

# Amortized Variational Inference for Simple Hierarchical Models

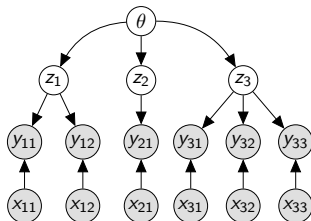
Abhinav Agrawal and Justin Domke

College of Information and Computer Science, UMass Amherst

## Paper in a slide

Hierarchical Branch Distribution (HBD)

$$p(\theta, z, y | x) = p(\theta) \prod_{i=1}^N p(z_i | \theta) p(y_i | \theta, z_i, x_i)$$



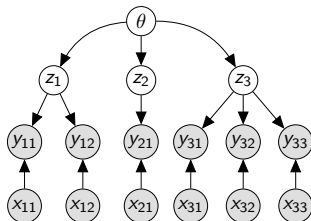
## Paper in a slide

Hierarchical Branch Distribution (HBD)

$$p(\theta, z, y | x) = p(\theta) \prod_{i=1}^N p(z_i | \theta) p(y_i | \theta, z_i, x_i)$$

Variational Inference (VI)

Approximate  $p(\theta, z | x, y)$  with  $q_\phi(\theta, z)$ .



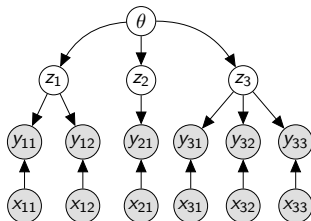
## Paper in a slide

Hierarchical Branch Distribution (HBD)

$$p(\theta, z, y | x) = p(\theta) \prod_{i=1}^N p(z_i | \theta) p(y_i | \theta, z_i, x_i)$$

Variational Inference (VI)

Approximate  $p(\theta, z | x, y)$  with  $q_\phi(\theta, z)$ .



**What we want?** Make sub-sampling *work* for large models.

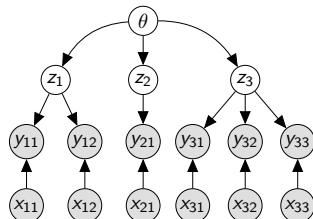
## Paper in a slide

Hierarchical Branch Distribution (HBD)

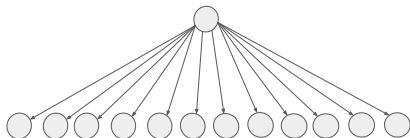
$$p(\theta, z, y | x) = p(\theta) \prod_{i=1}^N p(z_i | \theta) p(y_i | \theta, z_i, x_i)$$

Variational Inference (VI)

Approximate  $p(\theta, z | x, y)$  with  $q_\phi(\theta, z)$ .



**What we want?** Make sub-sampling *work* for large models.



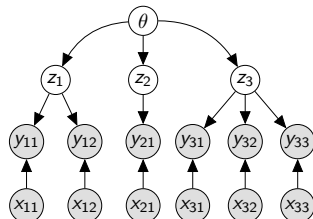
## Paper in a slide

Hierarchical Branch Distribution (HBD)

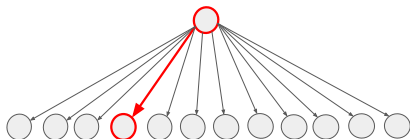
$$p(\theta, z, y | x) = p(\theta) \prod_{i=1}^N p(z_i | \theta) p(y_i | \theta, z_i, x_i)$$

Variational Inference (VI)

Approximate  $p(\theta, z | x, y)$  with  $q_\phi(\theta, z)$ .



**What we want?** Make sub-sampling *work* for large models.



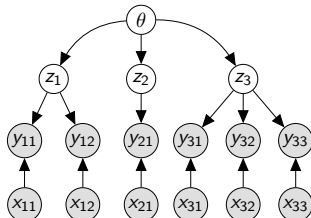
## Paper in a slide

Hierarchical Branch Distribution (HBD)

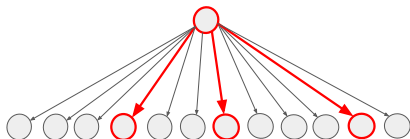
$$p(\theta, z, y | x) = p(\theta) \prod_{i=1}^N p(z_i | \theta) p(y_i | \theta, z_i, x_i)$$

Variational Inference (VI)

Approximate  $p(\theta, z | x, y)$  with  $q_\phi(\theta, z)$ .



**What we want?** Make sub-sampling *work* for large models.



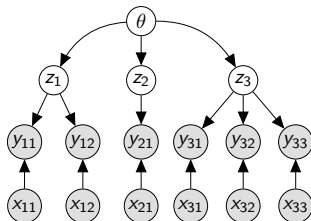
# Paper in a slide

Hierarchical Branch Distribution (HBD)

$$p(\theta, z, y|x) = p(\theta) \prod_{i=1}^N p(z_i|\theta) p(y_i|\theta, z_i, x_i)$$

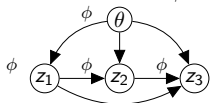
Variational Inference (VI)

Approximate  $p(\theta, z|x, y)$  with  $q_\phi(\theta, z)$ .



**What we want?** Make sub-sampling *work* for large models.

Joint Approach,  $q_\phi^{\text{Joint}}$



$$q_\phi^{\text{Joint}}(\theta, z) = q_\phi(\theta, z)$$

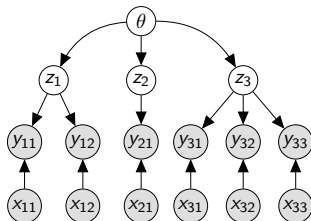
Optimize  $\phi$



# Paper in a slide

Hierarchical Branch Distribution (HBD)

$$p(\theta, z, y|x) = p(\theta) \prod_{i=1}^N p(z_i|\theta) p(y_i|\theta, z_i, x_i)$$

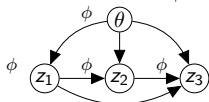


Variational Inference (VI)

Approximate  $p(\theta, z|x, y)$  with  $q_\phi(\theta, z)$ .

**What we want?** Make sub-sampling *work* for large models.

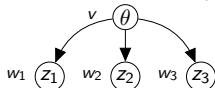
Joint Approach,  $q_\phi^{\text{Joint}}$



$$q_\phi^{\text{Joint}}(\theta, z) = q_\phi(\theta, z)$$

Optimize  $\phi$

Branch Approach,  $q_{v,w}^{\text{Branch}}$



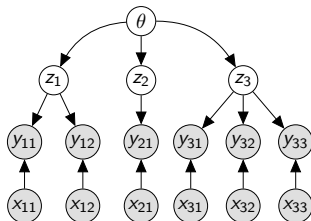
$$q_{v,w}^{\text{Branch}}(\theta, z) = q_v(\theta) \prod_{i=1}^3 q_{w_i}(z_i|\theta)$$

Optimize  $v, \{w_i\}_{i=1}^N$

# Paper in a slide

Hierarchical Branch Distribution (HBD)

$$p(\theta, z, y|x) = p(\theta) \prod_{i=1}^N p(z_i|\theta) p(y_i|\theta, z_i, x_i)$$

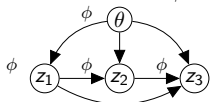


Variational Inference (VI)

Approximate  $p(\theta, z|x, y)$  with  $q_\phi(\theta, z)$ .

**What we want?** Make sub-sampling *work* for large models.

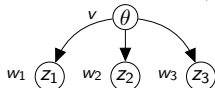
Joint Approach,  $q_\phi^{\text{Joint}}$



$$q_\phi^{\text{Joint}}(\theta, z) = q_\phi(\theta, z)$$

Optimize  $\phi$

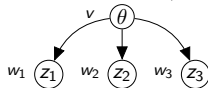
Branch Approach,  $q_{v,w}^{\text{Branch}}$



$$q_{v,w}^{\text{Branch}}(\theta, z) = q_v(\theta) \prod_{i=1}^3 q_{w_i}(z_i|\theta)$$

Optimize  $v, \{w_i\}_{i=1}^N$

Our Approach,  $q_{v,u}^{\text{Amort}}$



$$w_i = \text{net}_u(x_i, y_i)$$

Optimize  $v, u$

Why use  $q_{v,u}^{\text{Amort}}$ ?

# Why use $q_{v,u}^{\text{Amort}}$ ?

Faster Optimization

$q_{v,u}^{\text{Amort}}$  updates  $u$  with every step.

# Why use $q_{v,u}^{\text{Amort}}$ ?

Faster Optimization

$q_{v,u}^{\text{Amort}}$  updates  $u$  with every step.

Provably as accurate

Assuming sufficiently capable net $_u$ ,

$$\min_{v,u} KL(q_{v,u}^{\text{Amort}} \| p) \leq \min_{\phi} KL(q_{\phi}^{\text{Joint}} \| p).$$

# Why use $q_{v,u}^{\text{Amort}}$ ?

Faster Optimization

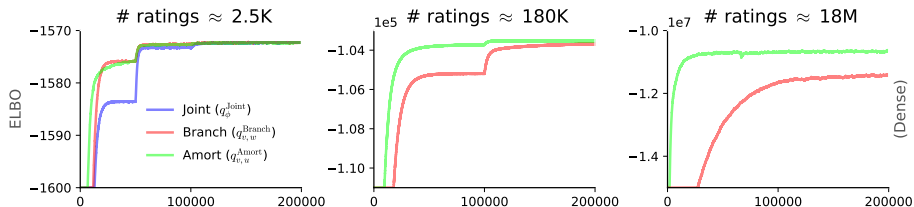
$q_{v,u}^{\text{Amort}}$  updates  $u$  with every step.

Provably as accurate

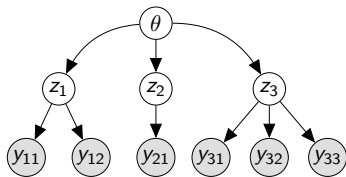
Assuming sufficiently capable  $\text{net}_u$ ,

$$\min_{v,u} KL(q_{v,u}^{\text{Amort}} \| p) \leq \min_{\phi} KL(q_{\phi}^{\text{Joint}} \| p).$$

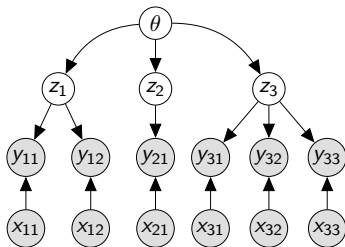
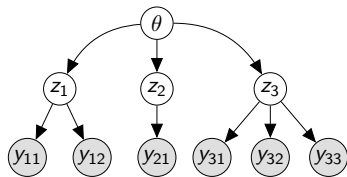
MovieLens dataset.  $\approx 160\text{K}$  users.  $\approx 18\text{M}$  ratings.  $\approx 1.6\text{M}$  variables.



# Hierarchical Branch Distributions

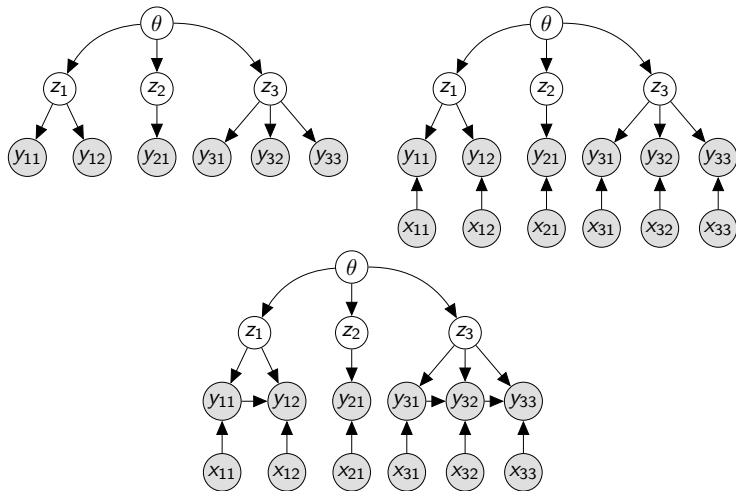


# Hierarchical Branch Distributions

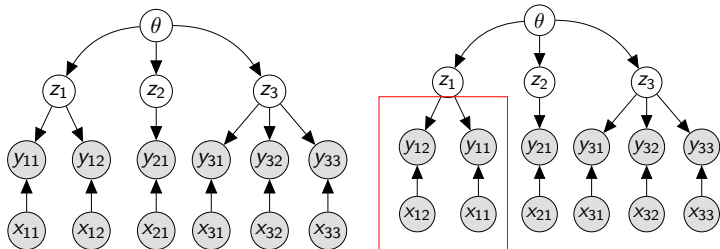




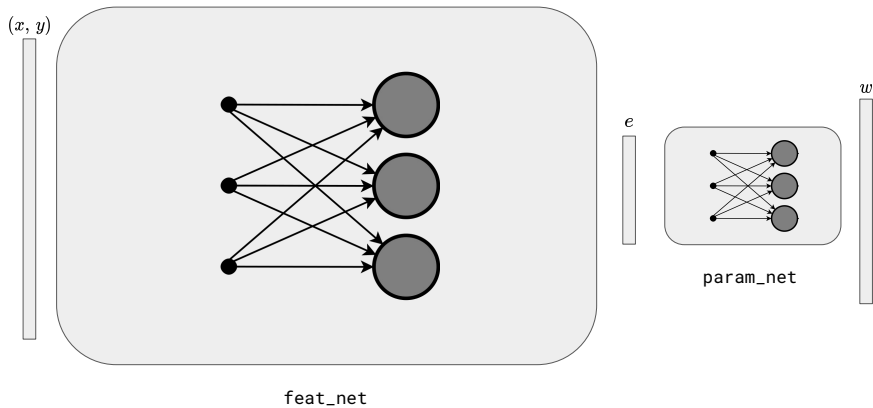
# Hierarchical Branch Distributions



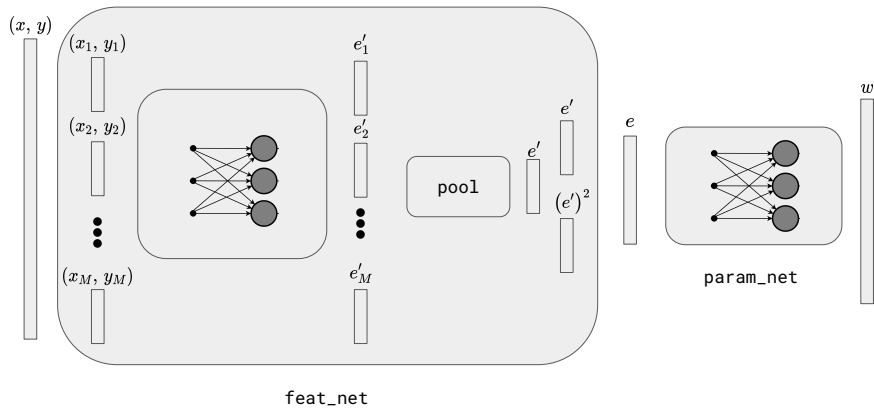
## Locally i.i.d



# Designing $\text{net}_u$



# Designing $\text{net}_u$



# Amortized Gaussian VI

## Dense Gaussian

---

Joint Approach

Branch Approach

Our Approach

---

---

# Amortized Gaussian VI

## Dense Gaussian

---

Joint Approach

Branch Approach

Our Approach

---

$$\mathcal{N}(\theta, z | \mu, \Sigma)$$

$$\phi = (\mu, \Sigma)$$

---

# Amortized Gaussian VI

## Dense Gaussian

---

Joint Approach

Branch Approach

Our Approach

---

$$\mathcal{N}(\theta, z | \mu, \Sigma) \quad \mathcal{N}(\theta | \mu_0, \Sigma_0) \prod_{i=1}^N \mathcal{N}(z_i | \mu_i + A_i \theta, \Sigma_i)$$

$$\phi = (\mu, \Sigma)$$

$$v = (\mu_0, \Sigma_0), w = \{\mu_i, A_i, \Sigma_i\}_{i=1}^N$$

---

# Amortized Gaussian VI

## Dense Gaussian

Joint Approach	Branch Approach	Our Approach
$\mathcal{N}(\theta, z \mu, \Sigma)$	$\mathcal{N}(\theta \mu_0, \Sigma_0) \prod_{i=1}^N \mathcal{N}(z_i \mu_i + A_i\theta, \Sigma_i)$	$\mathcal{N}(\theta \mu_0, \Sigma_0) \prod_{i=1}^N \mathcal{N}(z_i \mu_i + A_i\theta, \Sigma_i)$ where, $(\mu_i, A_i, \Sigma_i) = \text{net}_u(x_i, y_i)$
$\phi = (\mu, \Sigma)$	$v = (\mu_0, \Sigma_0), w = \{\mu_i, A_i, \Sigma_i\}_{i=1}^N$	$v = (\mu_0, \Sigma_0), u$



# Amortized Gaussian VI

## Dense Gaussian

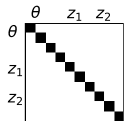
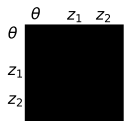
Joint Approach	Branch Approach	Our Approach
$\mathcal{N}(\theta, z \mu, \Sigma)$	$\mathcal{N}(\theta \mu_0, \Sigma_0) \prod_{i=1}^N \mathcal{N}(z_i \mu_i + A_i\theta, \Sigma_i)$	$\mathcal{N}(\theta \mu_0, \Sigma_0) \prod_{i=1}^N \mathcal{N}(z_i \mu_i + A_i\theta, \Sigma_i)$ where, $(\mu_i, A_i, \Sigma_i) = \text{net}_u(x_i, y_i)$
$\phi = (\mu, \Sigma)$	$v = (\mu_0, \Sigma_0), w = \{\mu_i, A_i, \Sigma_i\}_{i=1}^N$	$v = (\mu_0, \Sigma_0), u$

Scalable

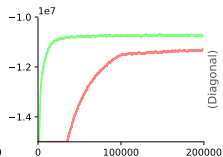
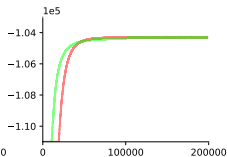
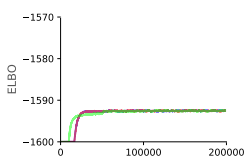
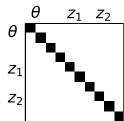
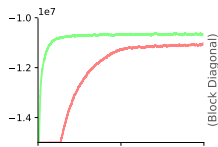
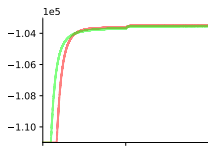
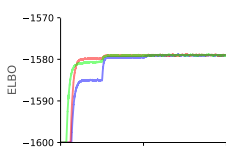
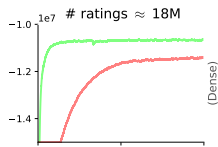
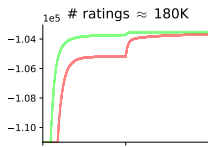
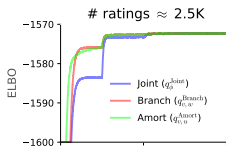
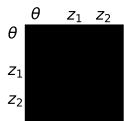
Faster

Accurate

# More results



# More results



# Insights

When  $N$  is large,  $\theta$  collapses.

No need to condition on it.

