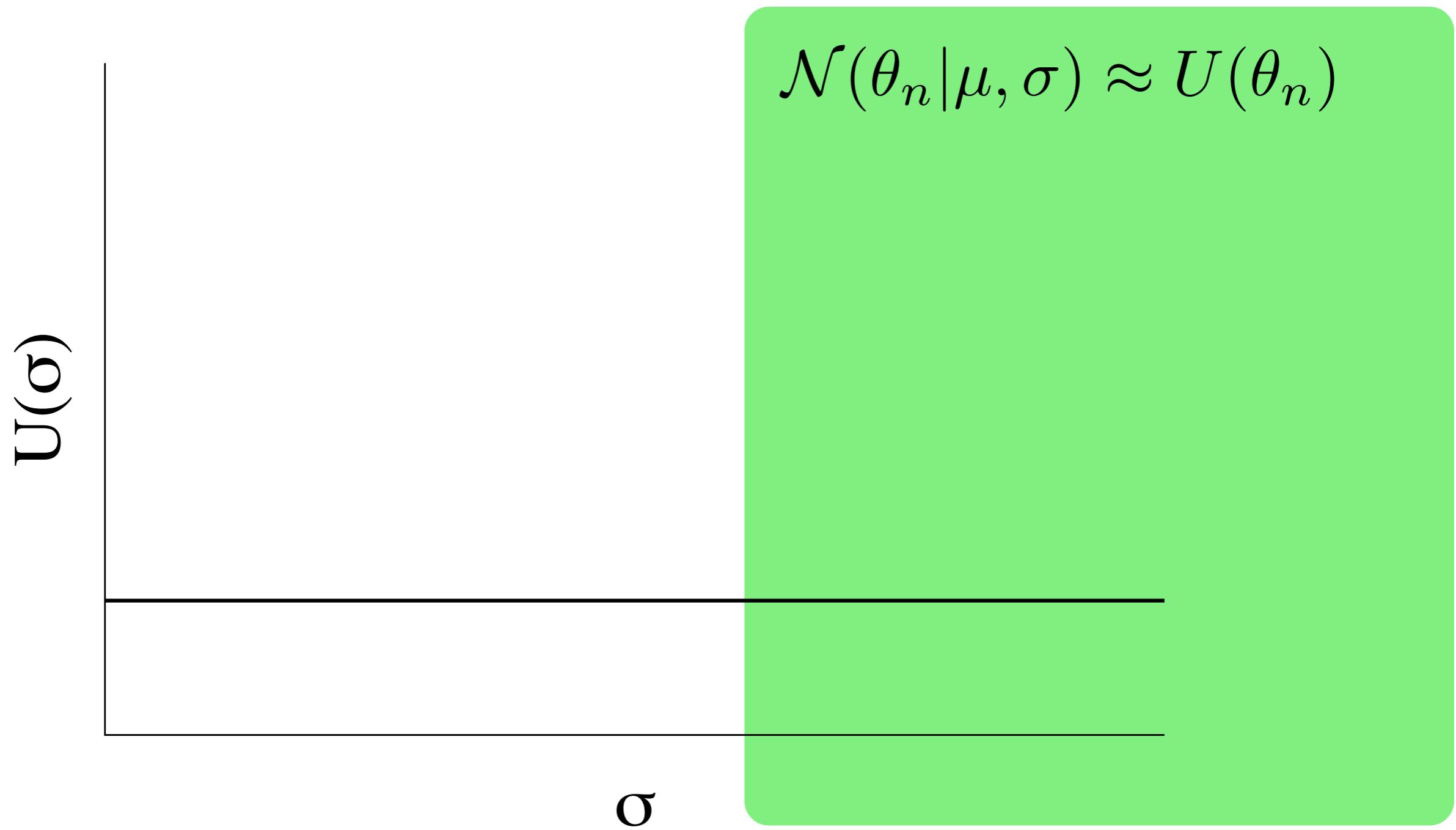
A uniform prior on the population deviance places too much probability at high values and no pooling.



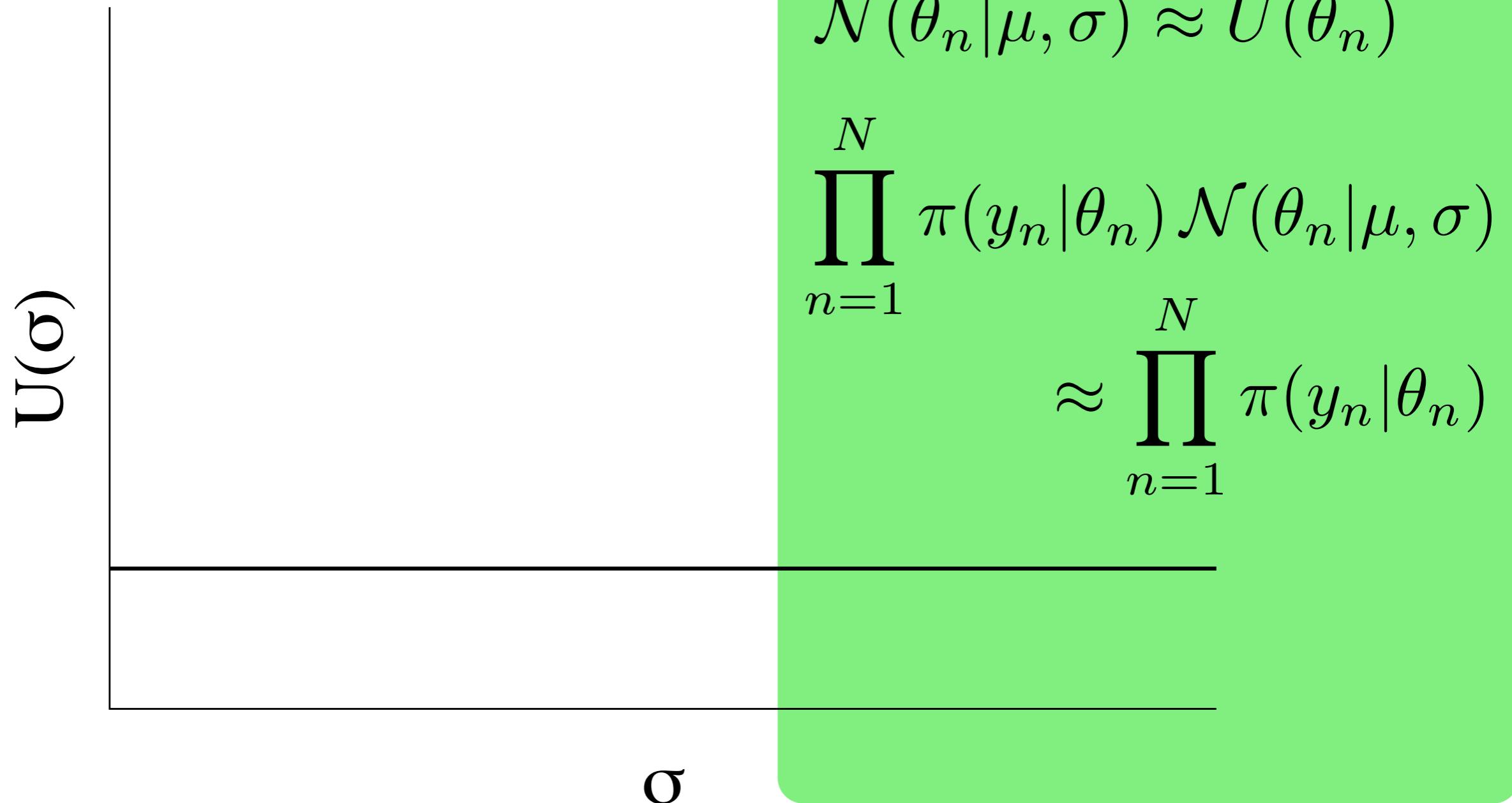
A uniform prior on the population deviance places too much probability at high values and no pooling.



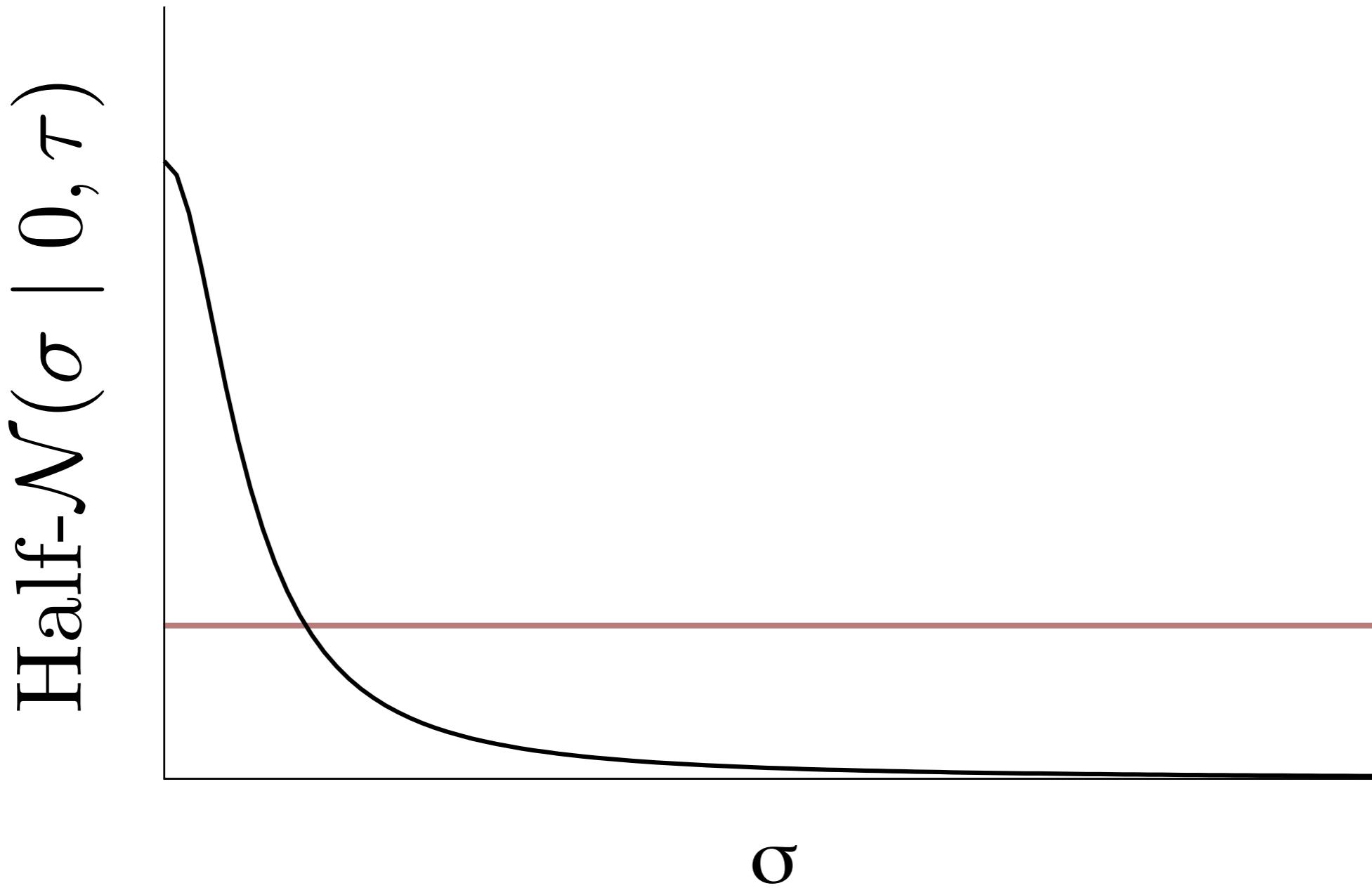
A uniform prior on the population deviance places too much probability at high values and no pooling.



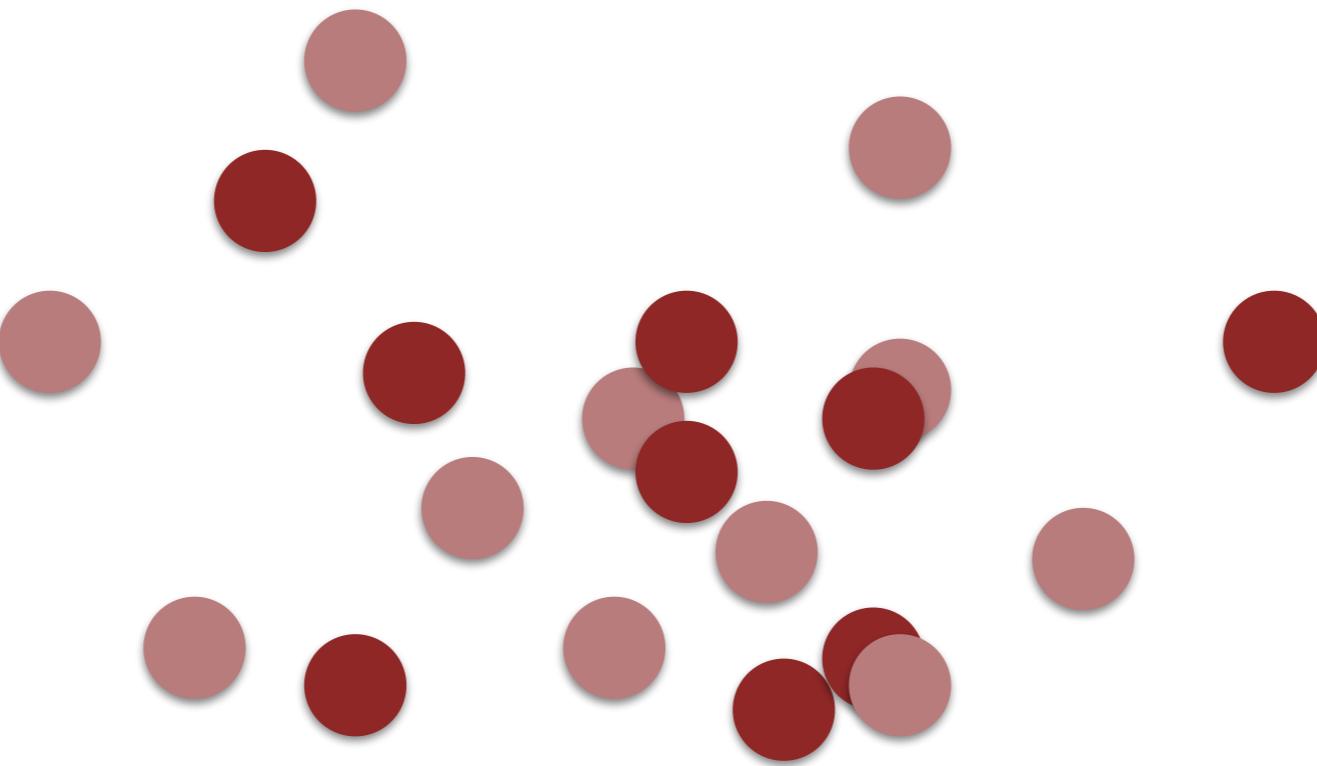
A uniform prior on the population deviance places too much probability at high values and no pooling.



For the model to be able to partially pool we need a weakly informative prior that concentrates around zero.



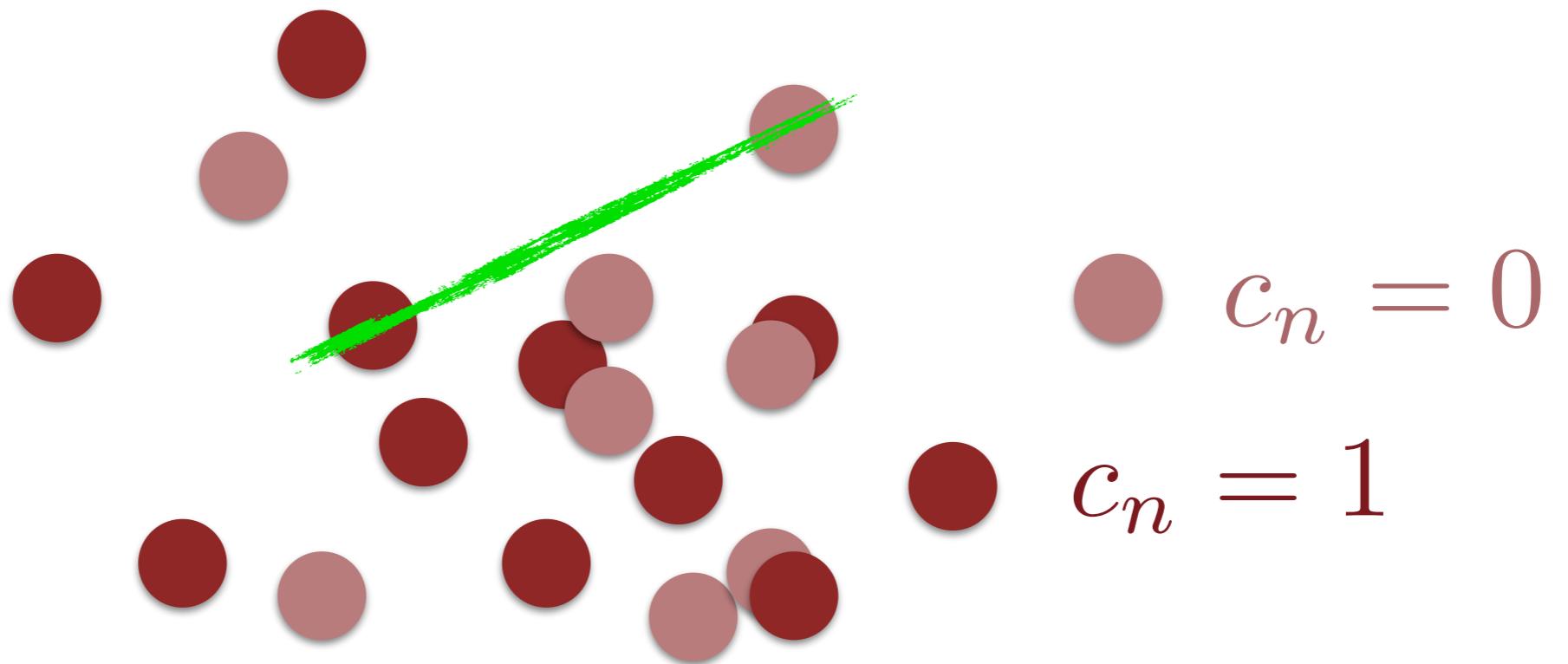
# Multilevel Models



Given labels that discriminate between individuals,  
the individuals are no longer exchangeable.



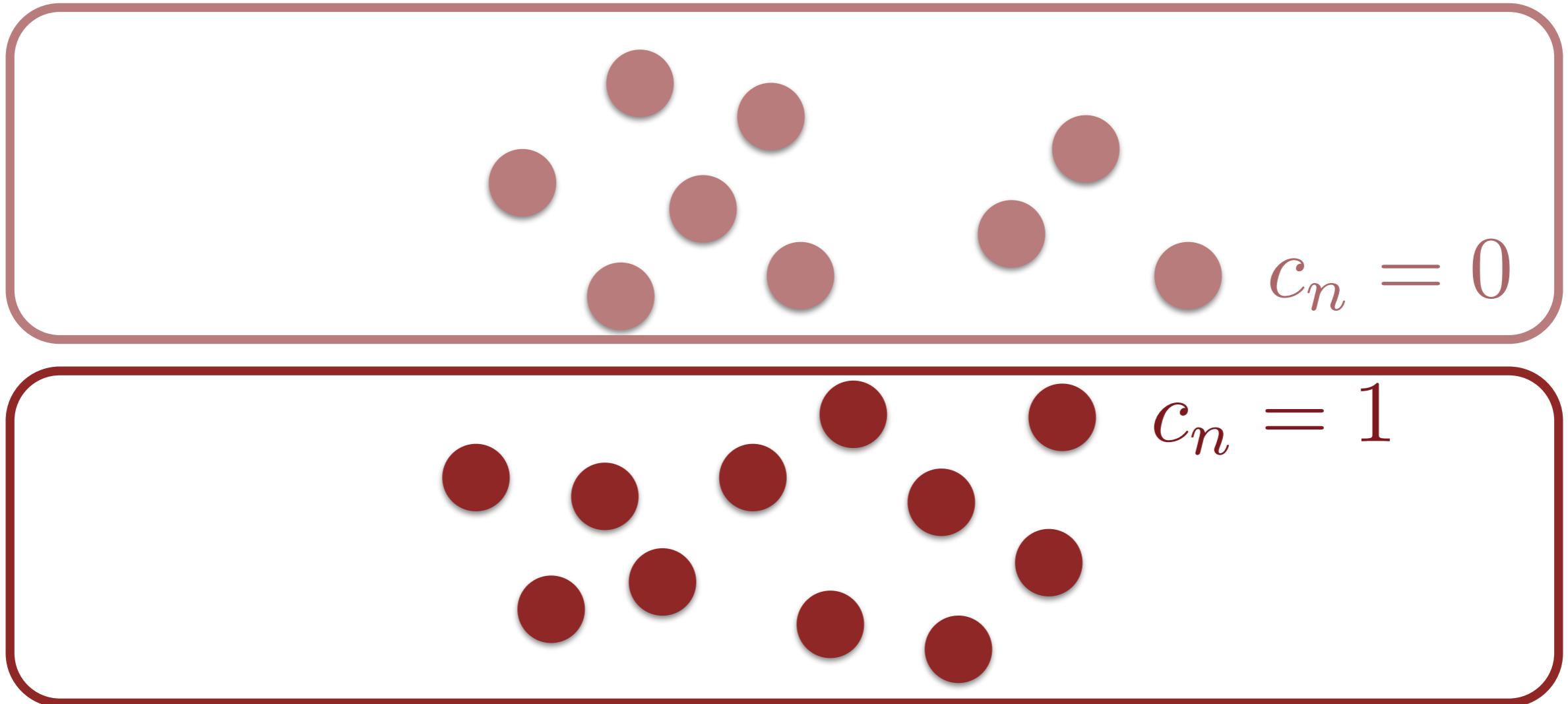
Given labels that discriminate between individuals,  
the individuals are no longer exchangeable.



The partitions of the population, however, define exchangeable groups within a metapopulation.



The partitions of the population, however, define exchangeable groups within a metapopulation.



The partitions of the population, however, define exchangeable groups within a metapopulation.

$$c_n = 0$$

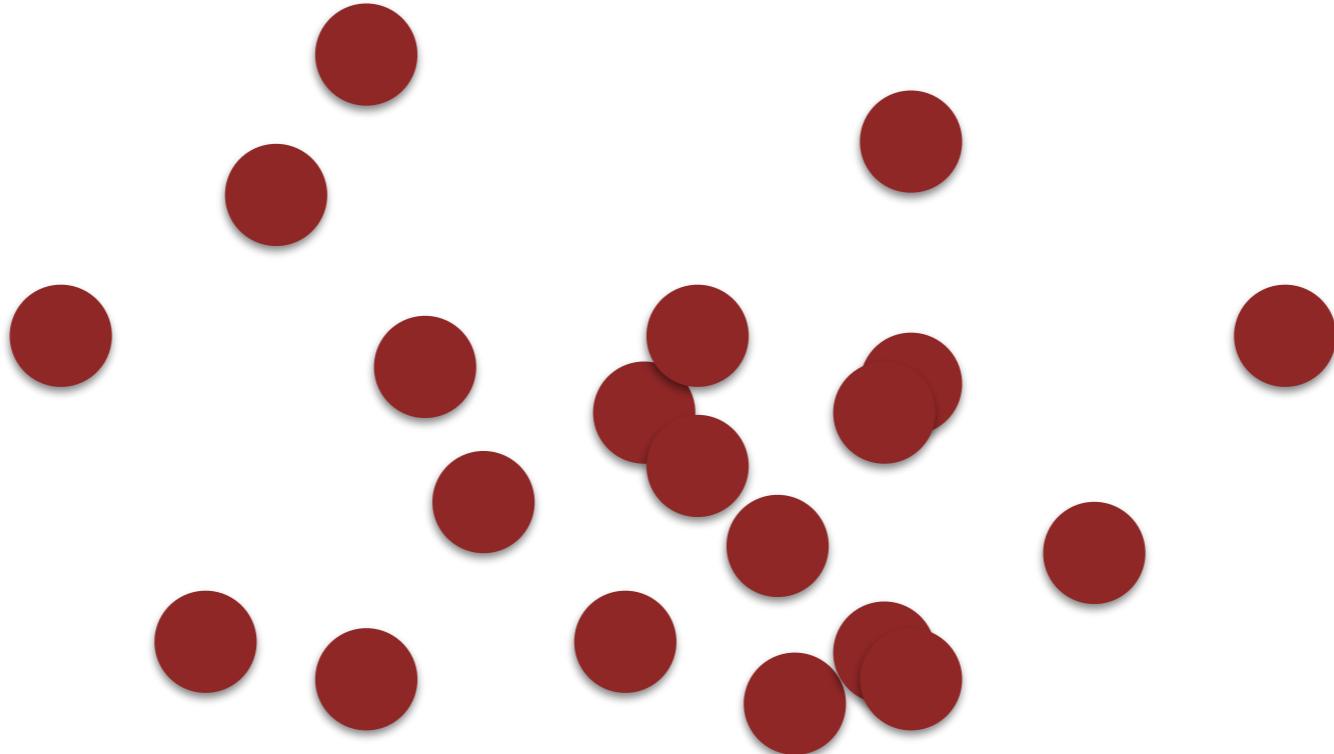
$$c_n = 1$$

The partitions of the population, however, define exchangeable groups within a metapopulation.

$$c_n = 0$$

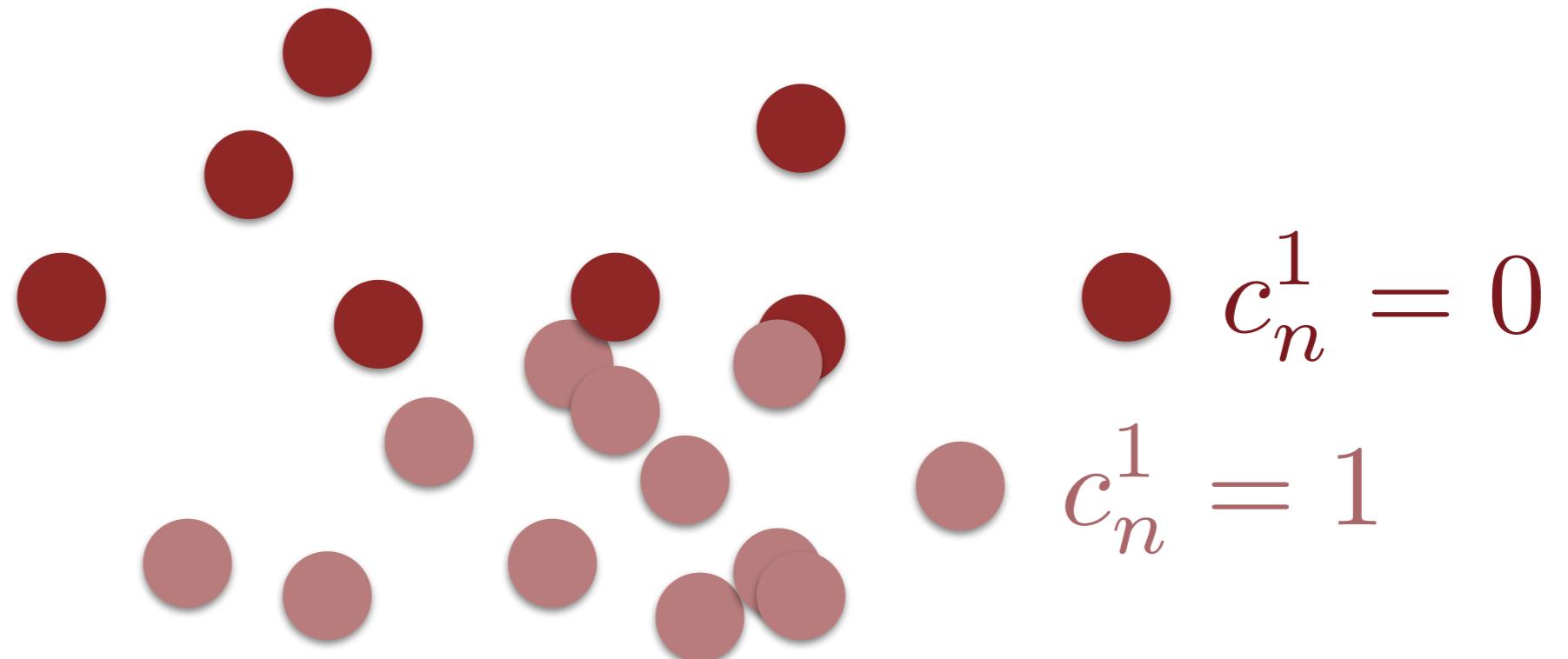
$$c_n = 1$$

The partitions of the population, however, define exchangeable groups within a metapopulation.



$$\{c_n^1, c_n^2, c_n^3\} \in \{0, 1\}$$

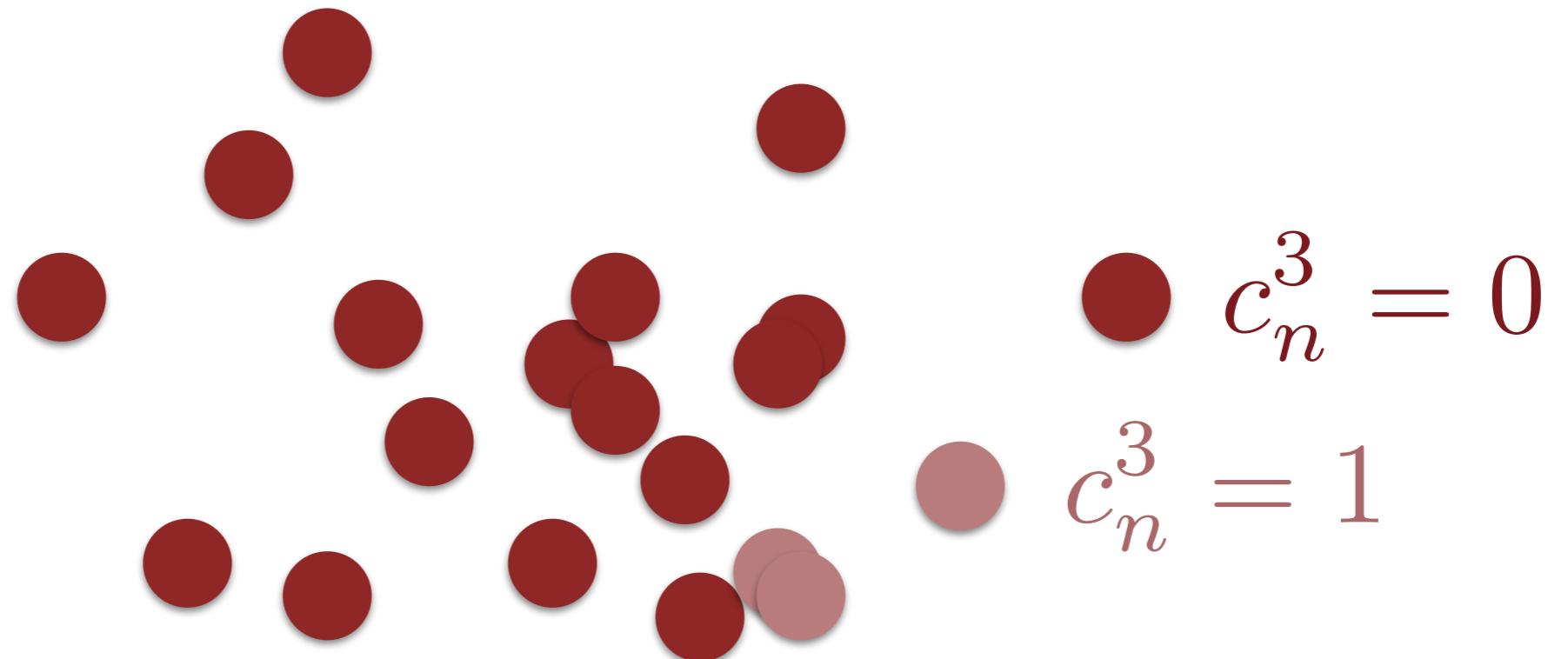
Every discriminating label defines new, exchangeable groups and new metapopulations.



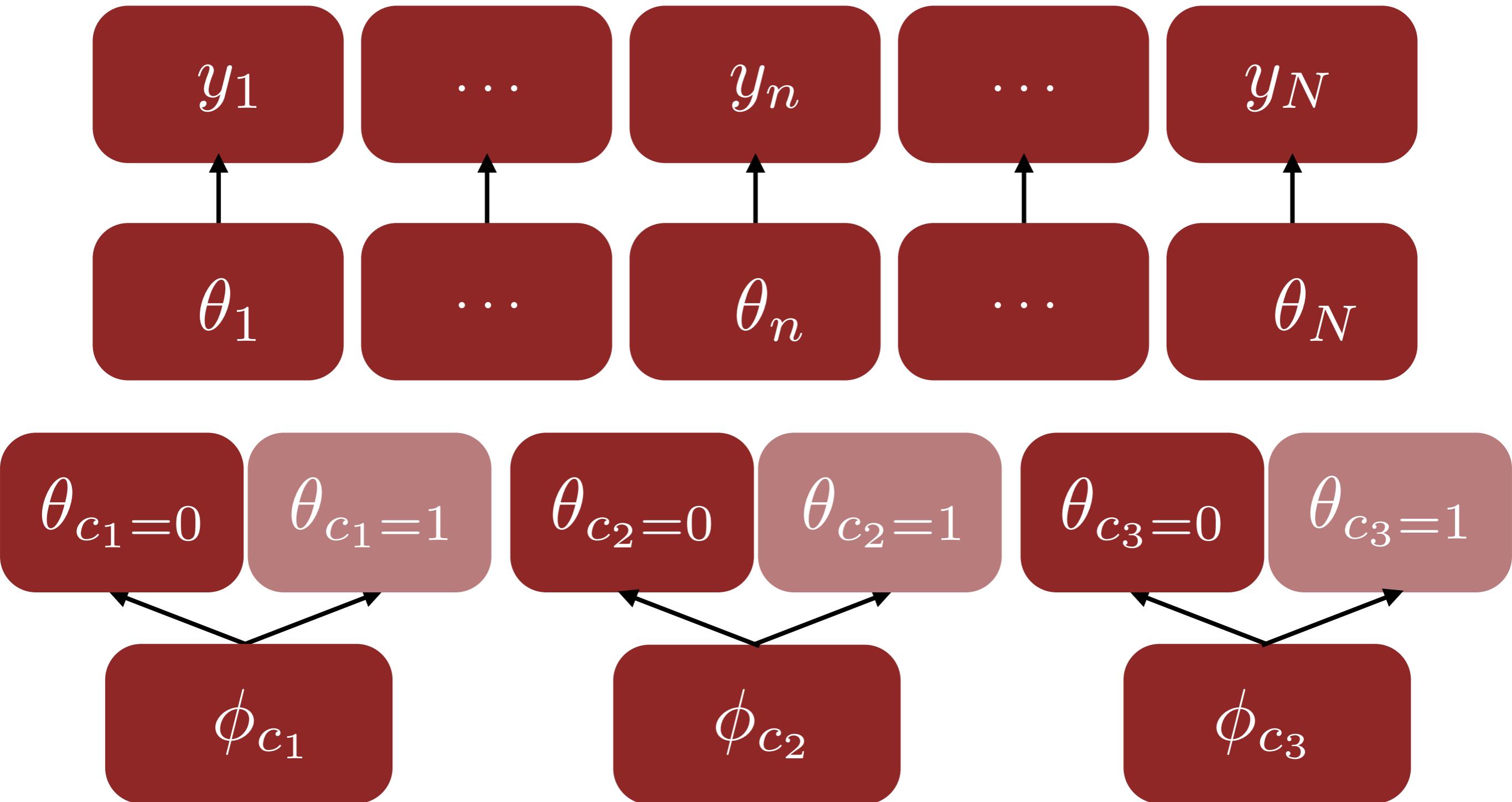
Every discriminating label defines new, exchangeable groups and new metapopulations.



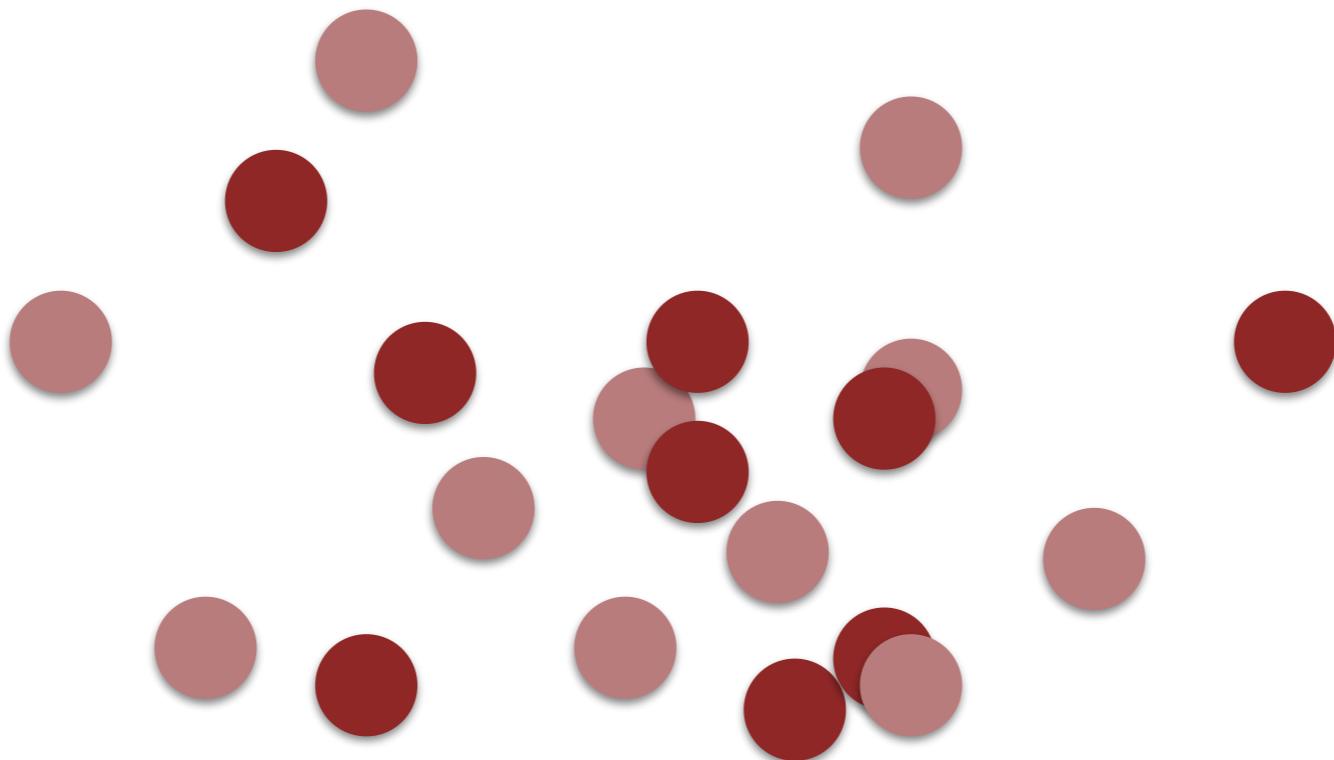
Every discriminating label defines new, exchangeable groups and new metapopulations.



In a *multilevel* model linear parameters vary across multiple hierarchies.



# General Linear Hierarchical Models



General linear models are naturally suited to exchangeability and hierarchical models.

$$\pi(y|g(\mathbf{X}^T\boldsymbol{\beta} + \boldsymbol{\alpha}), \boldsymbol{\theta})$$

General linear models are naturally suited to exchangeability and hierarchical models.

$$\pi(y_n | g(\mathbf{X}_n^T \boldsymbol{\beta}_n + \alpha_n), \theta)$$

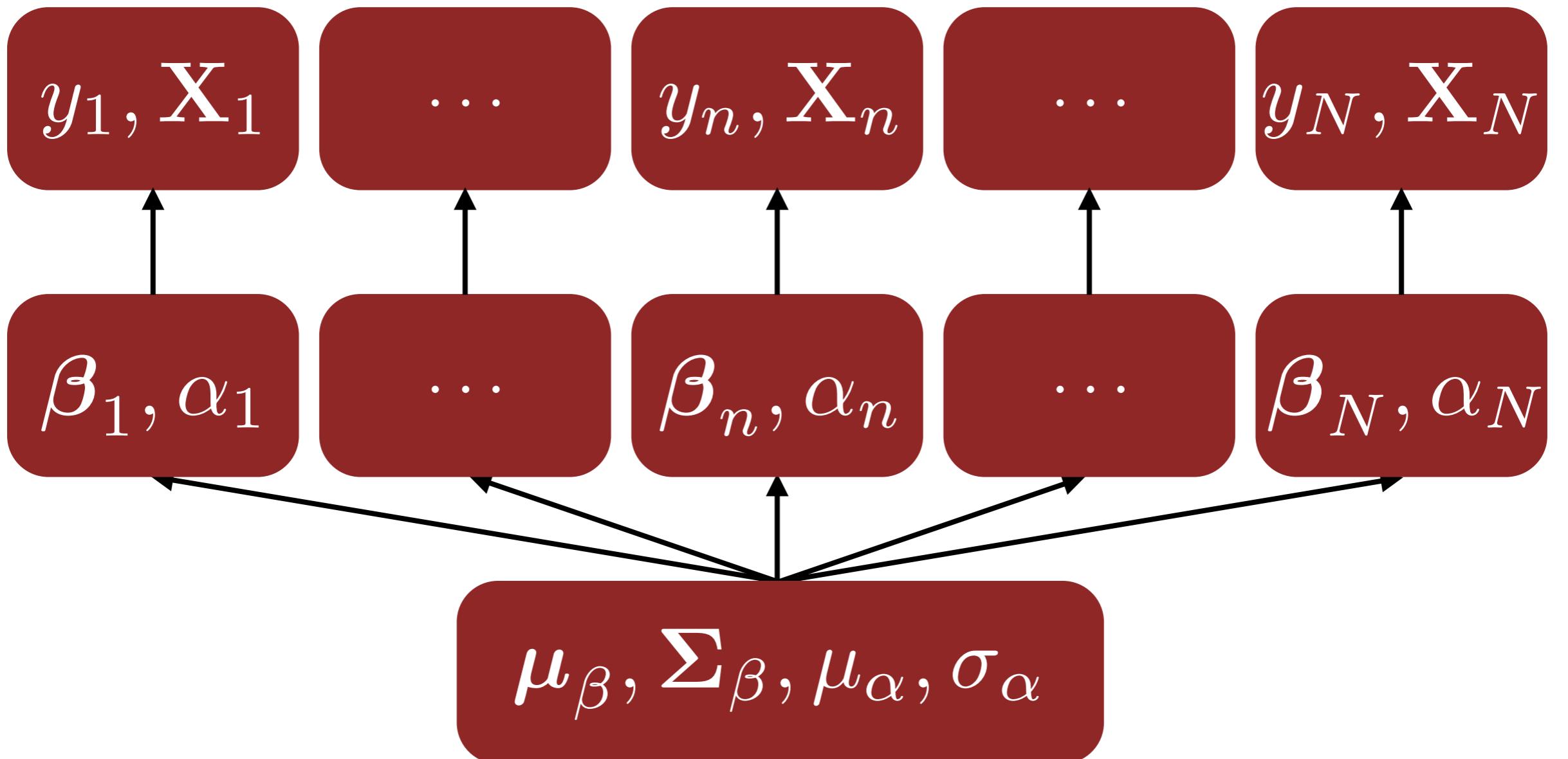
General linear models are naturally suited to exchangeability and hierarchical models.

$$\pi(y_n | g(\mathbf{X}_n^T \boldsymbol{\beta}_n + \alpha_n), \theta)$$

$$\boldsymbol{\beta}_n \sim \mathcal{N}(\boldsymbol{\mu}_{\beta}, \boldsymbol{\Sigma}_{\beta})$$

$$\alpha_n \sim \mathcal{N}(\mu_{\alpha}, \sigma_{\alpha})$$

General linear models are naturally suited to exchangeability and hierarchical models.



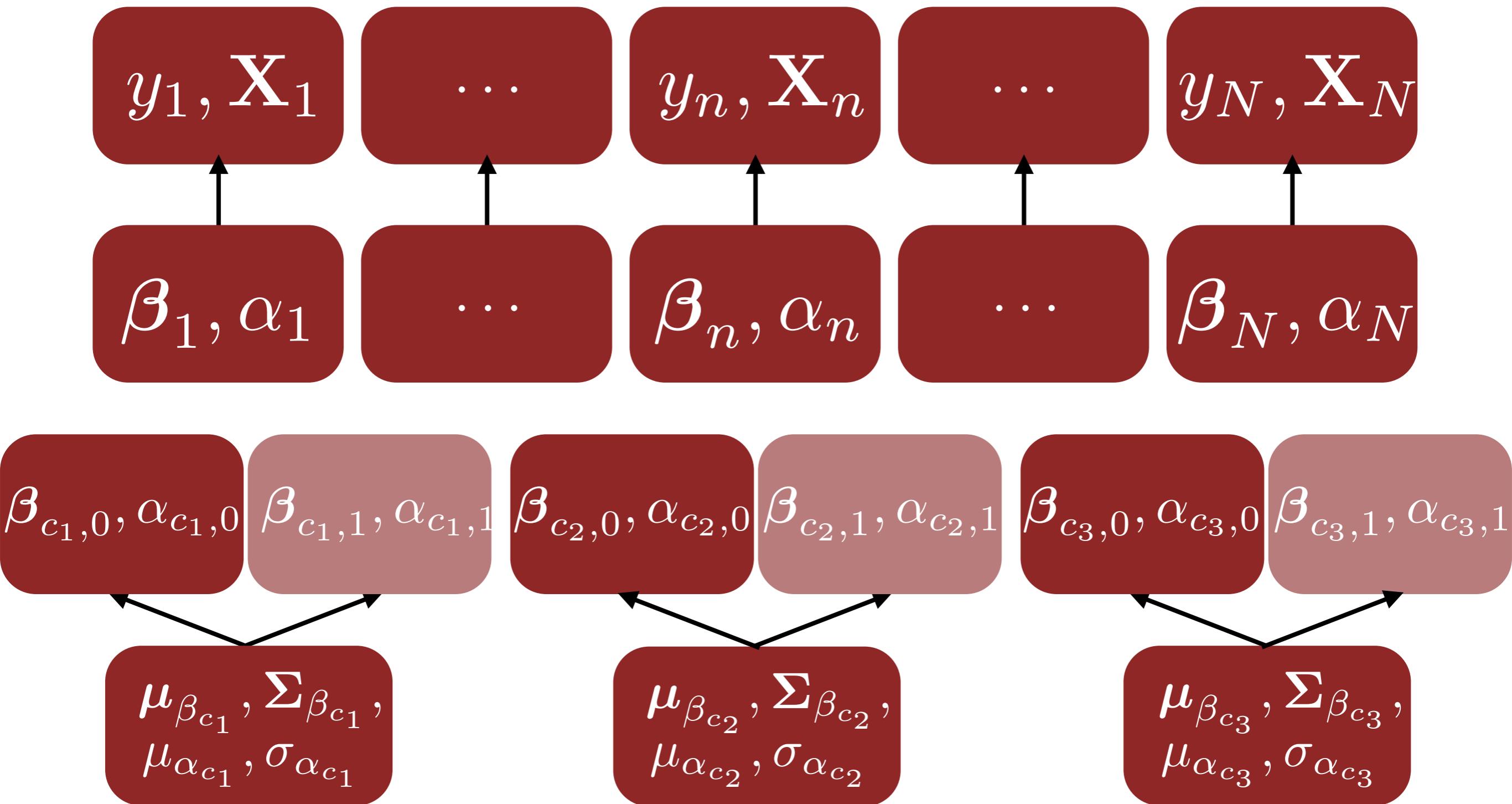
In a *multilevel* model the linear parameters vary across multiple hierarchies.

$$\pi(y_n | g(\mathbf{X}_n^T \boldsymbol{\beta}_n + \alpha_n), \theta)$$

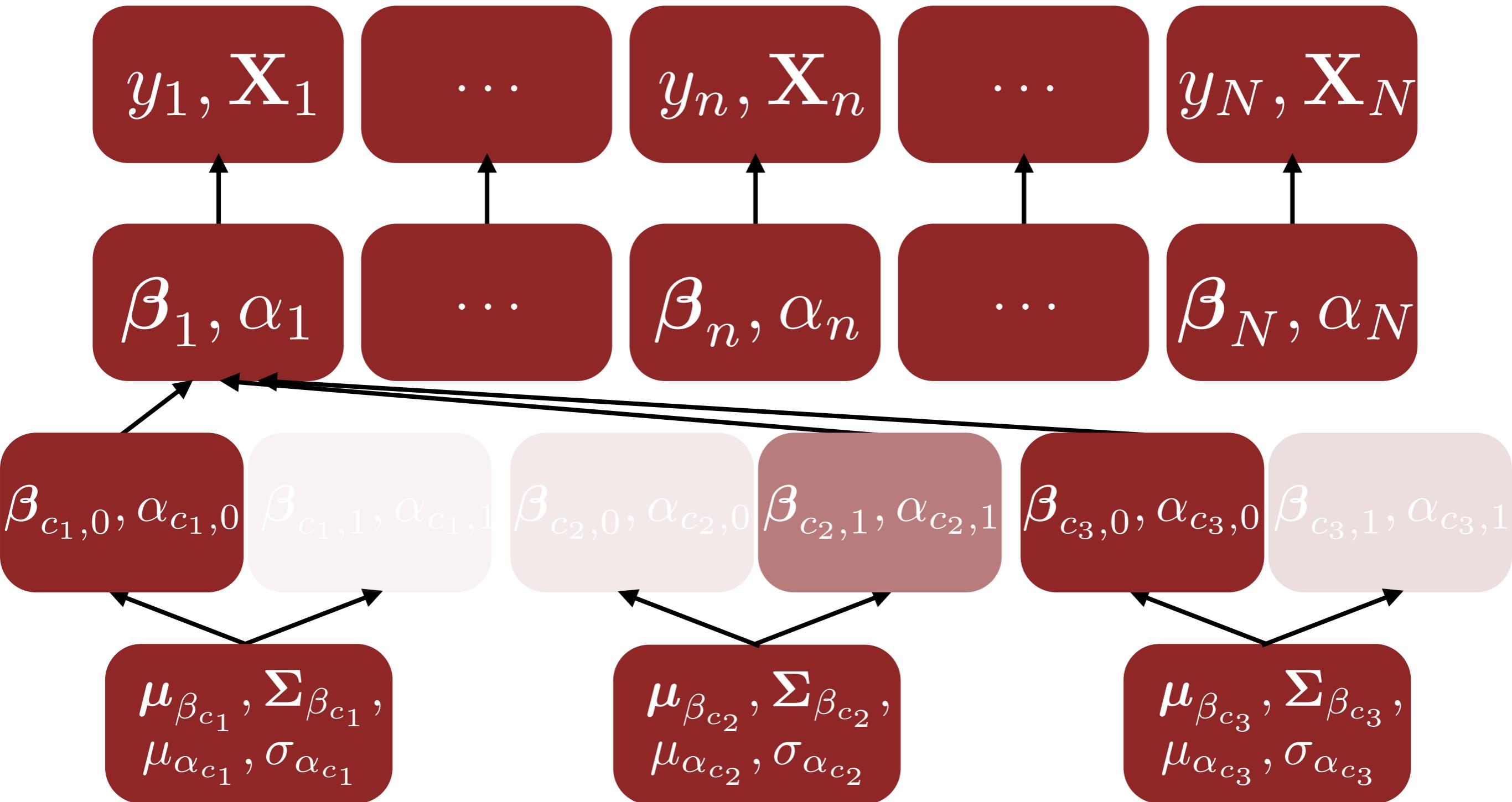
$$\alpha_n = \sum_{k=1}^{N_{\text{groups}}} \alpha_{k,j(n)}$$

$$\alpha_{k,j} \sim \mathcal{N}(\mu_{\alpha_k}, \sigma_{\alpha_k})$$

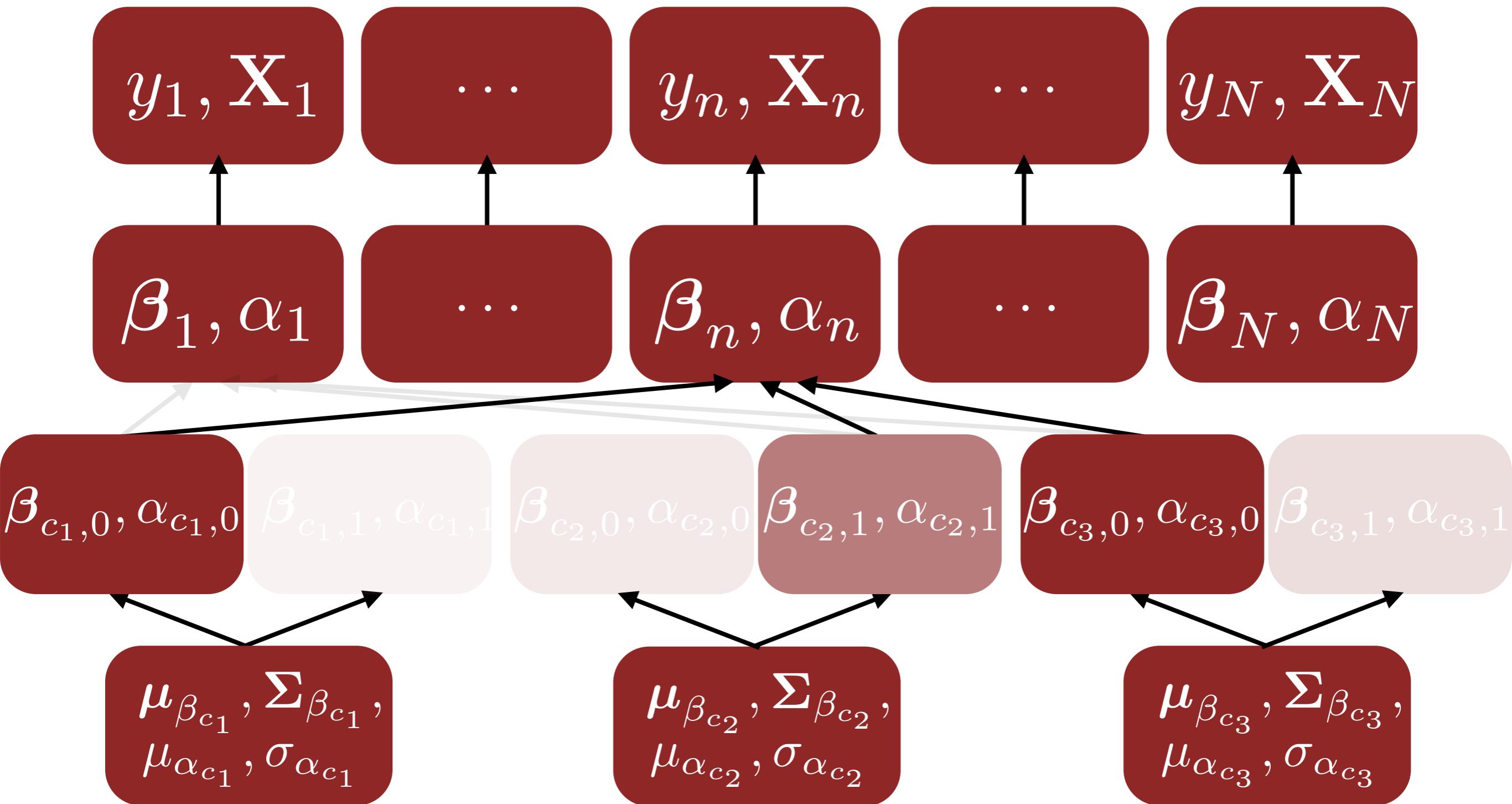
In a *multilevel* model the linear parameters vary across multiple hierarchies.



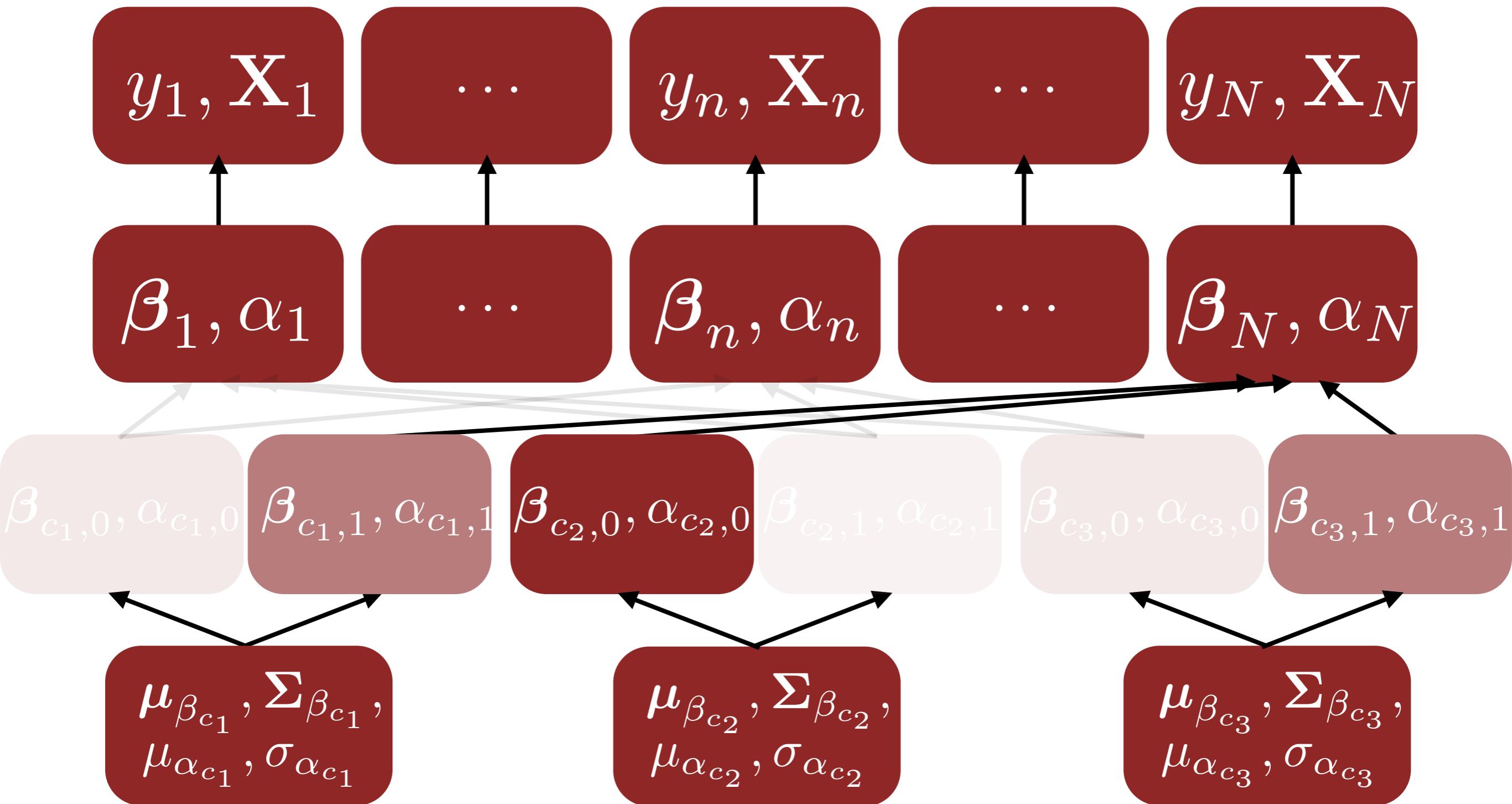
In a *multilevel* model the linear parameters vary across multiple hierarchies.



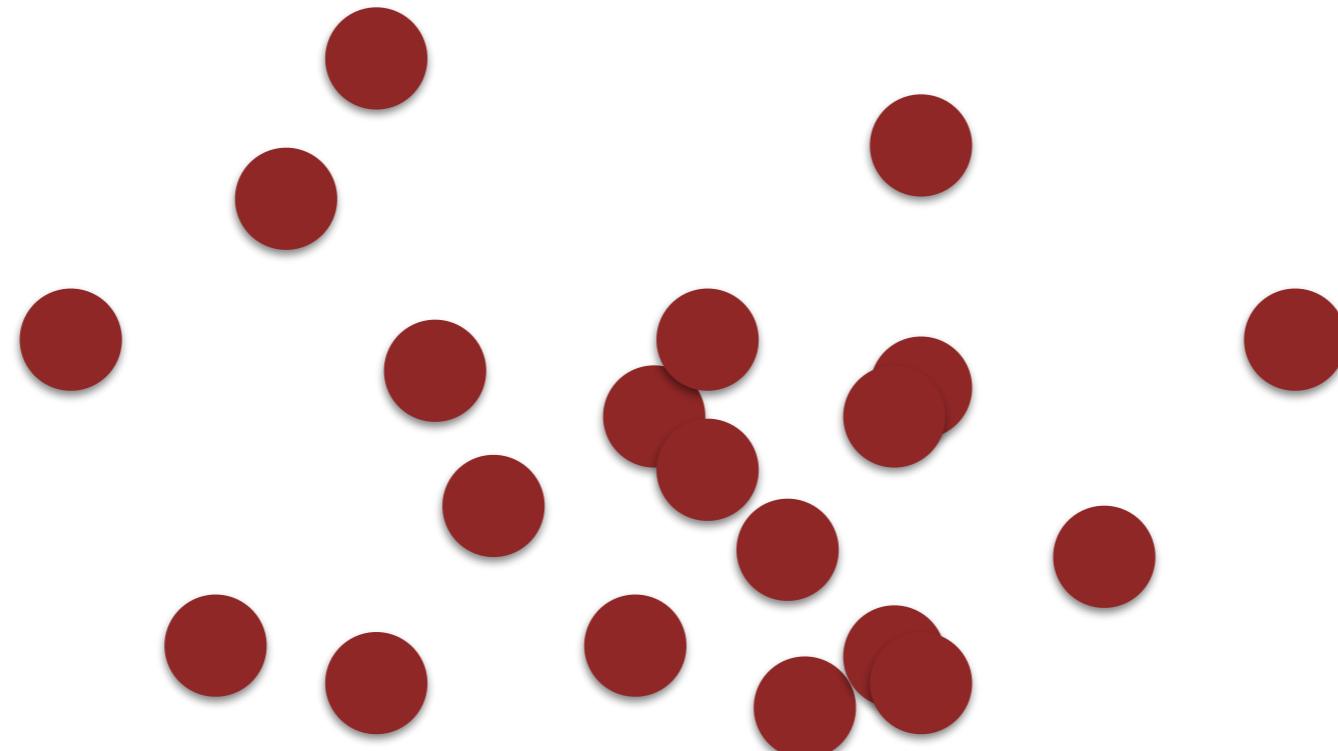
In a *multilevel* model the linear parameters vary across multiple hierarchies.



In a *multilevel* model the linear parameters vary across multiple hierarchies.

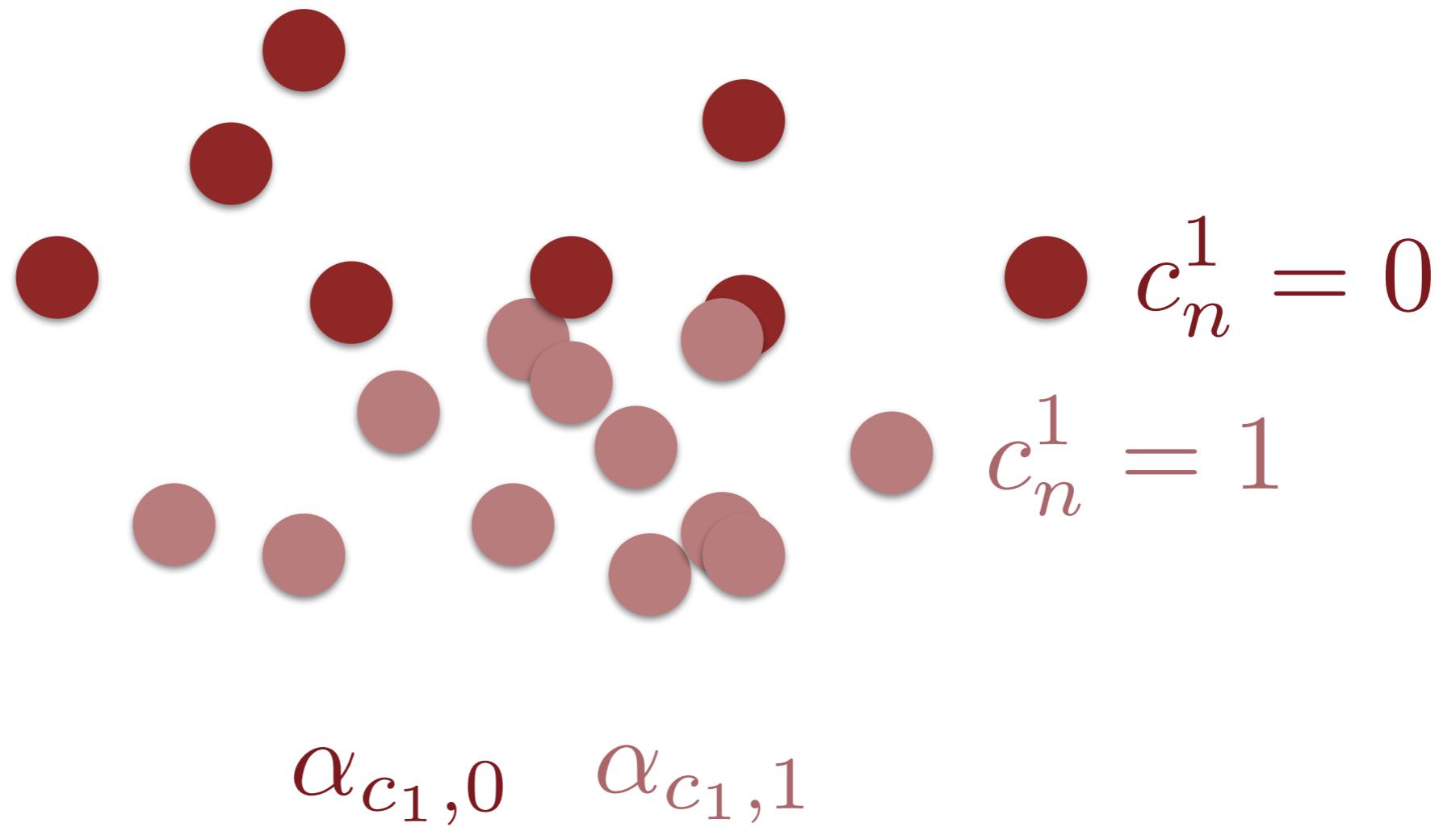


Each level contributes its own slope and intercept.



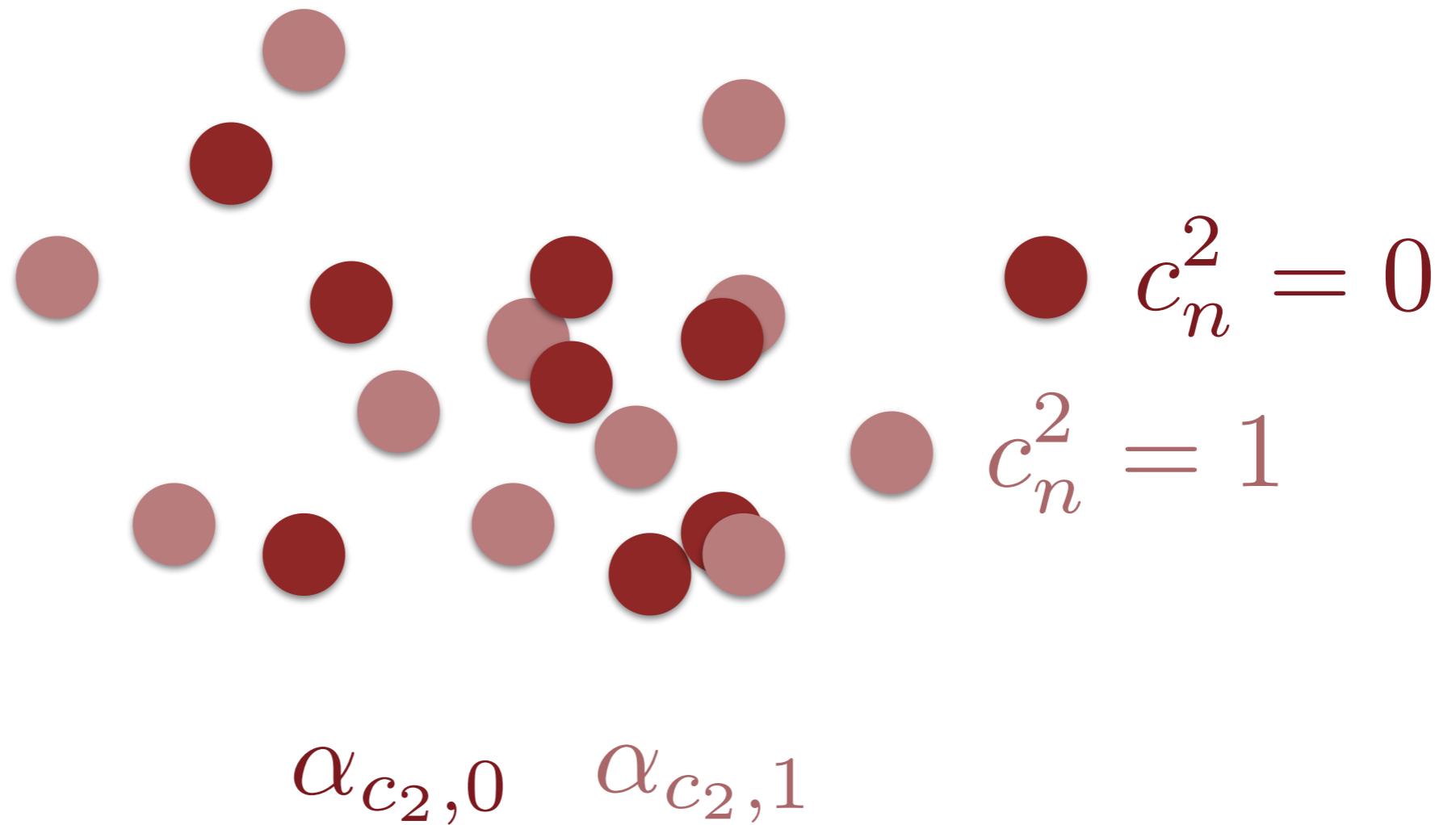
$$\{c_n^1, c_n^2, c_n^3\} \in \{0, 1\}$$

Each level contributes its own slope and intercept.

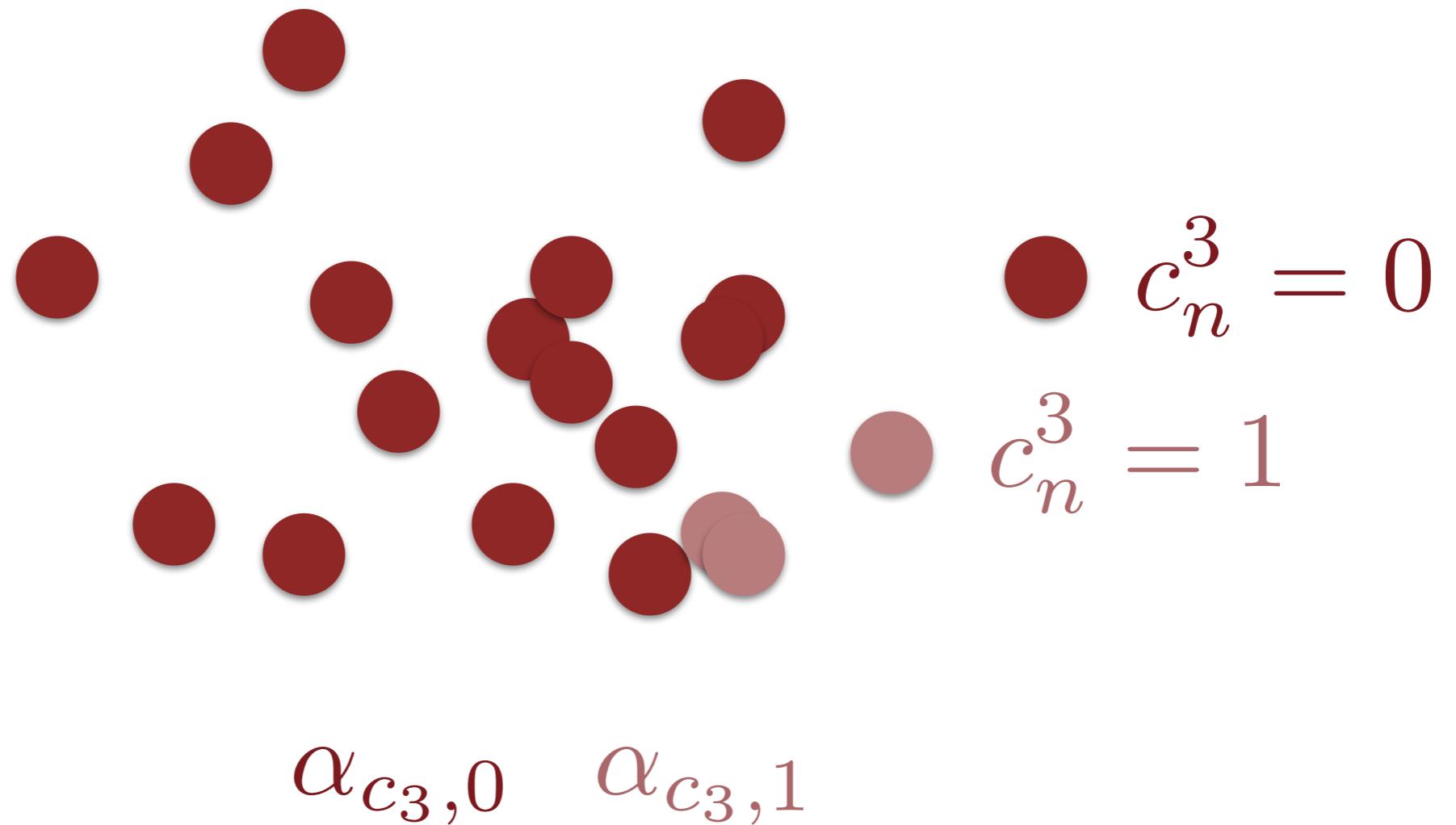


$$\alpha_{c_1,j} \sim \mathcal{N}(\mu_{\alpha_{c_1}}, \sigma_{\alpha_{c_1}})$$

Each level contributes its own slope and intercept.



Each level contributes its own slope and intercept.



$$\alpha_{c_3,j} \sim \mathcal{N}(\mu_{\alpha_{c_3}}, \sigma_{\alpha_{c_3}})$$

Each level contributes its own slope and intercept.



$$c_n^1 = 0$$



$$c_n^2 = 1$$



$$c_n^3 = 0$$

$$\alpha_n =$$

Each level contributes its own slope and intercept.



$$c_n^1 = 0$$



$$c_n^2 = 1$$



$$c_n^3 = 0$$

$$\alpha_n = \alpha_{c_1, 0}$$

Each level contributes its own slope and intercept.



$$c_n^1 = 0$$



$$c_n^2 = 1$$



$$c_n^3 = 0$$

$$\alpha_n = \alpha_{c_1,0} + \alpha_{c_2,1}$$

Each level contributes its own slope and intercept.



$$c_n^1 = 0$$



$$c_n^2 = 1$$



$$c_n^3 = 0$$

$$\alpha_n = \alpha_{c_1,0} + \alpha_{c_2,1} + \alpha_{c_3,0}$$

# General Linear Hierarchical Models and Markov Chain Monte Carlo

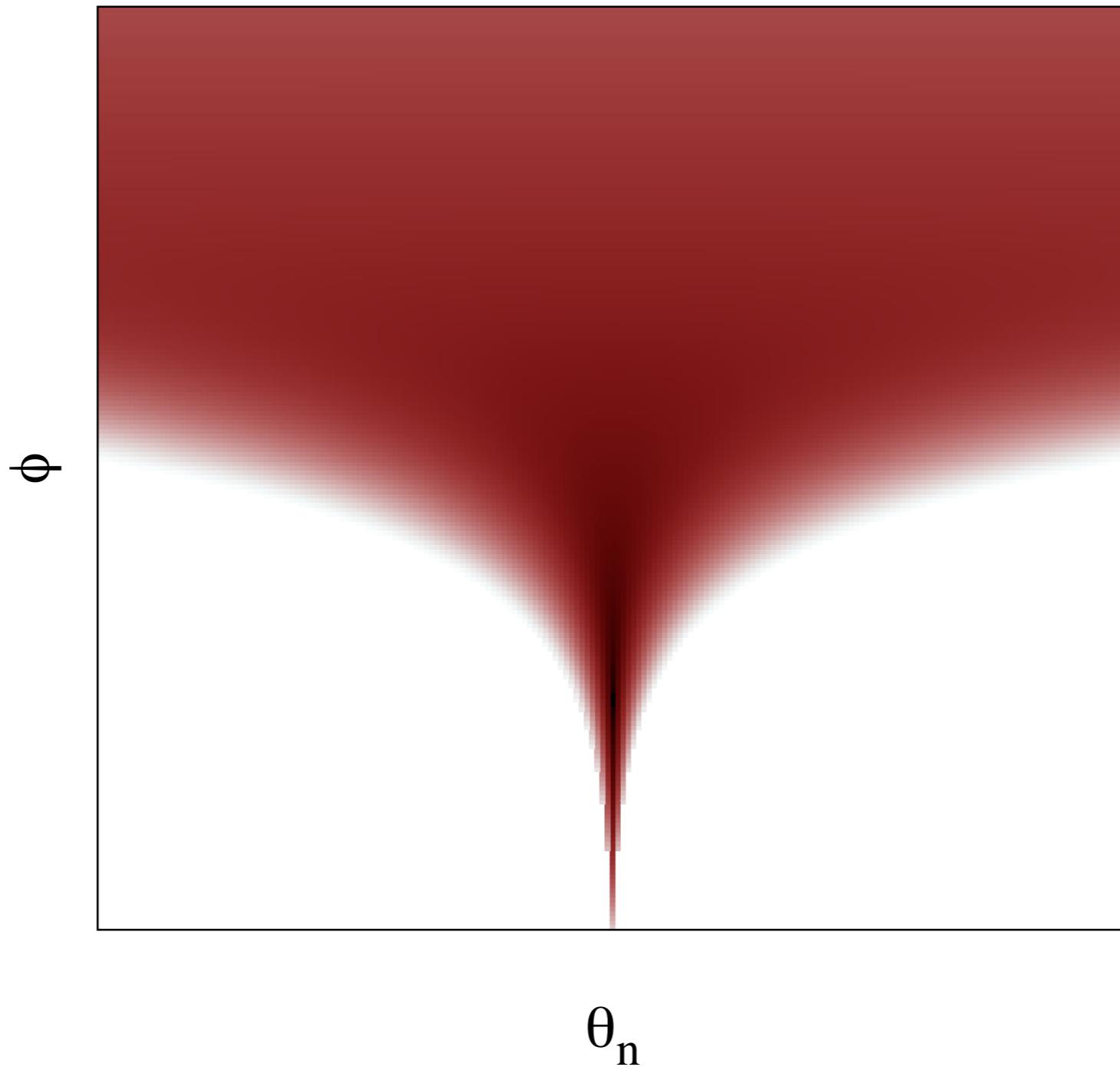
$\theta$

$\theta_n$

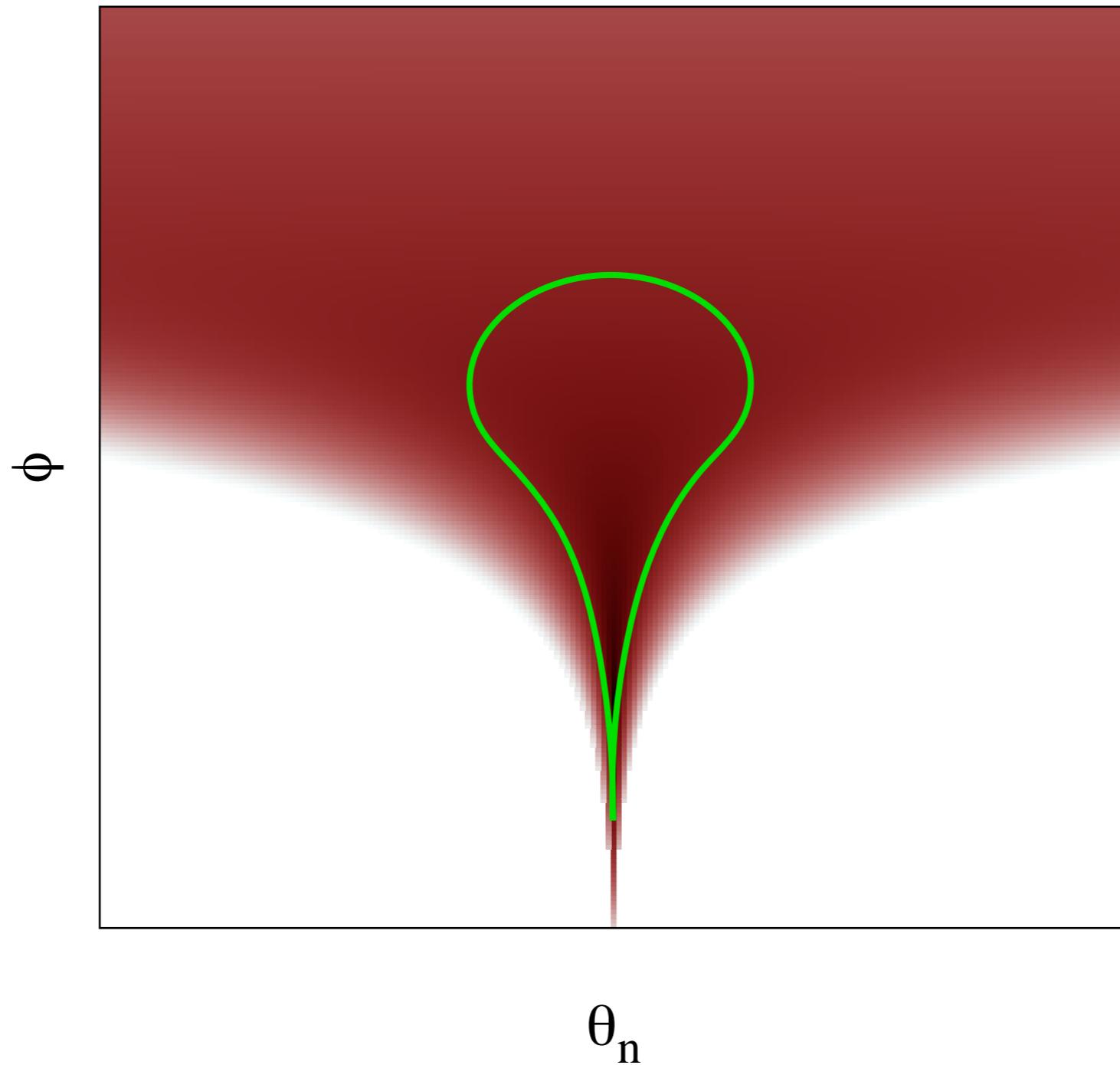
Neal's funnel distribution illustrates some of the problems with sampling from hierarchical distributions.

$$\pi(\theta|\phi) = \prod_{n=1}^N \mathcal{N}\left(\theta_n|0, e^{\frac{\phi}{2}}\right) \mathcal{N}(\phi|0, 3)$$

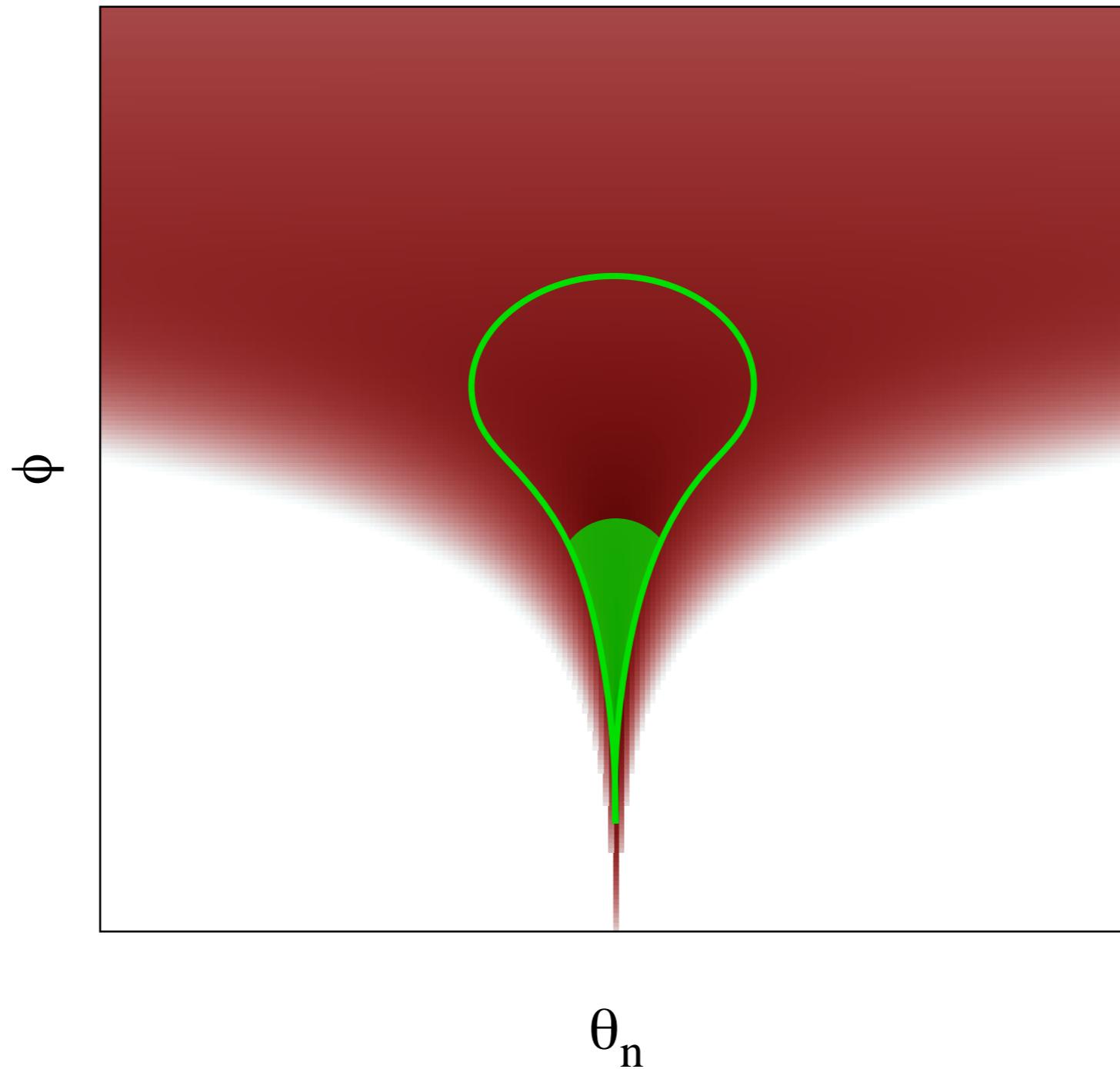
Neal's funnel distribution illustrates some of the problems with sampling from hierarchical priors.



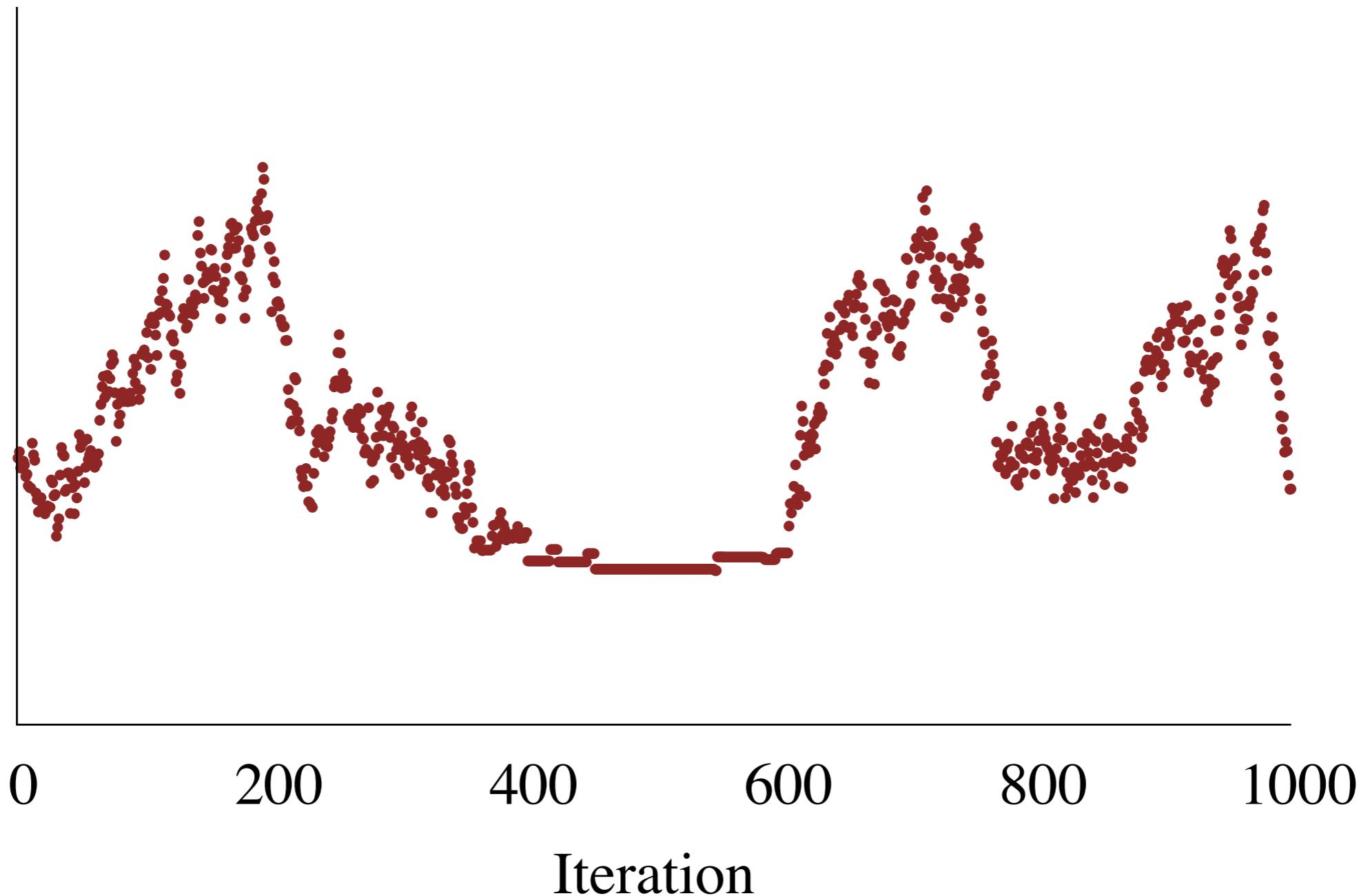
Neal's funnel distribution illustrates some of the problems with sampling from hierarchical priors.



Neal's funnel distribution illustrates some of the problems with sampling from hierarchical priors.



Neal's funnel distribution illustrates some of the problems with sampling from hierarchical priors.



The geometry of the funnel, however, dramatically improves when we sample from auxiliary parameters.

$$\pi(\theta|\phi) = \prod_{n=1}^N \mathcal{N}\left(\theta_n|0, e^{\frac{\phi}{2}}\right) \mathcal{N}(\phi|0, 3)$$

The geometry of the funnel, however, dramatically improves when we sample from auxiliary parameters.

$$\pi(\theta|\phi) = \prod_{n=1}^N \mathcal{N}\left(\theta_n|0, e^{\frac{\phi}{2}}\right) \mathcal{N}(\phi|0, 3)$$

$$\tilde{\theta}_n \sim \mathcal{N}(0, 1)$$

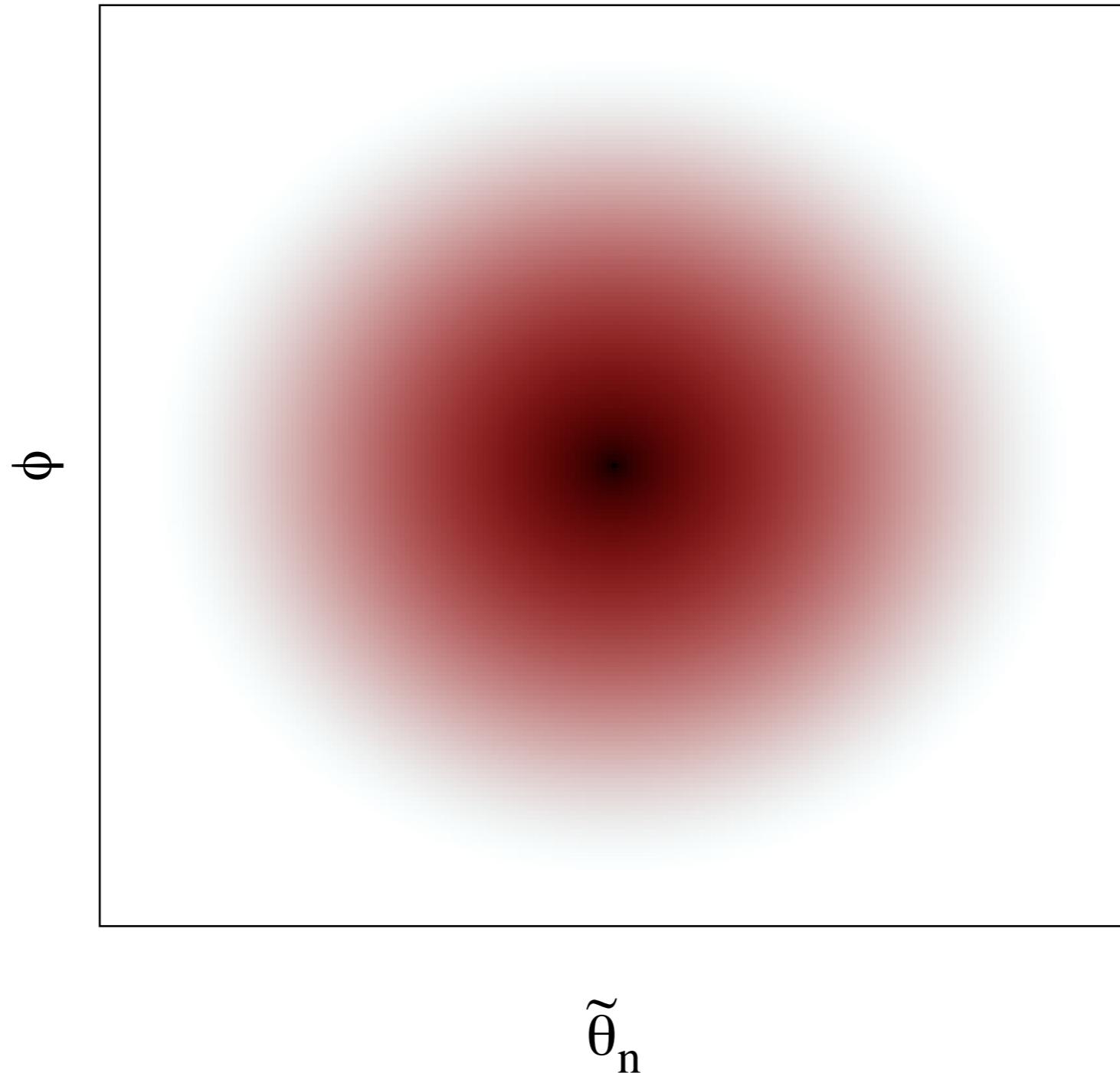
The geometry of the funnel, however, dramatically improves when we sample from auxiliary parameters.

$$\pi(\theta|\phi) = \prod_{n=1}^N \mathcal{N}\left(\theta_n|0, e^{\frac{\phi}{2}}\right) \mathcal{N}(\phi|0, 3)$$

$$\tilde{\theta}_n \sim \mathcal{N}(0, 1)$$

$$\theta_n = 0 + \tilde{\theta}_n e^{\frac{\phi}{2}}$$

The geometry of the funnel, however, dramatically improves when we sample from auxiliary parameters.



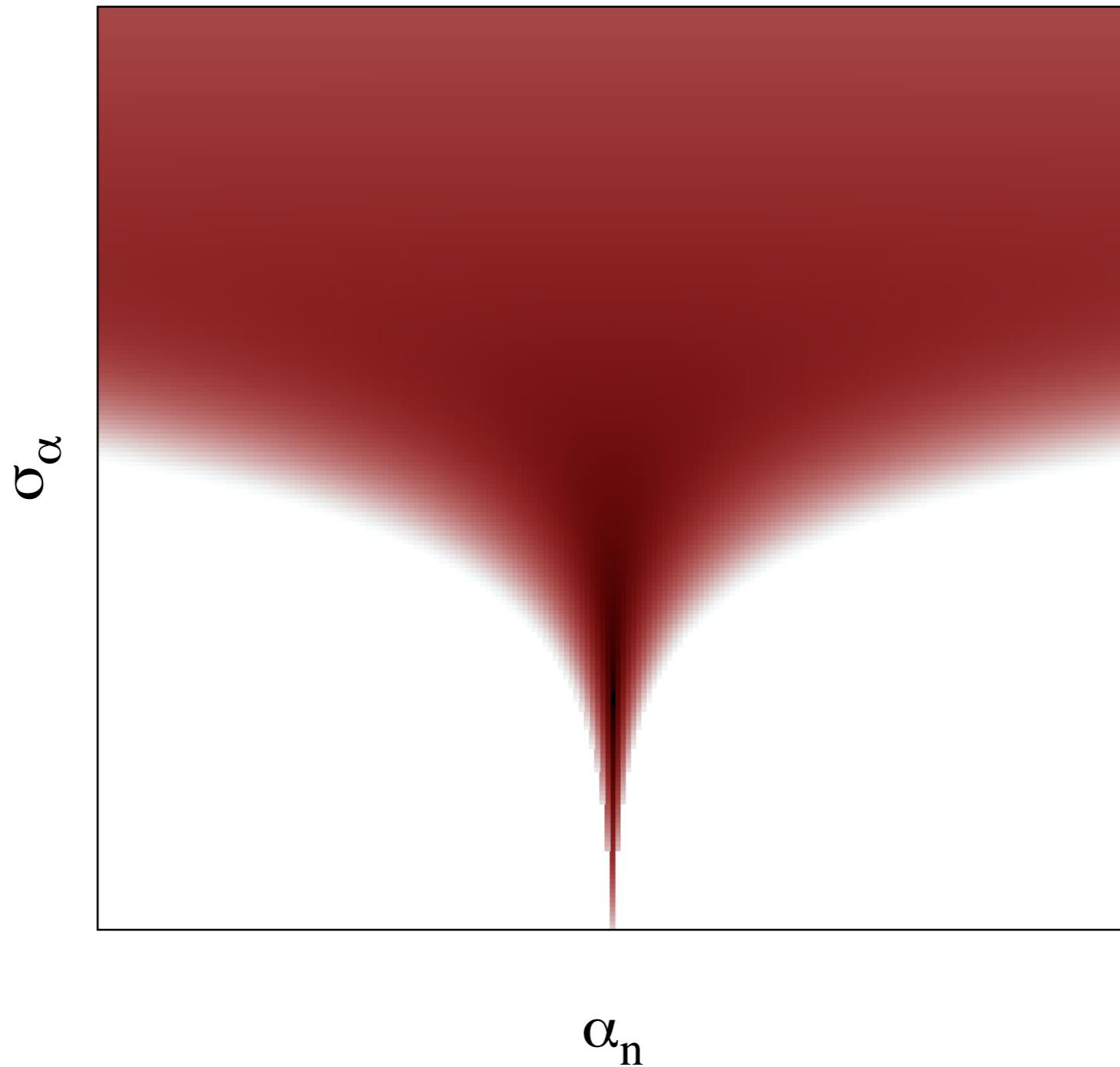
The most natural representation of a general linear hierarchical model is a *centered parameterization*.

$$\pi(y_n | g(\mathbf{X}_n^T \boldsymbol{\beta}_n + \alpha_n), \theta)$$

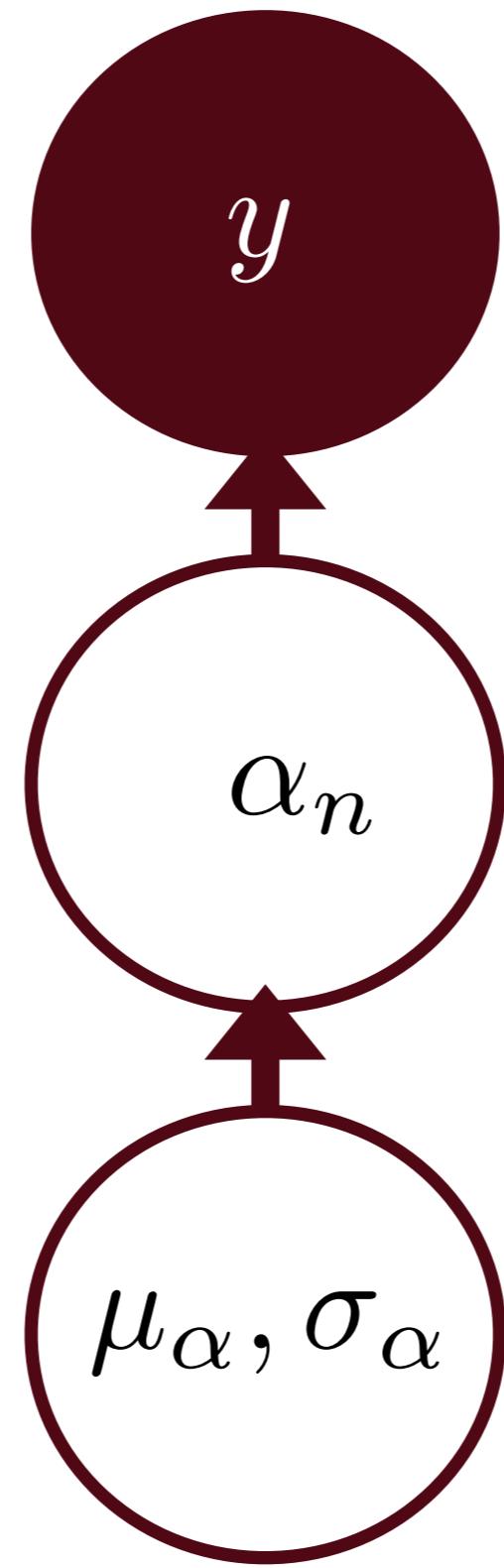
$$\boldsymbol{\beta}_n \sim \mathcal{N}(\boldsymbol{\mu}_{\beta}, \boldsymbol{\Sigma}_{\beta})$$

$$\alpha_n \sim \mathcal{N}(\mu_{\alpha}, \sigma_{\alpha})$$

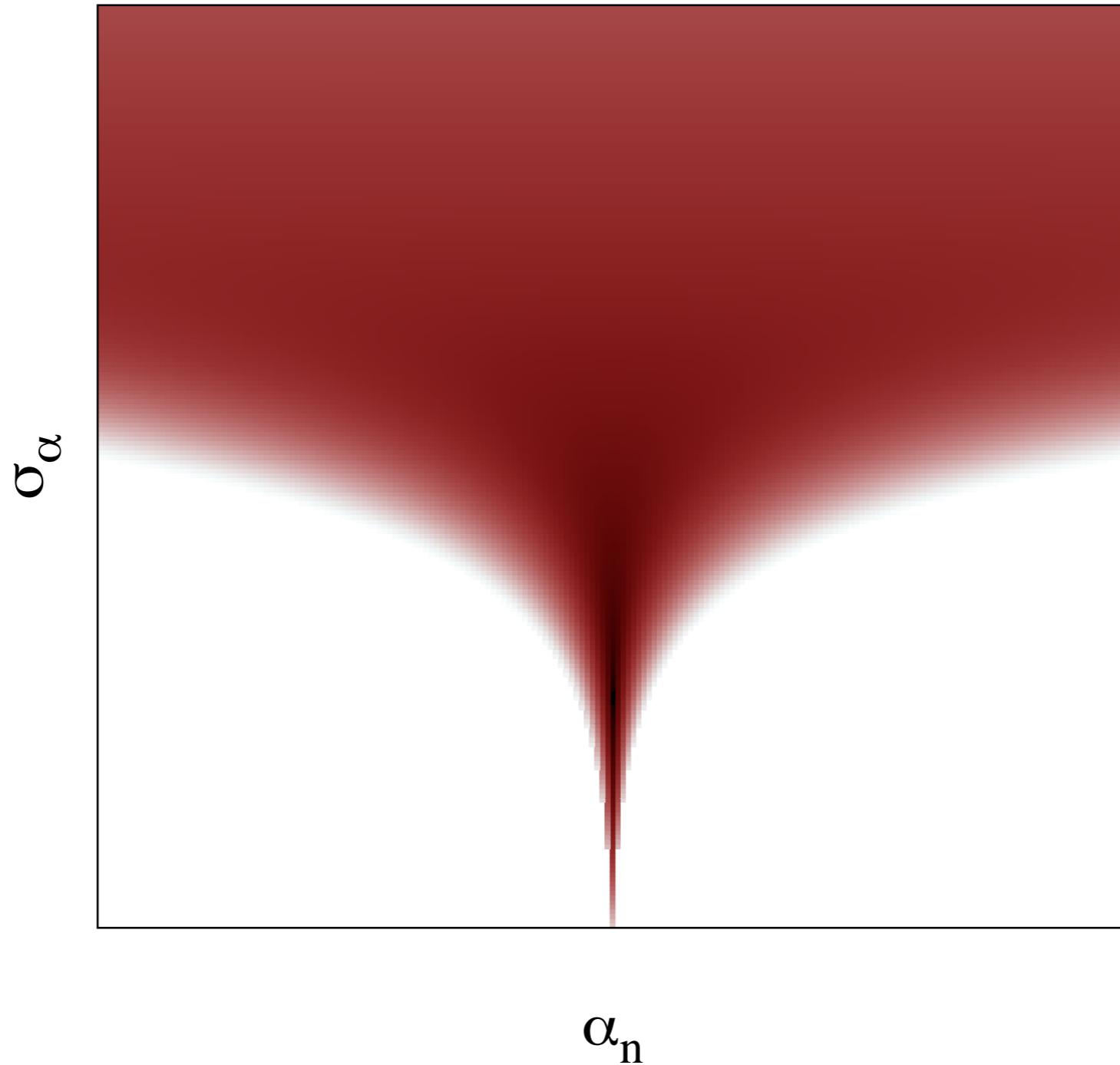
When data are sparse and uninformative, the centered parameterization of the posterior resembles the funnel.



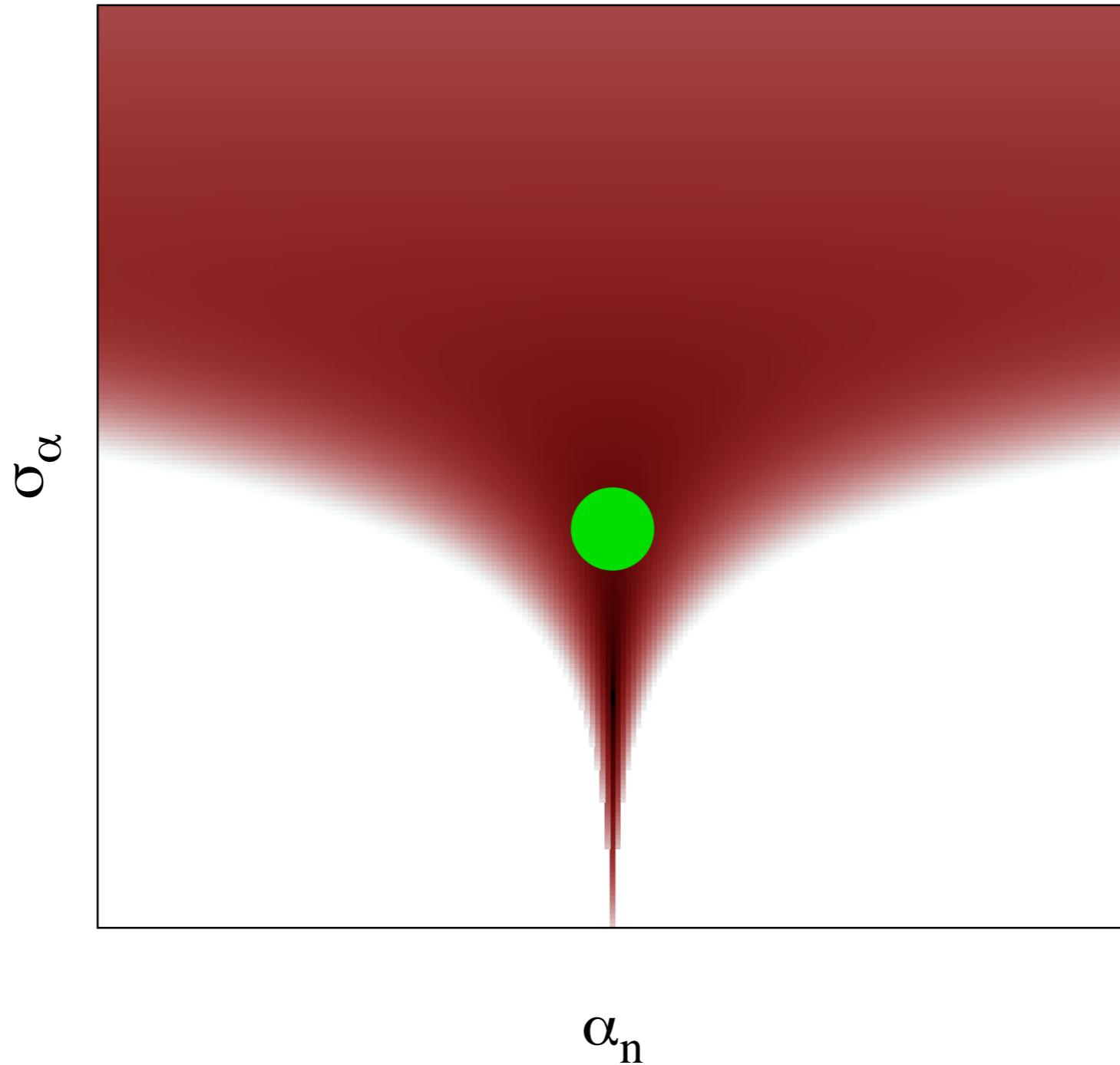
When data are sparse and uninformative, the centered parameterization of the posterior resembles the funnel.



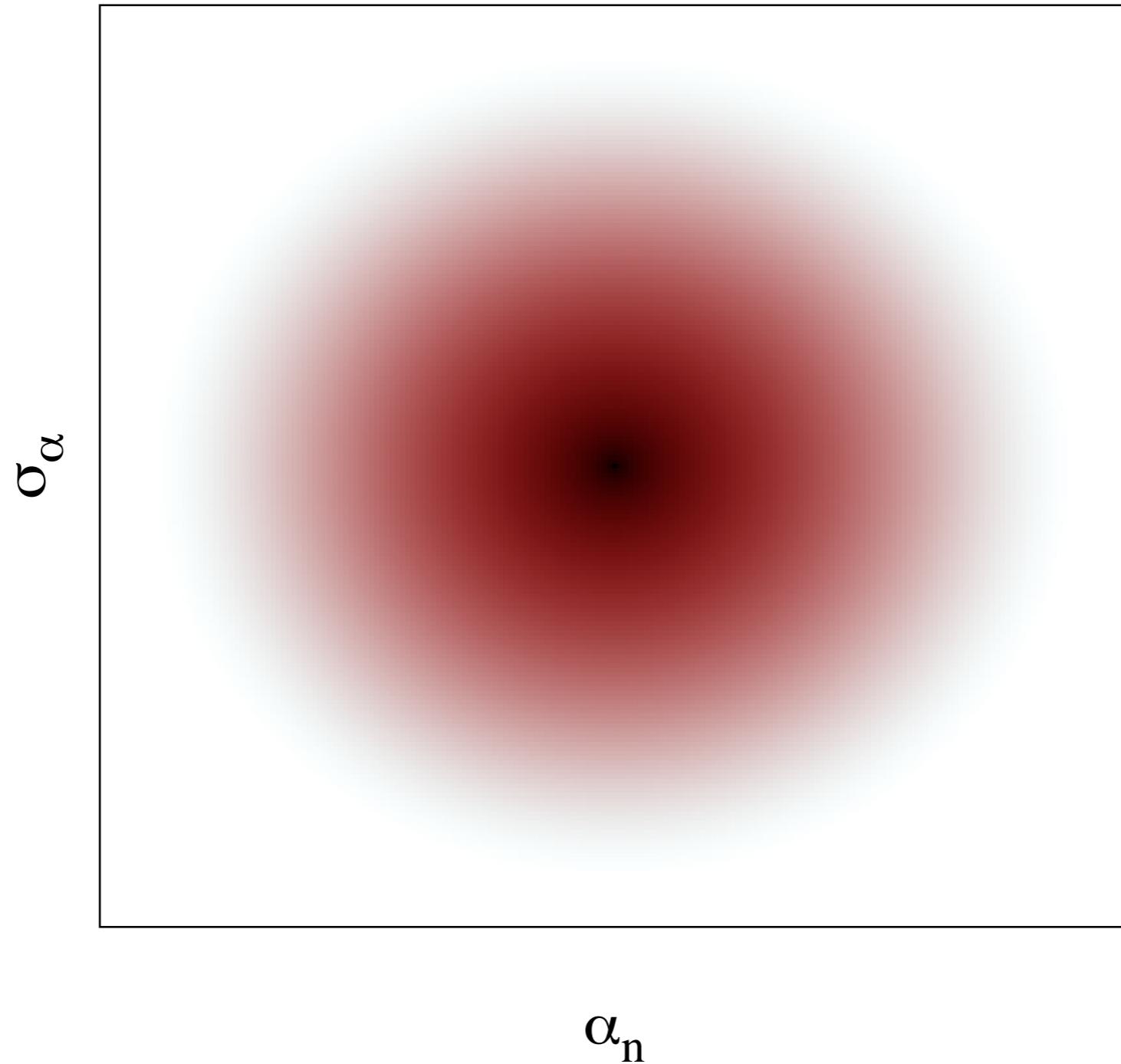
But when the data are more informative  
the posterior geometry is much nicer.



But when the data are more informative  
the posterior geometry is much nicer.



But when the data are more informative  
the posterior geometry is much nicer.



The same general linear hierarchical model can also be implemented with a *non-centered parameterization*.

$$\pi(y_n | g(\mathbf{X}_n^T \tilde{\boldsymbol{\beta}}_n + \alpha_n), \theta)$$

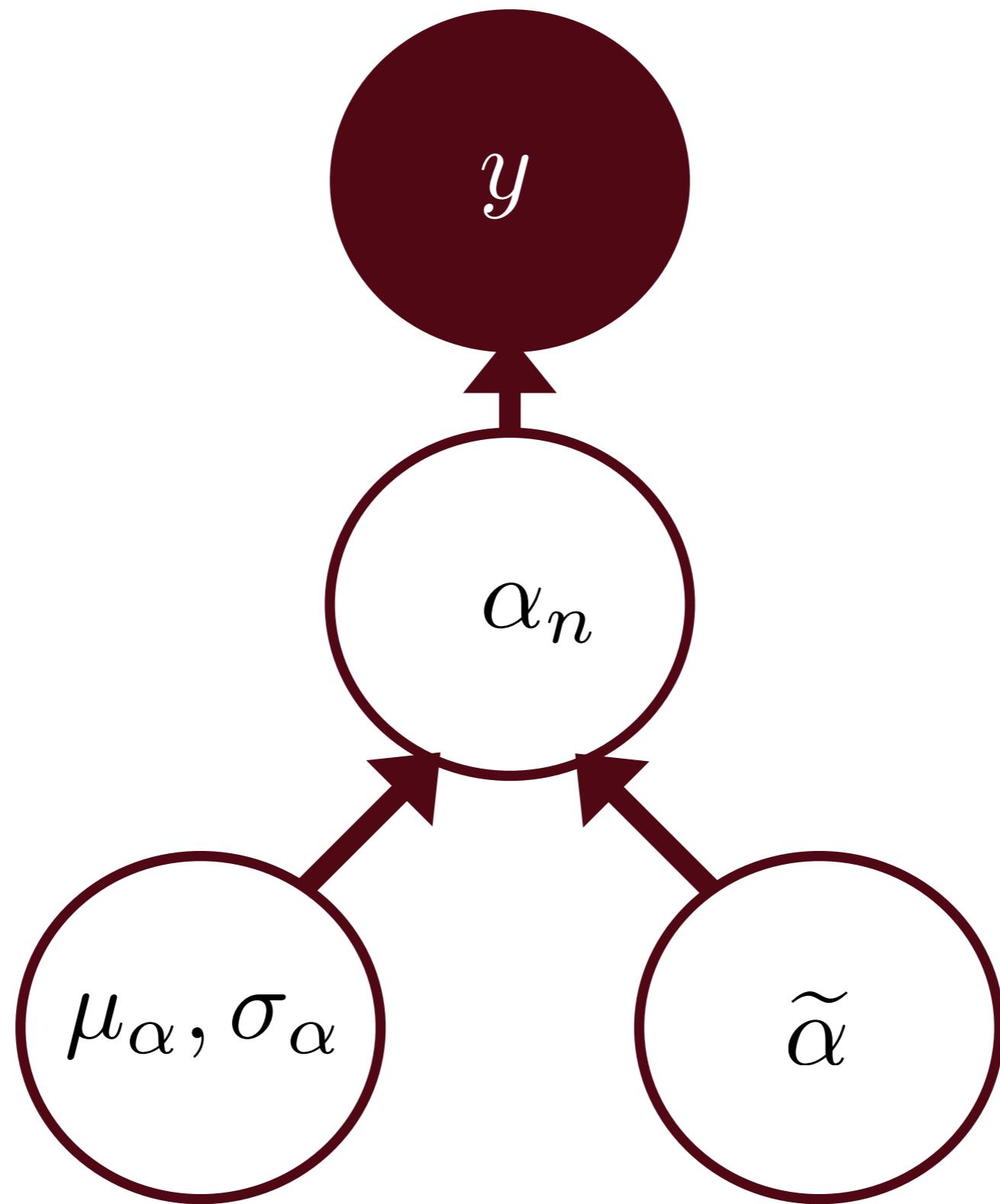
$$\tilde{\boldsymbol{\beta}}_n \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$\boldsymbol{\beta}_n = \boldsymbol{\mu}_{\beta} + \mathbf{L}\tilde{\boldsymbol{\beta}}_n, \quad \mathbf{L}\mathbf{L}^T = \boldsymbol{\Sigma}_{\beta}$$

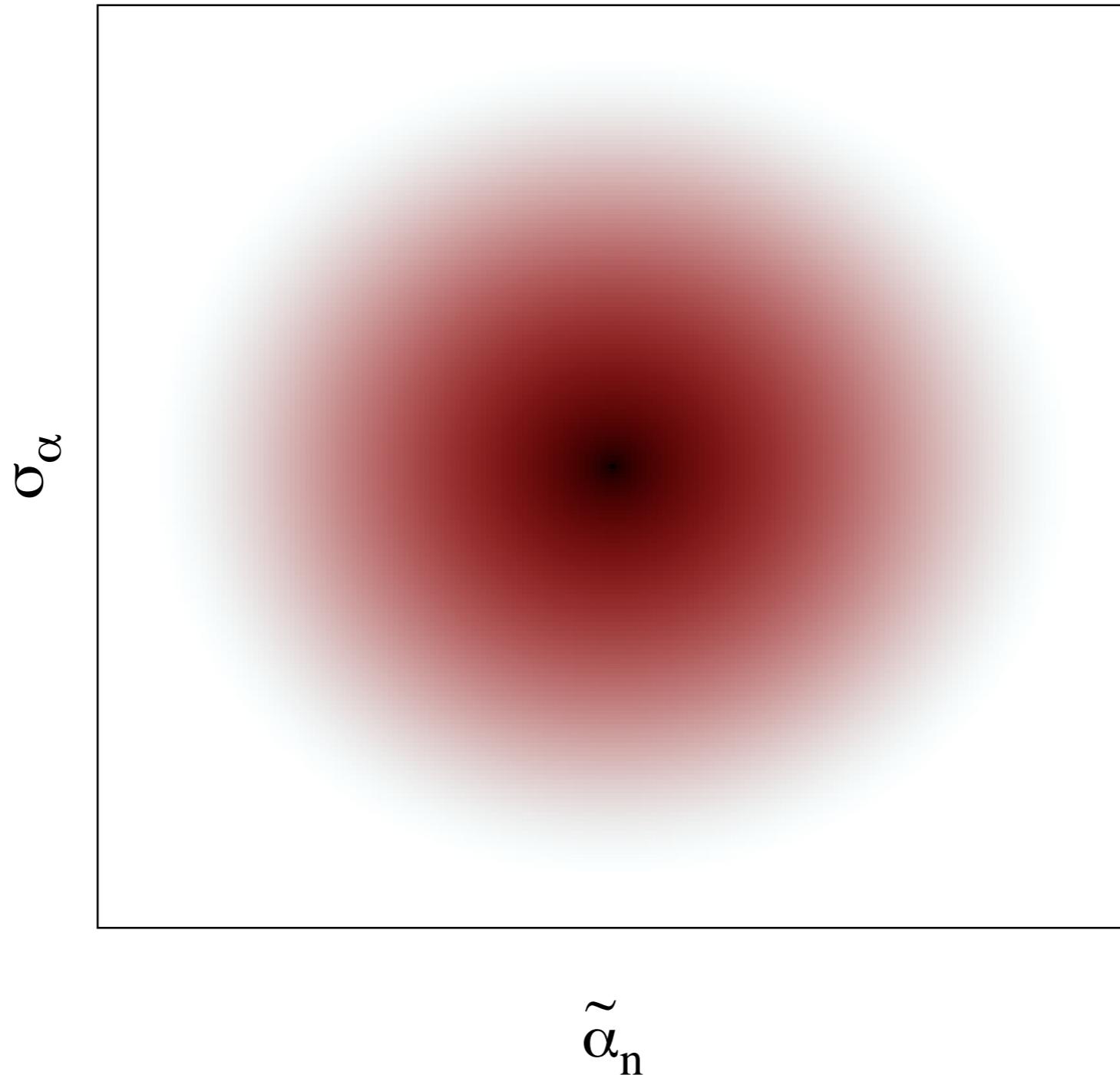
$$\tilde{\alpha}_n \sim \mathcal{N}(0, 1)$$

$$\alpha_n = \mu_{\alpha} + \sigma_{\alpha}\tilde{\alpha}_n$$

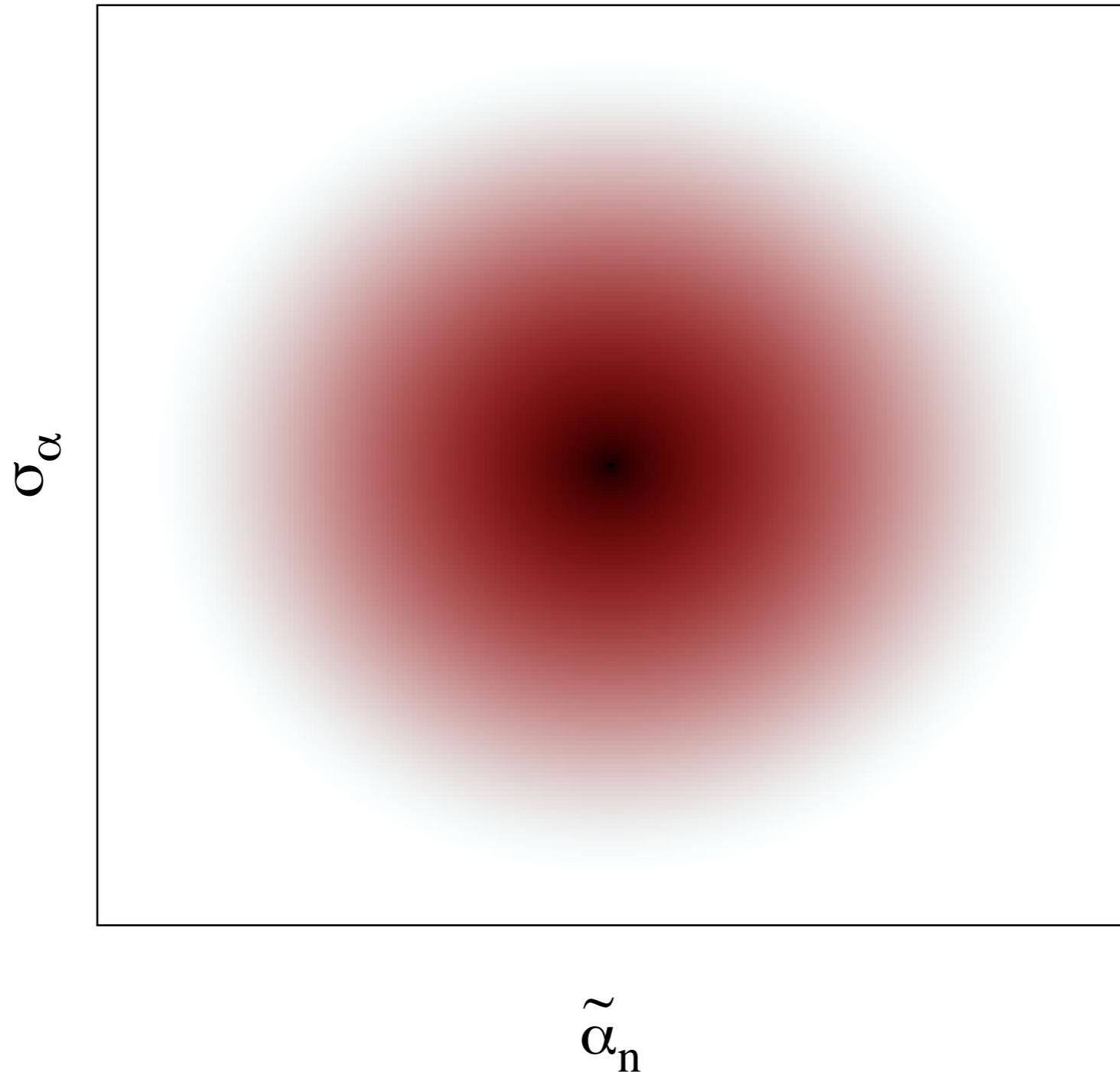
When data are sparse, the non-centered parameterization of the posterior yields a favorable geometry.



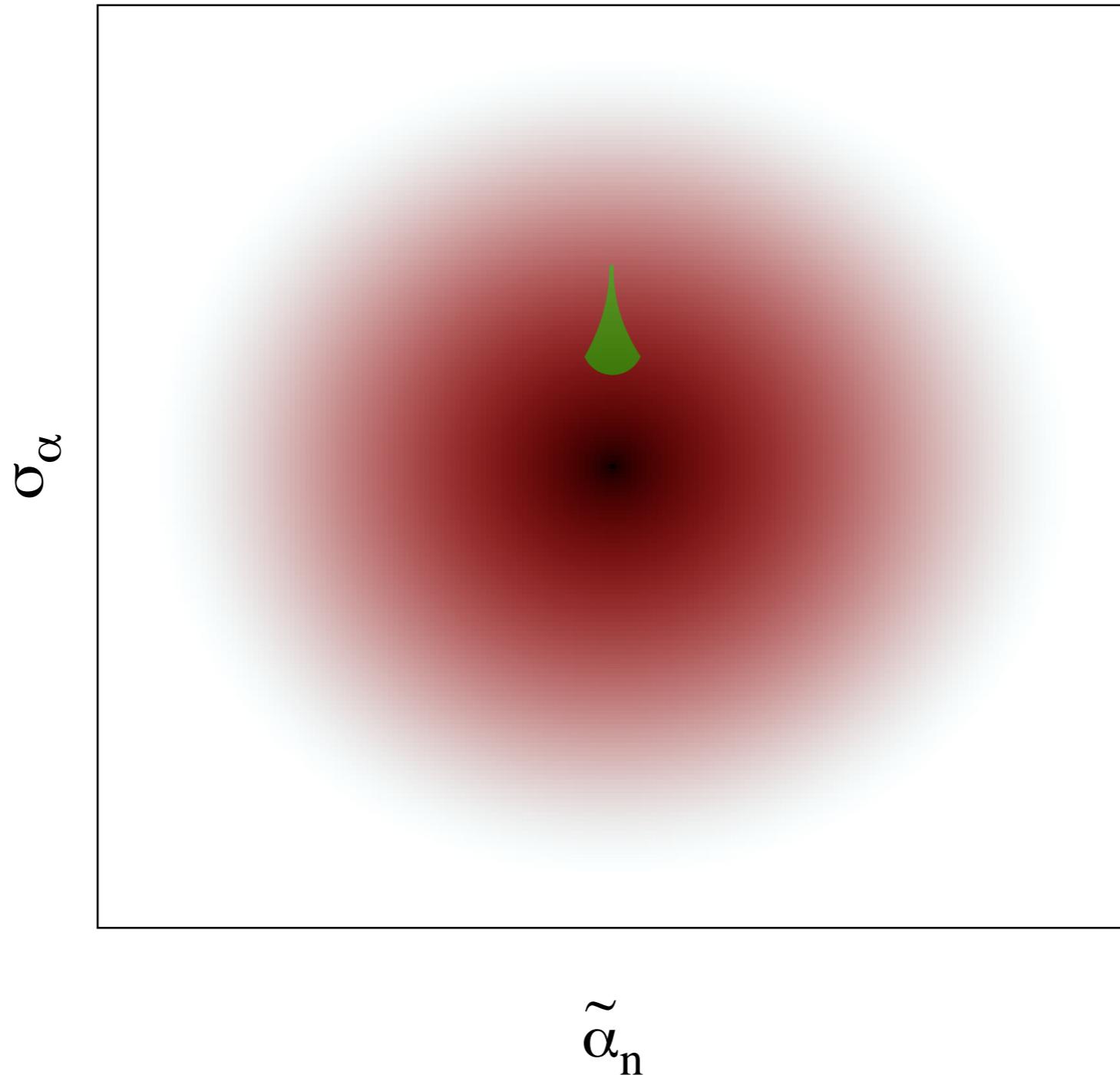
When data are sparse, the non-centered parameterization  
of the posterior yields a favorable geometry.



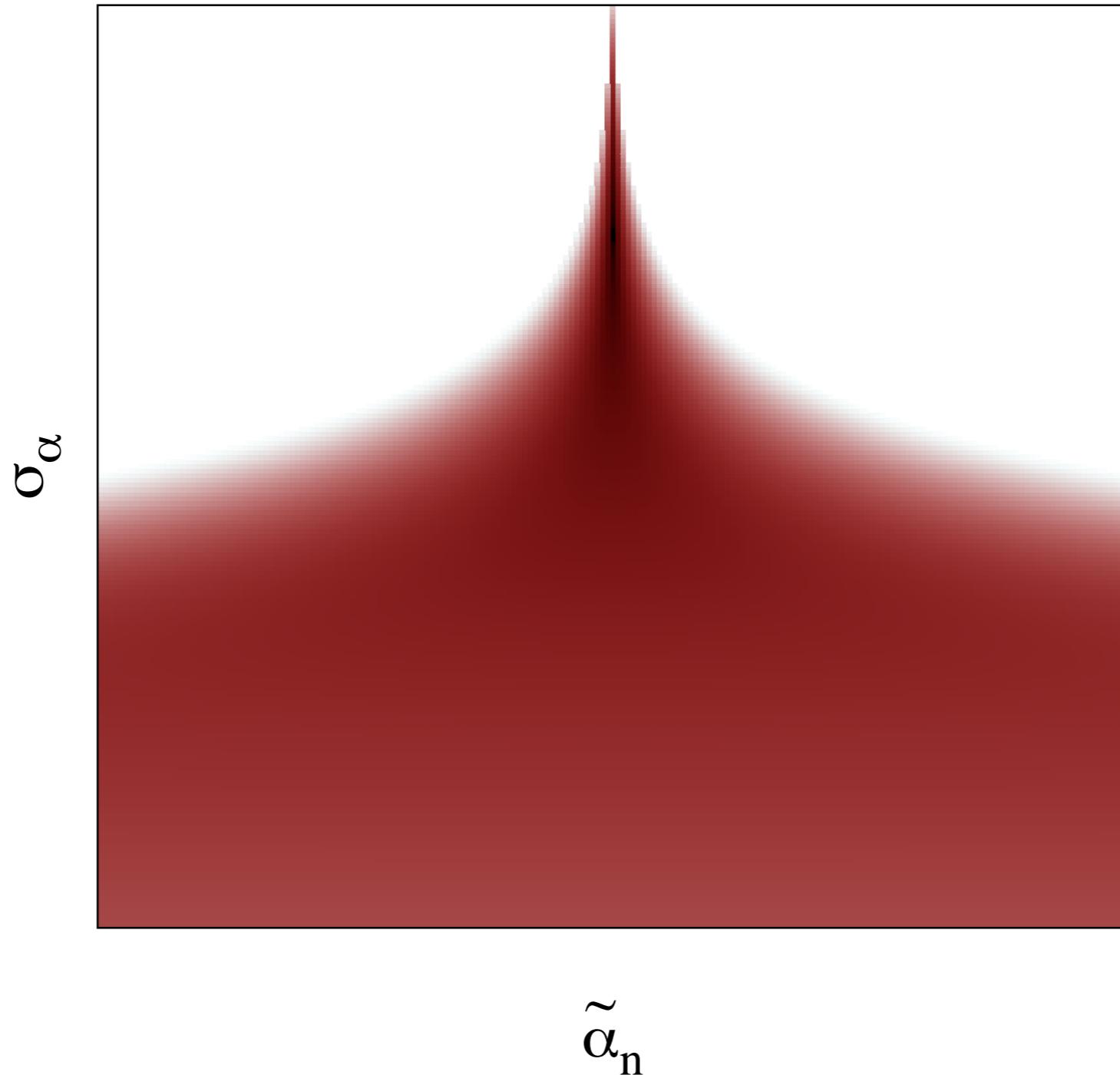
The constraints imposed by the data, however, are more complex, yielding pathological posteriors with more data.



The constraints imposed by the data, however, are more complex, yielding pathological posteriors with more data.



The constraints imposed by the data, however, are more complex, yielding pathological posteriors with more data.



Consequently the two parameterizations are each advantageous in different circumstances.

Centered  
Parameterization

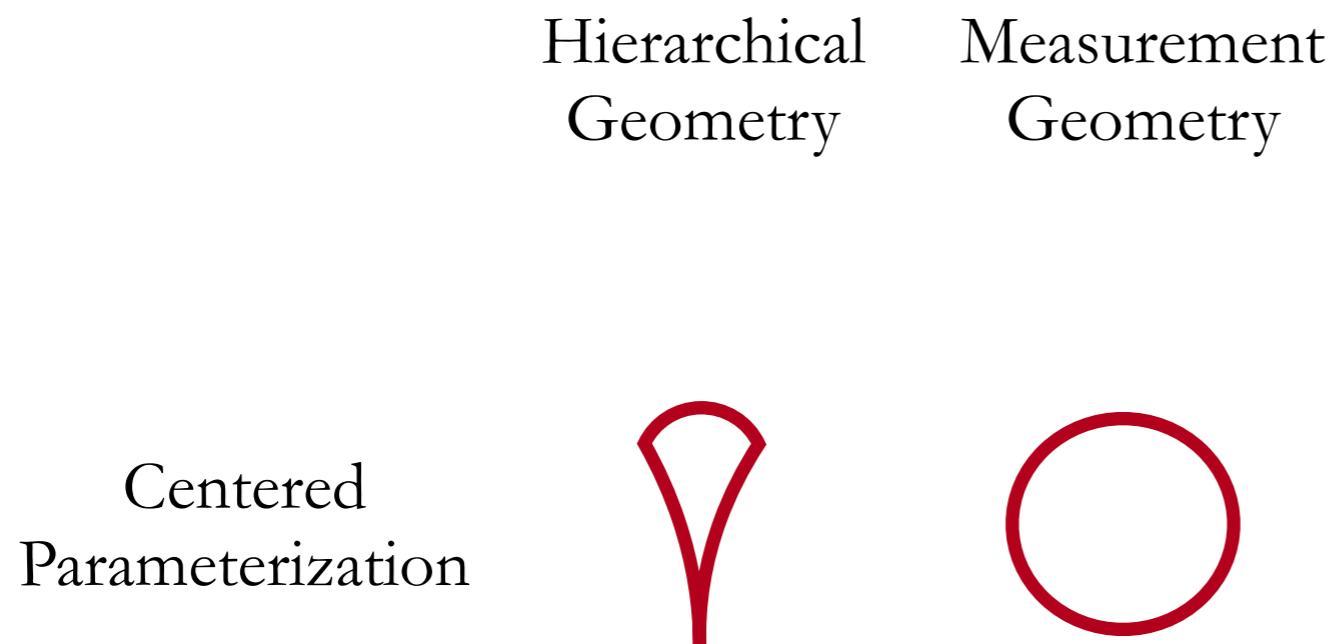
Consequently the two parameterizations are each advantageous in different circumstances.

Hierarchical  
Geometry

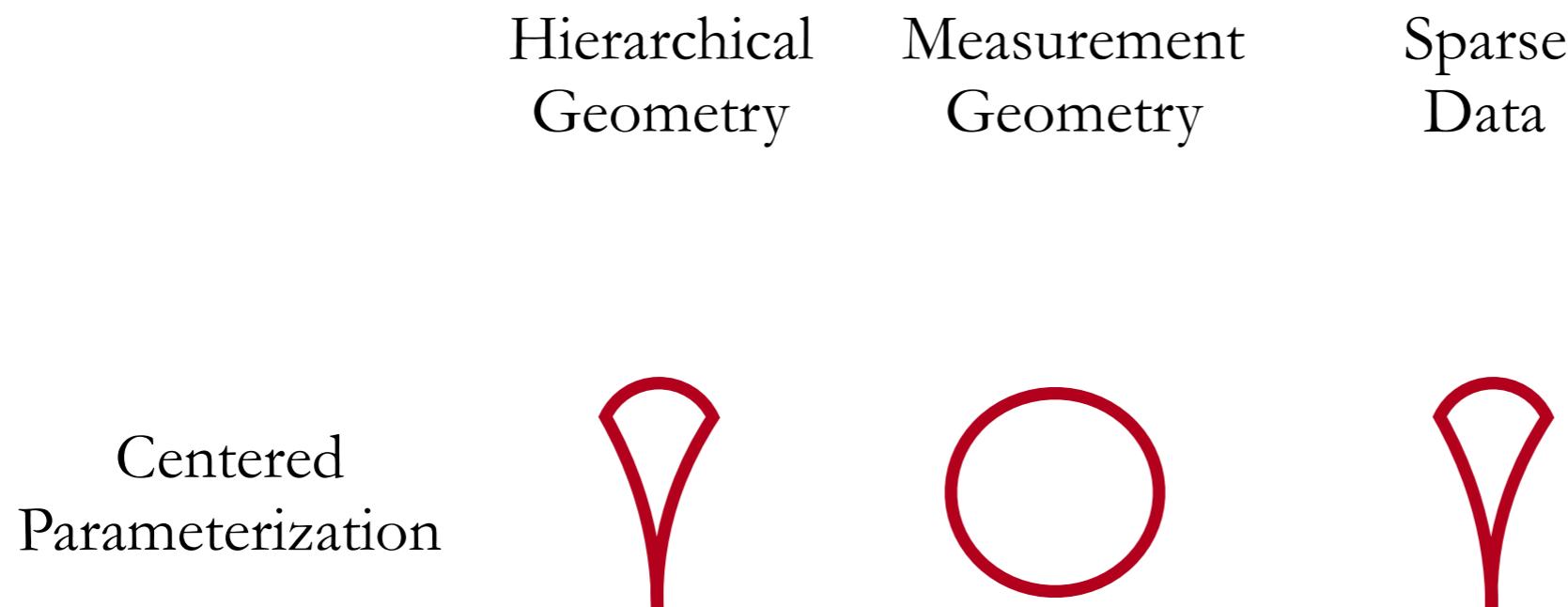
Centered  
Parameterization



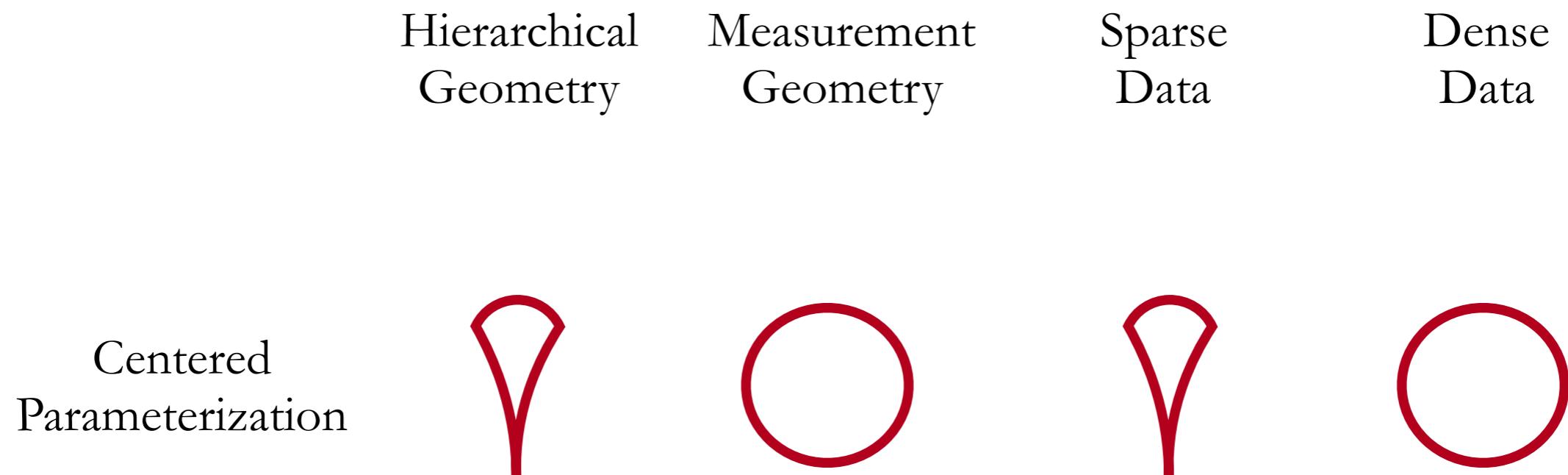
Consequently the two parameterizations are each advantageous in different circumstances.



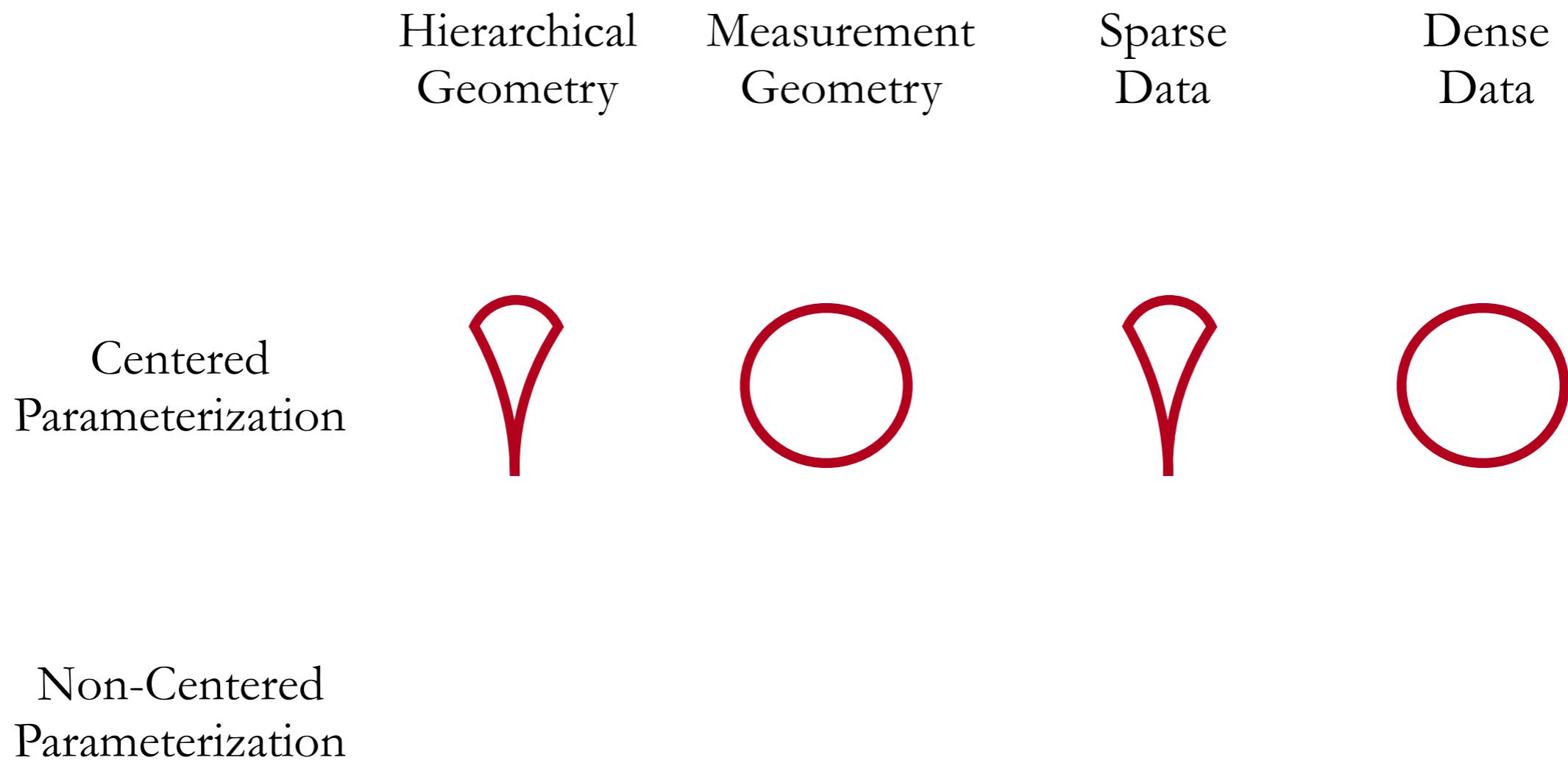
Consequently the two parameterizations are each advantageous in different circumstances.



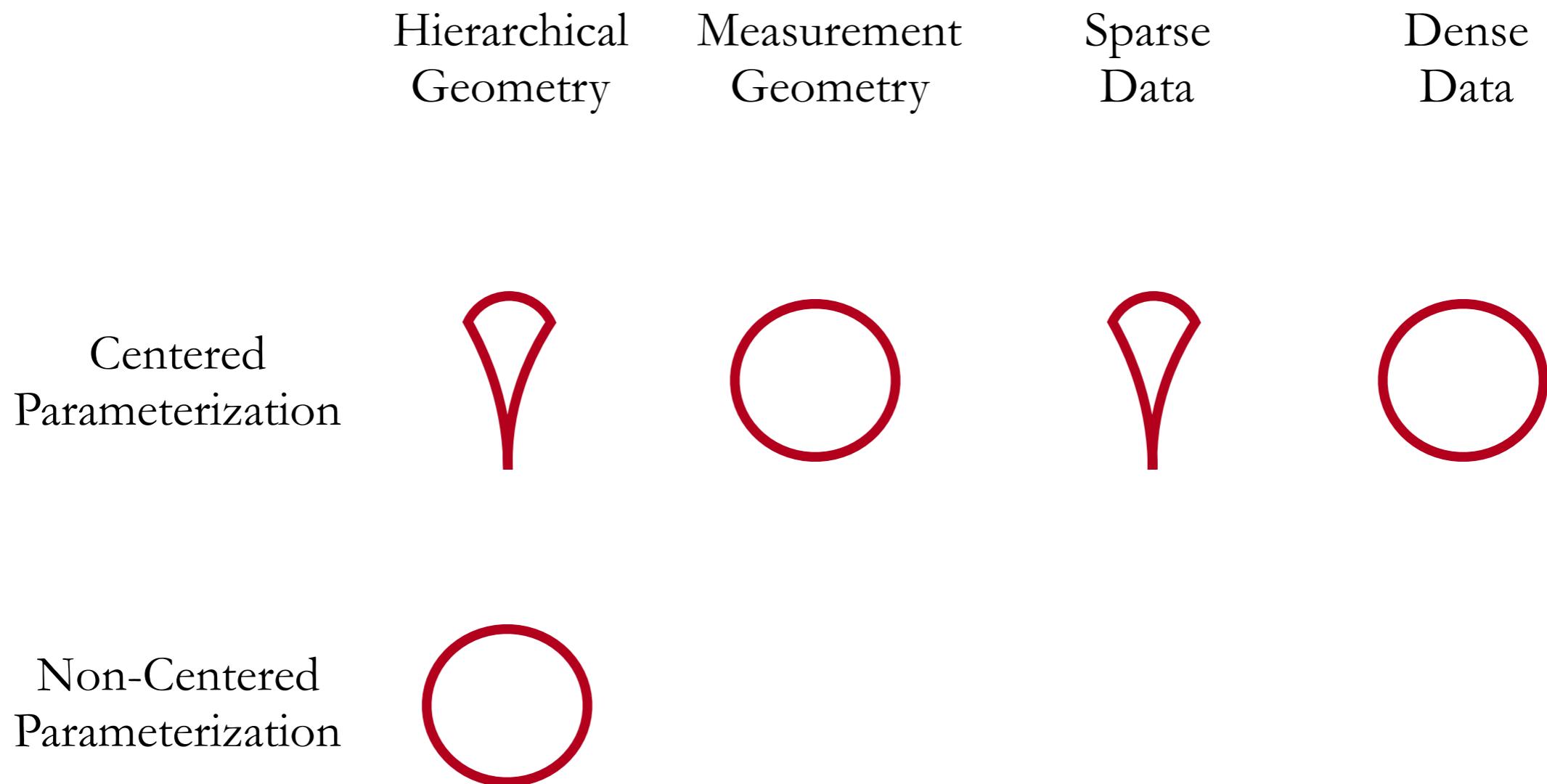
Consequently the two parameterizations are each advantageous in different circumstances.



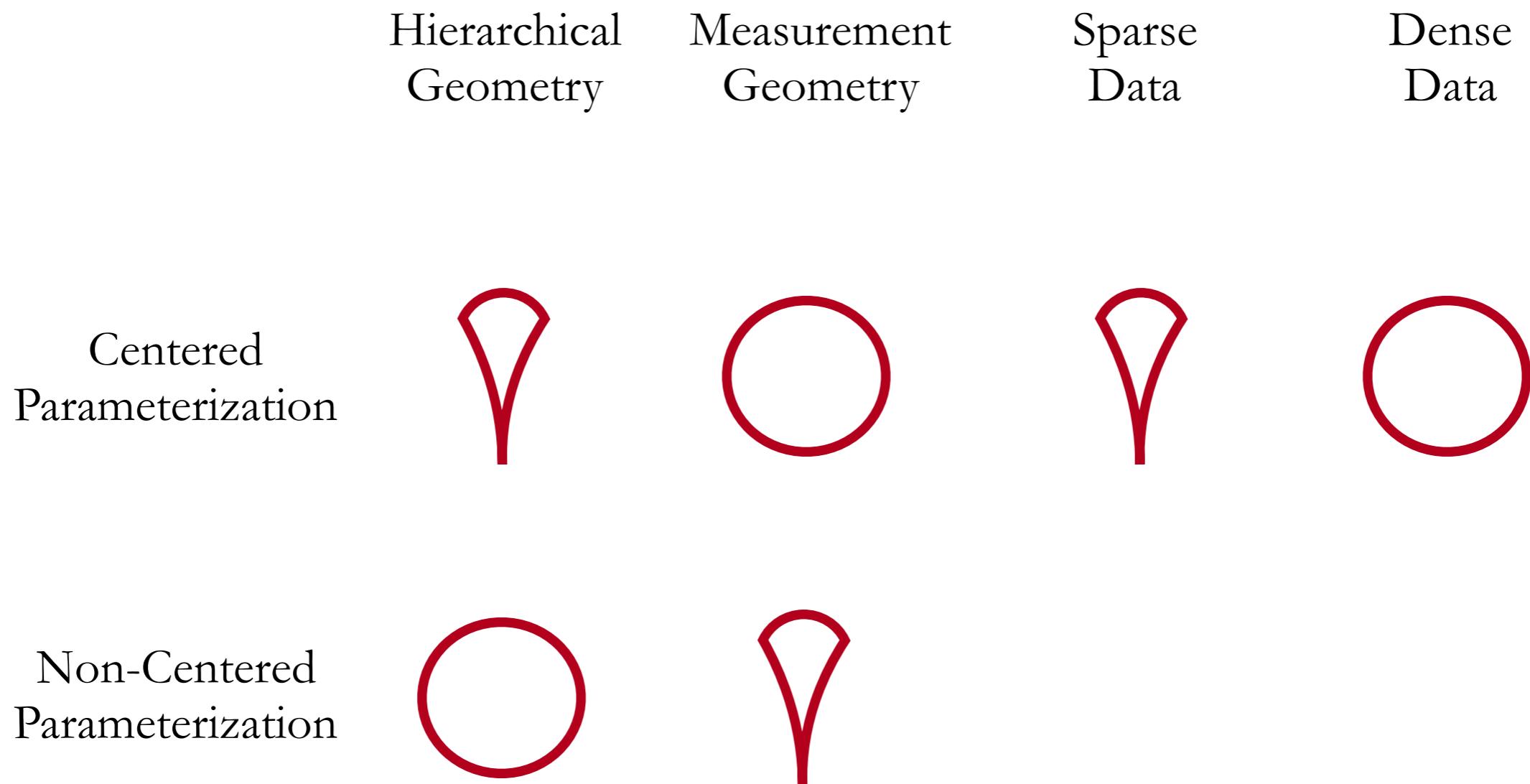
Consequently the two parameterizations are each advantageous in different circumstances.



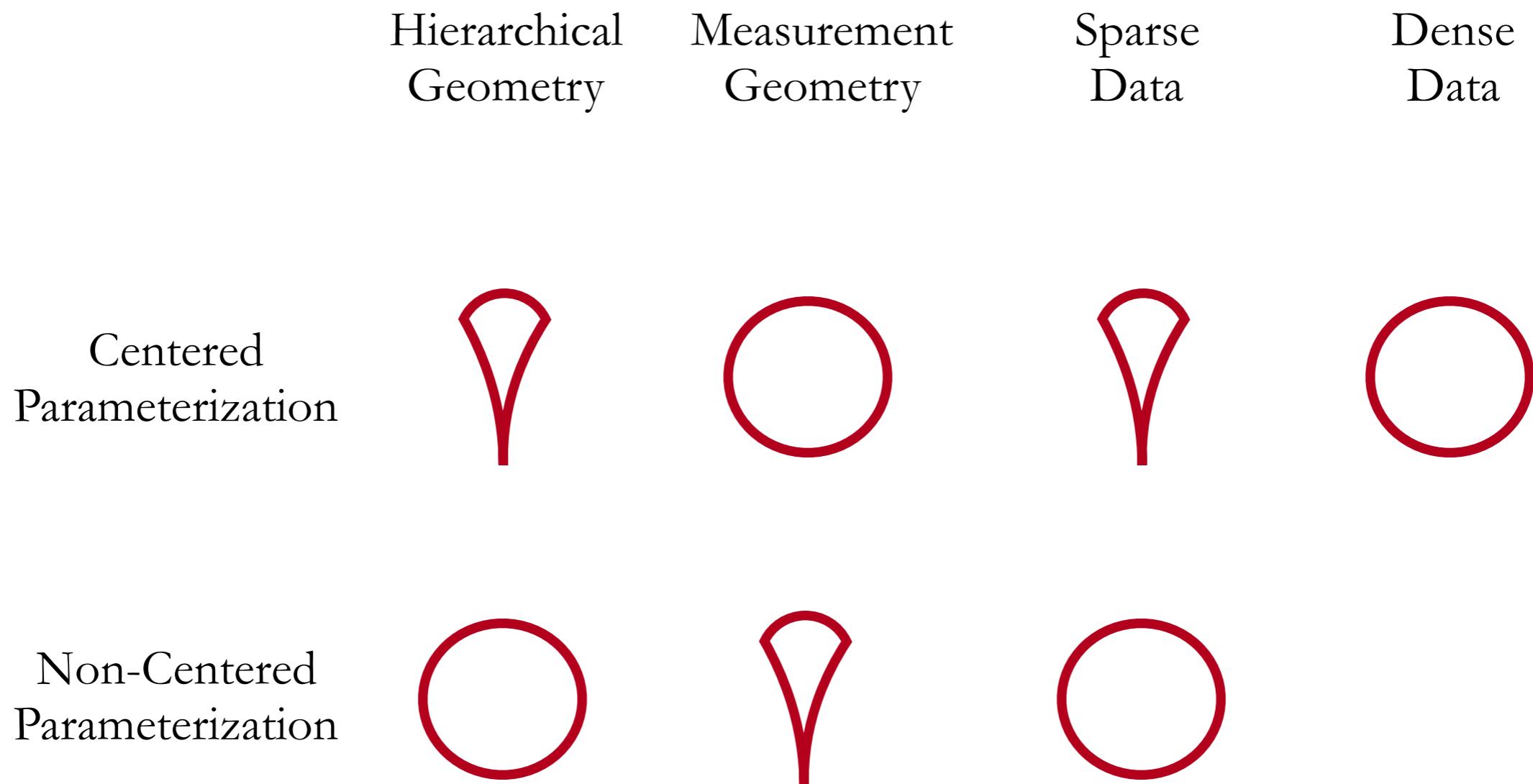
Consequently the two parameterizations are each advantageous in different circumstances.



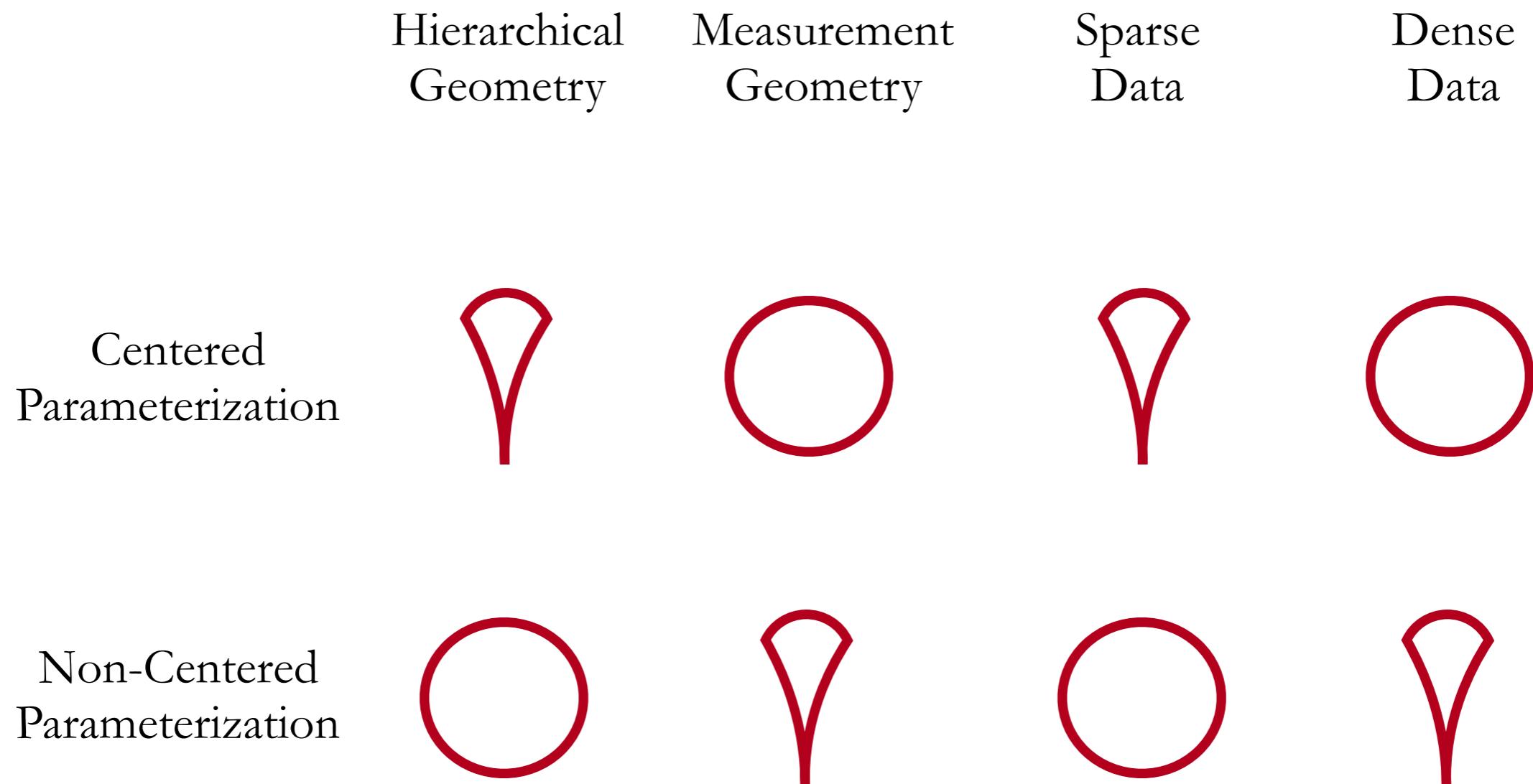
Consequently the two parameterizations are each advantageous in different circumstances.



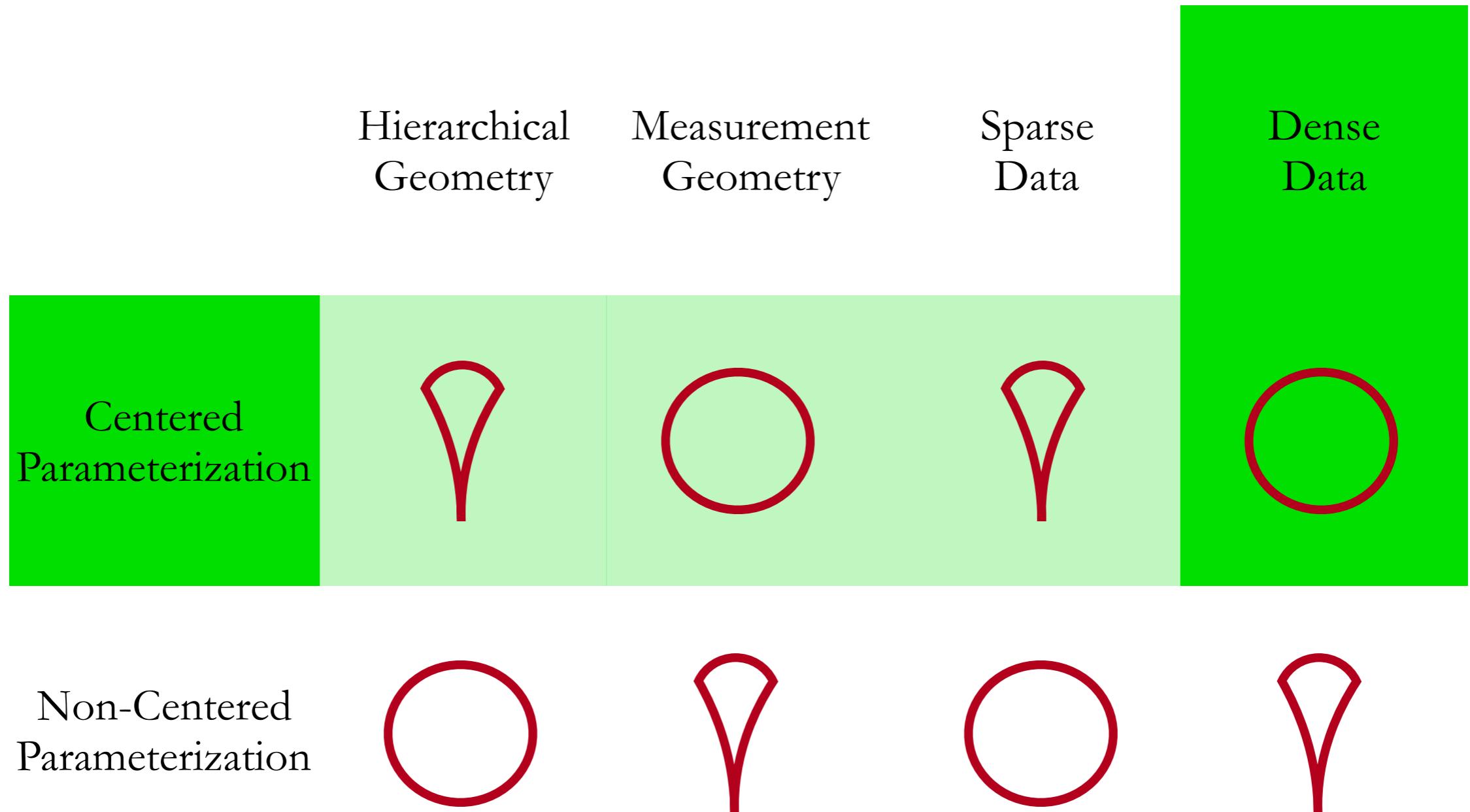
Consequently the two parameterizations are each advantageous in different circumstances.



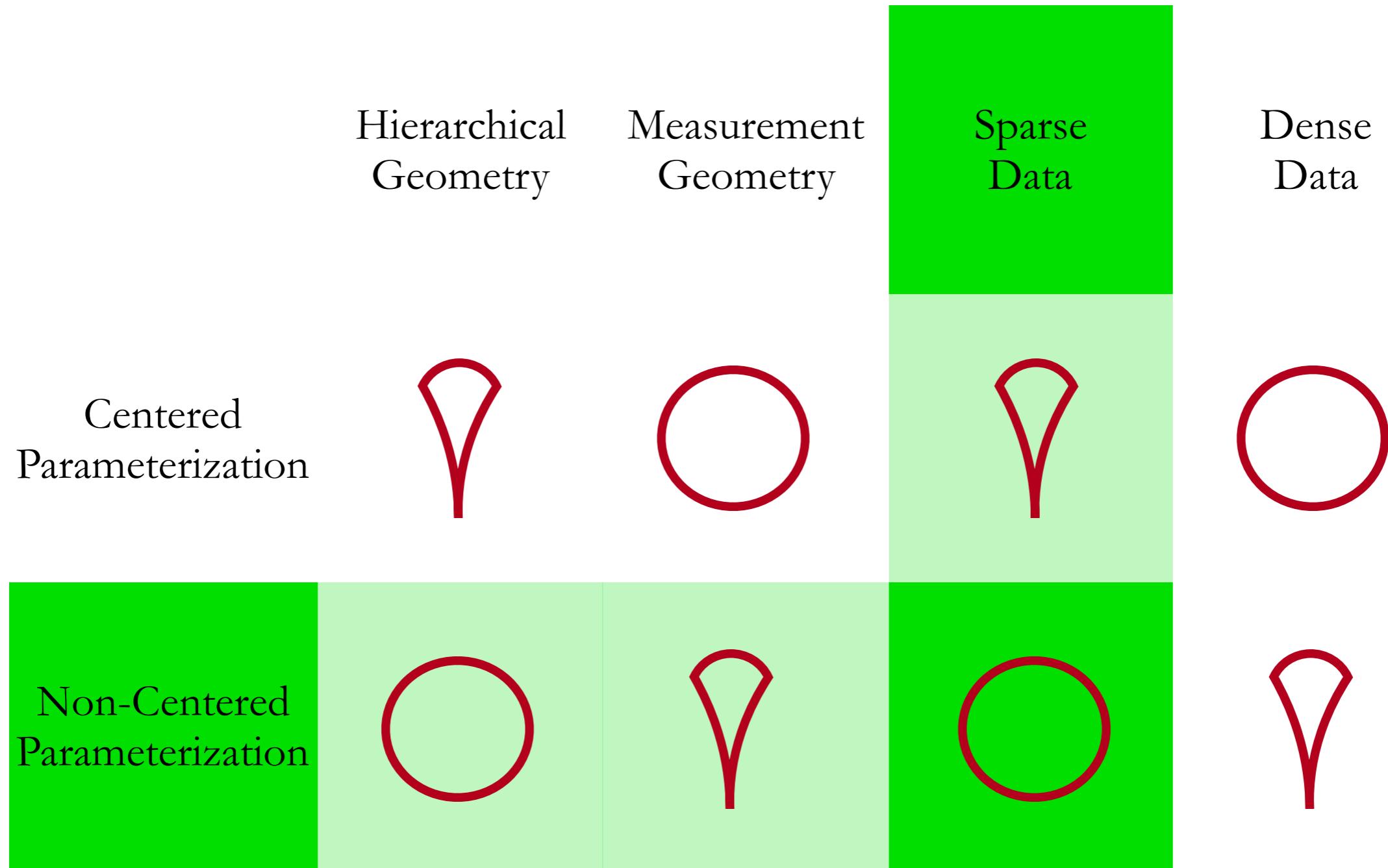
Consequently the two parameterizations are each advantageous in different circumstances.



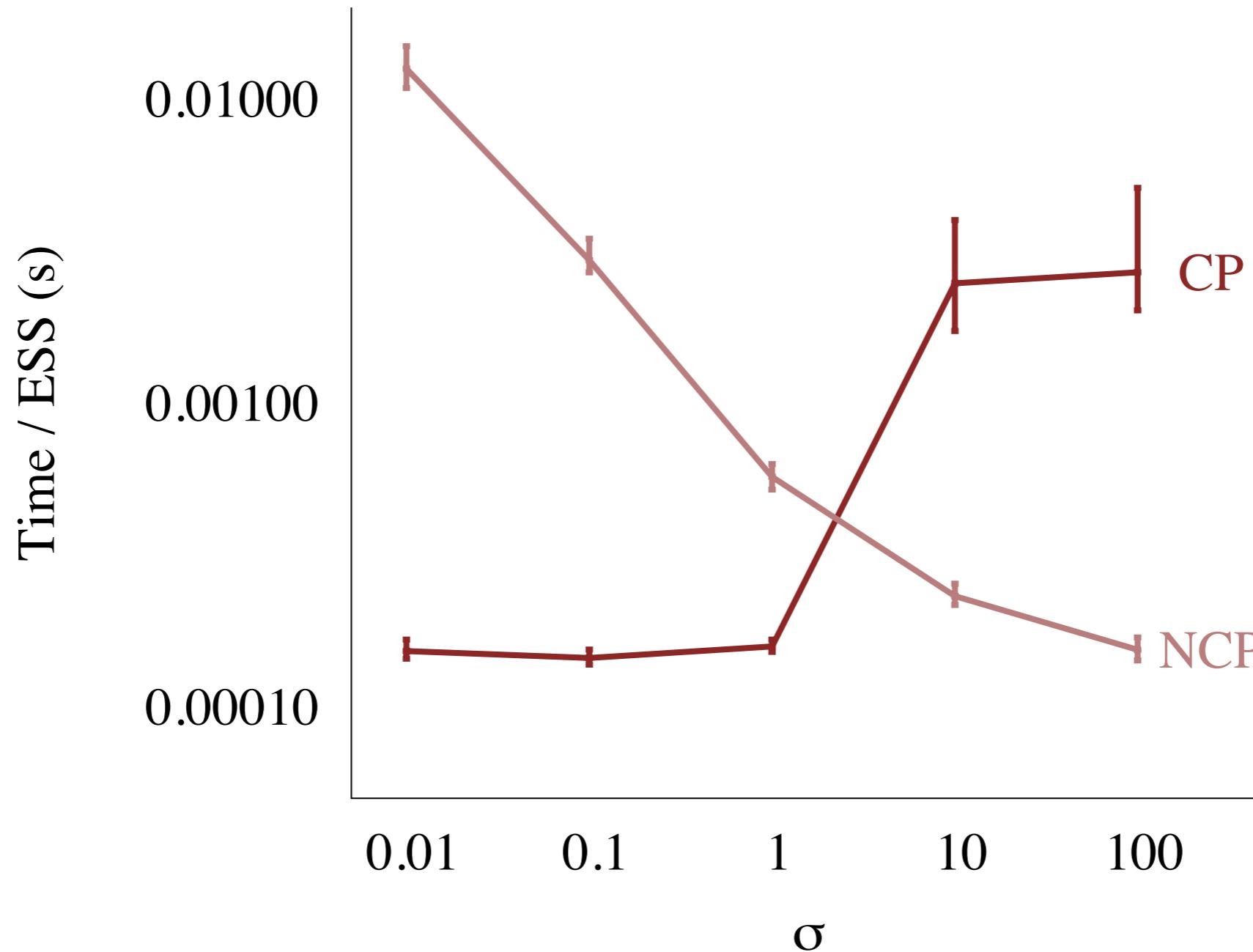
Consequently the two parameterizations are each advantageous in different circumstances.



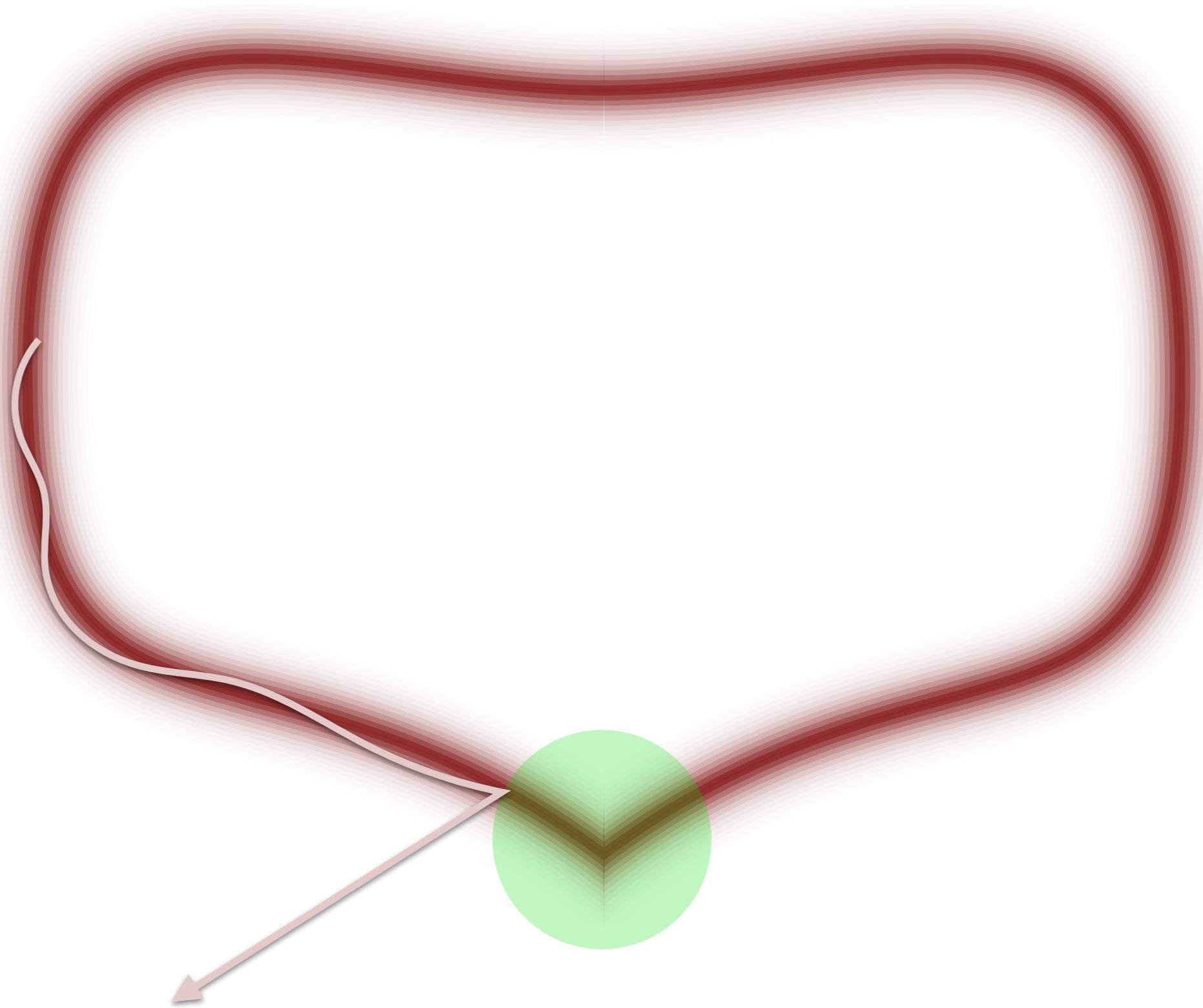
Consequently the two parameterizations are each advantageous in different circumstances.



Consequently the two parameterizations are each advantageous in different circumstances.

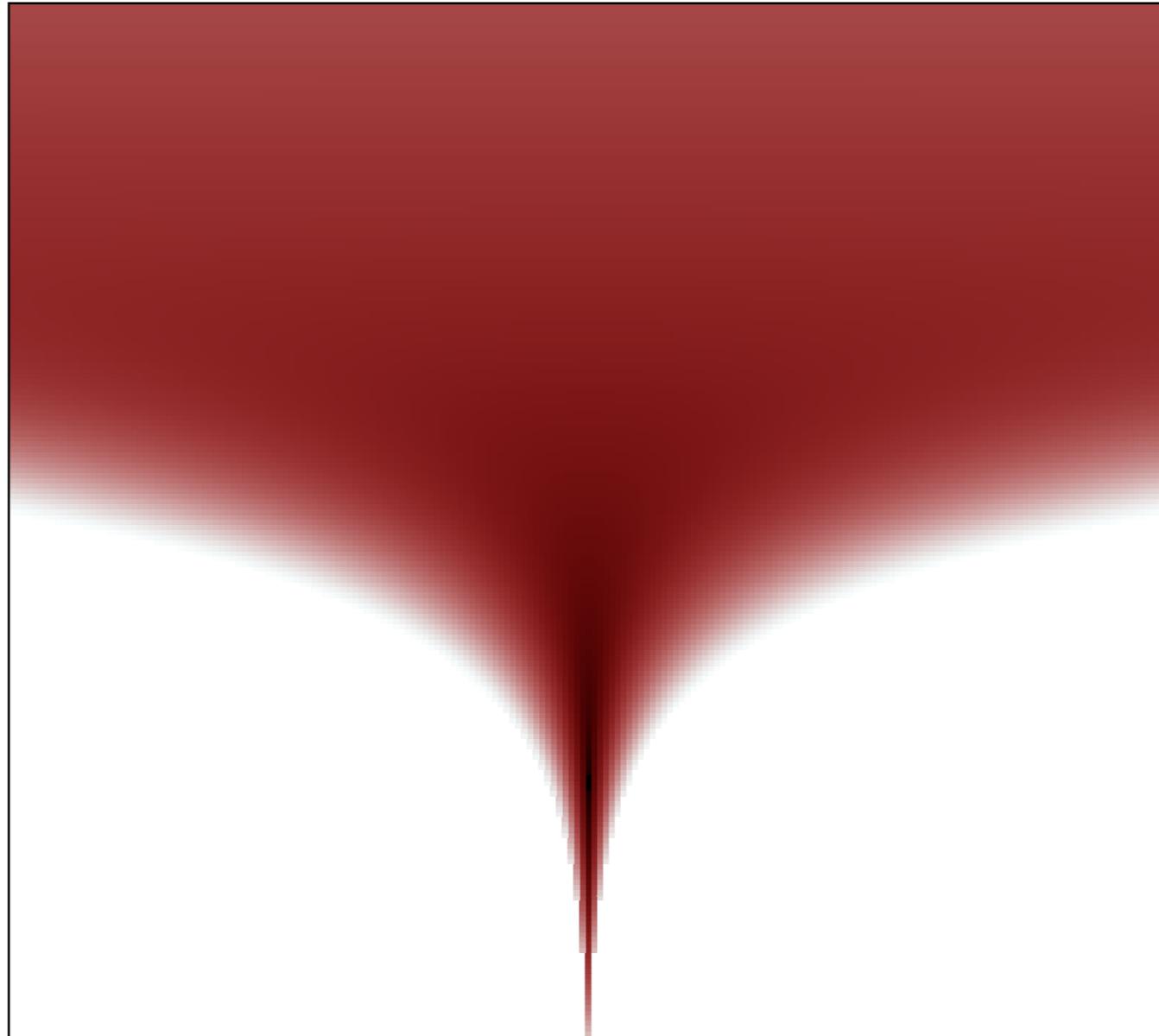


In practice, careful consideration of MCMC diagnostics help to identify poorly chosen parameterizations.



$\phi$

$\theta_n$



# StanCon 2018!

Wed Jan 10, 2018 - Fri Jan 12, 2018

Asilomar, California

Susan Holmes (Department of Statistics, Stanford University)

Frank Harrell (School of Medicine and Department of Biostatistics, Vanderbilt University)

Sophia Rabe-Hesketh (Educational Statistics and Biostatistics, University of California, Berkeley)

Sean Taylor and Ben Letham (Facebook Core Data Science)

Manuel Rivas (Department of Biomedical Data Science, Stanford University)

Talia Weiss (Department of Physics, Massachusetts Institute of Technology)

<http://mc-stan.org/events/stancon2018/>