

# Scaling up natural gradient by factorizing Fisher information

Roger Grosse

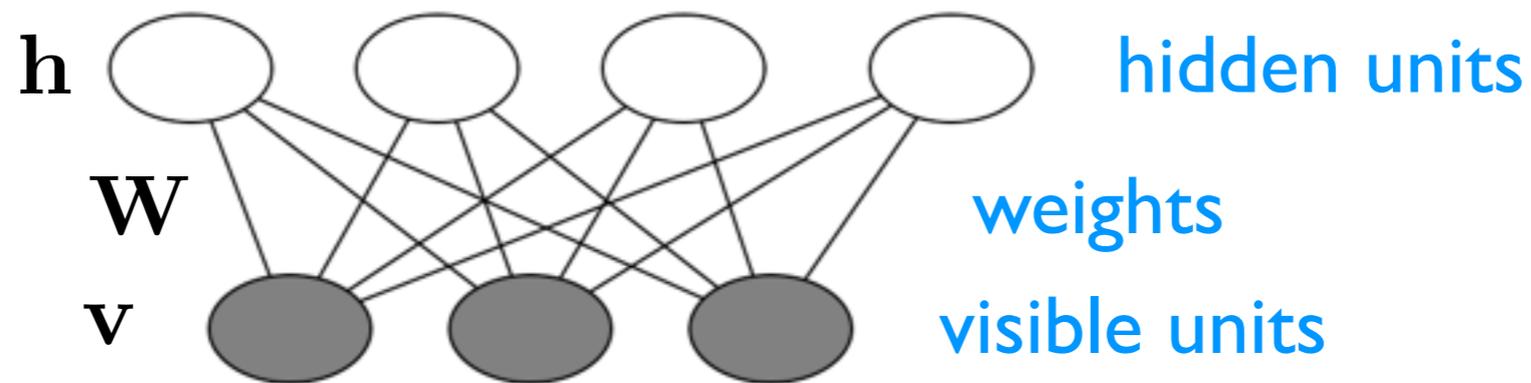


# Introduction

- Training RBMs is still a black art for several reasons
  - We can't evaluate the likelihood since this requires the intractable partition function
  - We can't compute the gradient exactly
  - The optimization problem has ill-conditioned curvature since small changes to the model parameters can dramatically change the distribution it represents
- I will focus on the third issue

# Introduction

- RBMs are an undirected graphical model:



$$p(\mathbf{v}, \mathbf{h}) = \frac{1}{Z} \exp(\mathbf{v}^T \mathbf{W} \mathbf{h} + \mathbf{a}^T \mathbf{v} + \mathbf{b}^T \mathbf{h})$$

$Z$  is the partition function, which is intractable to compute exactly

# Introduction

- The log-likelihood gradient is a difference of two expectations:

$$\frac{\partial}{\partial w_{ij}} \log p(\mathbf{x}) = \underbrace{\mathbb{E}_{p(\mathbf{h}|\mathbf{x})} [x_i h_j]}_{\substack{\text{conditional distribution} \\ \text{given data} \\ \text{(easy)}}} - \underbrace{\mathbb{E}_{p(\mathbf{v}, \mathbf{h})} [v_i h_j]}_{\substack{\text{model distribution} \\ \text{(hard)}}$$

- Many clever approximations to the model statistics, e.g.
  - contrastive divergence
  - persistent contrastive divergence
  - tempered transitions

# Introduction

- It's expensive to approximate the gradient well, so we'd like to accomplish more per iteration
- Second-order optimization methods find better updates by accounting for curvature
- Usually the practical ones either:
  - use a cheap, tractable approximation (e.g. diagonal or block diagonal)
  - require an expensive iterative procedure (e.g. Hessian-free optimization of neural nets)
- I will propose a second-order optimization method for RBMs which approximates the curvature well, but is only slightly more expensive than SGD

# Outline

- Information geometry background
  - exponential families
  - Fisher information matrix
- Factorized Natural Gradient (FANG) for RBM training
  - understanding the curvature of maximum likelihood for RBMs
  - approximating the curvature using a Gaussian graphical model
  - how to efficiently invert the approximate curvature
  - experiments on training RBMs as generative models

# Background: (Linear) exponential families

$$p(\mathbf{x}) = \frac{1}{\mathcal{Z}(\boldsymbol{\eta})} h(\mathbf{x}) \exp(\boldsymbol{\eta}^T \mathbf{g}(\mathbf{x}))$$

sufficient statistics  $\mathbf{g}(\mathbf{x})$

natural parameters  $\boldsymbol{\eta}$

moments  $\mathbf{s} = \mathbb{E}[\mathbf{g}(\mathbf{x})]$

partition function  $\mathcal{Z}(\boldsymbol{\eta})$

# Background: (Linear) exponential families

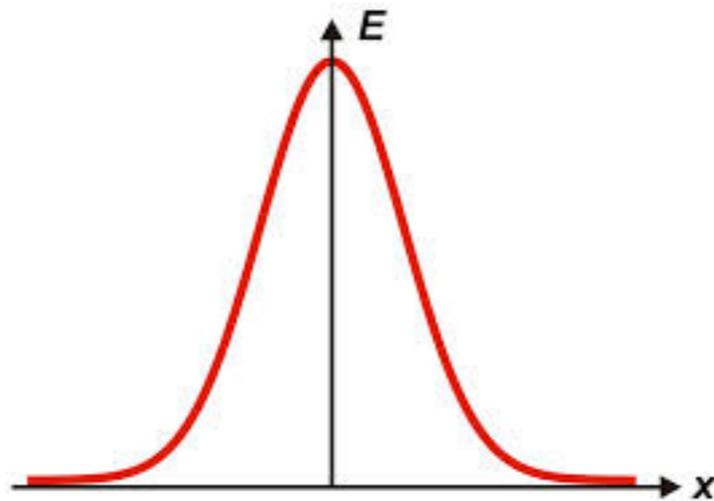
$$p(\mathbf{x}) = \frac{1}{\mathcal{Z}(\boldsymbol{\eta})} h(\mathbf{x}) \exp(\boldsymbol{\eta}^T \mathbf{g}(\mathbf{x}))$$



## Bernoulli

$$\mathbf{g}(x) = x$$

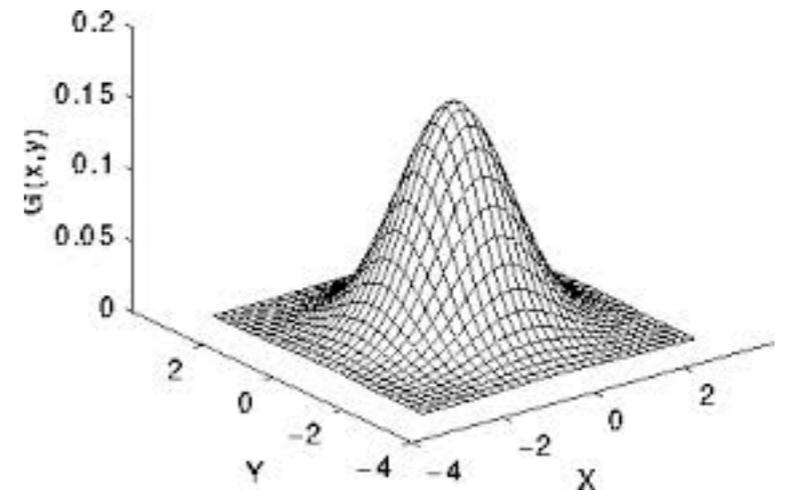
$$\boldsymbol{\eta} = \log \mu - \log(1 - \mu)$$



## Gaussian

$$\mathbf{g}(x) = \left( x, -\frac{1}{2}x^2 \right)$$

$$\boldsymbol{\eta} = (h, \lambda)$$

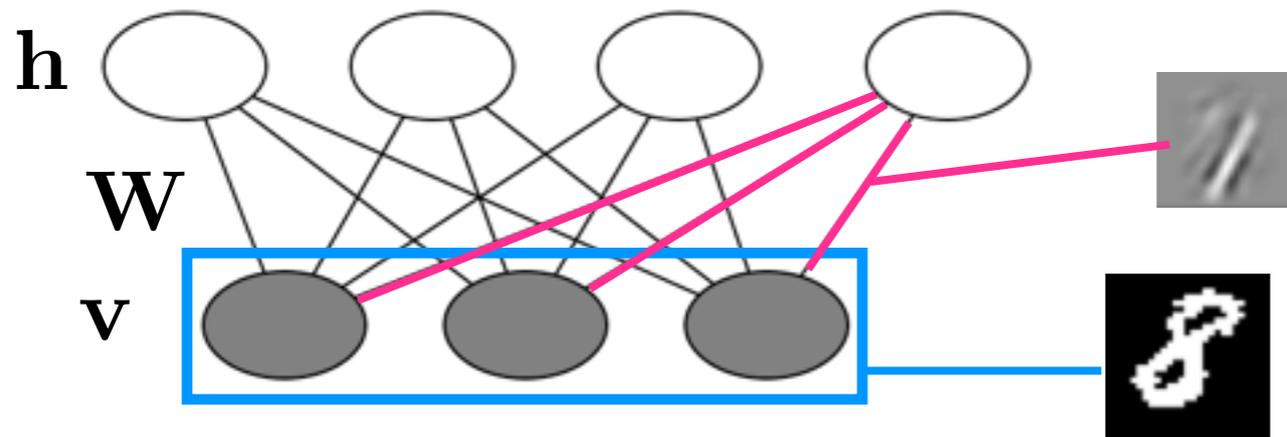


## Multivariate Gaussian

$$\mathbf{g}(\mathbf{x}) = \begin{pmatrix} \mathbf{x} \\ -\frac{1}{2} \text{vec}(\mathbf{x}\mathbf{x}^T) \end{pmatrix}$$

$$\boldsymbol{\eta} = \begin{pmatrix} \mathbf{h} \\ \text{vec}(\boldsymbol{\Lambda}) \end{pmatrix}$$

# Background: RBMs as exponential families



**joint distribution:**

$$p(\mathbf{v}, \mathbf{h}) \propto \exp(\mathbf{a}^T \mathbf{v} + \mathbf{b}^T \mathbf{h} + \mathbf{v}^T \mathbf{W} \mathbf{h})$$

**sufficient statistics:**

$$\mathbf{g}(\mathbf{v}, \mathbf{h}) = \begin{pmatrix} \mathbf{v} \\ \mathbf{h} \\ \text{vec}(\mathbf{v} \mathbf{h}^T) \end{pmatrix}$$

**Two parameterizations:**

$$\boldsymbol{\eta} = \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \\ \text{vec}(\mathbf{W}) \end{pmatrix} \begin{array}{c} \xrightarrow{\text{inference}} \\ \xleftarrow{\text{learning}} \\ \text{(fully observed)} \end{array}$$

**natural parameters**

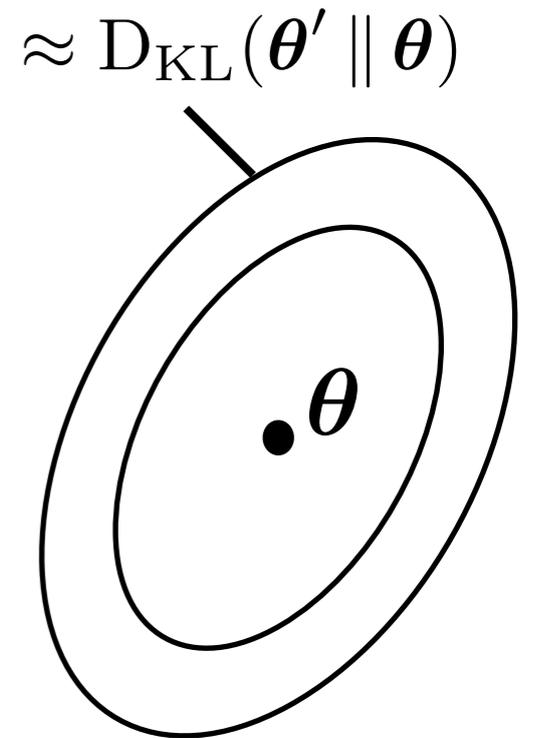
$$\mathbf{s} = \mathbb{E}[\mathbf{g}(\mathbf{v}, \mathbf{h})] = \begin{pmatrix} \mathbb{E}[\mathbf{v}] \\ \mathbb{E}[\mathbf{h}] \\ \text{vec}(\mathbb{E}[\mathbf{v} \mathbf{h}^T]) \end{pmatrix}$$

**moments**

# Background: Fisher information matrix

$$\mathbf{G}_\theta \triangleq \nabla_{\theta'}^2 D_{\text{KL}}(\theta' \parallel \theta) \Big|_{\theta'=\theta}$$

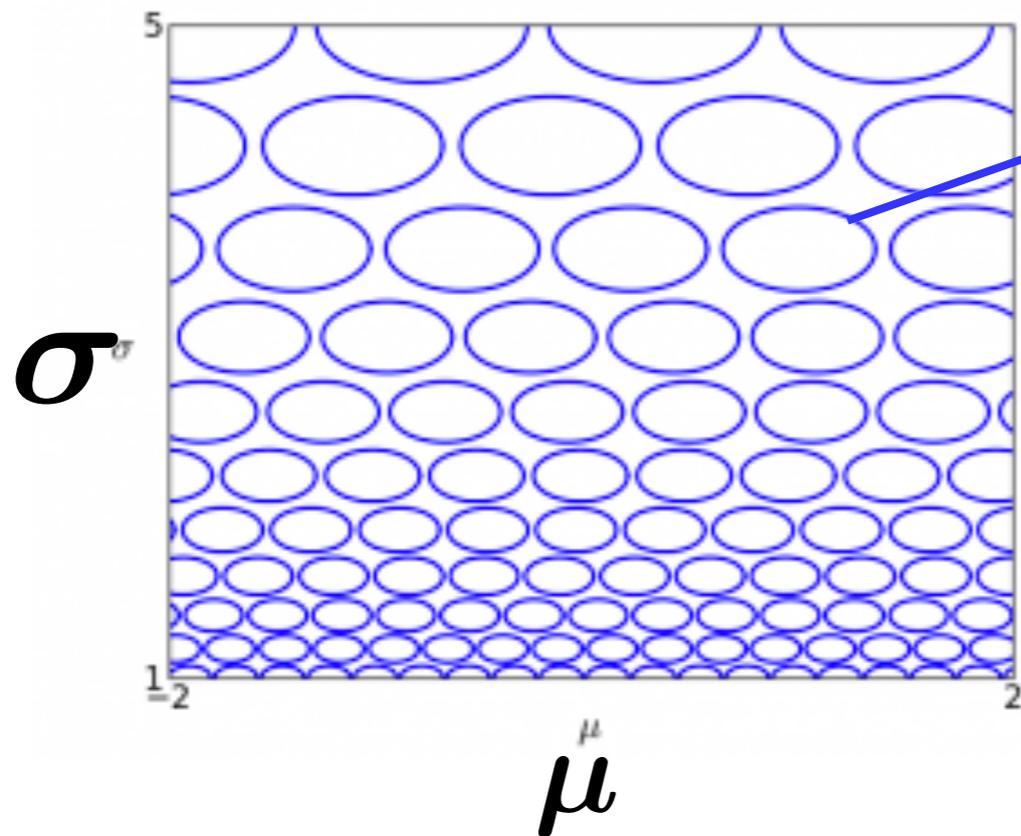
$$D_{\text{KL}}(\theta' \parallel \theta) \approx \frac{1}{2} (\theta' - \theta)^T \mathbf{G}_\theta (\theta' - \theta)$$



Aside: KL divergence is like *squared* distance

# Background: Fisher information matrix

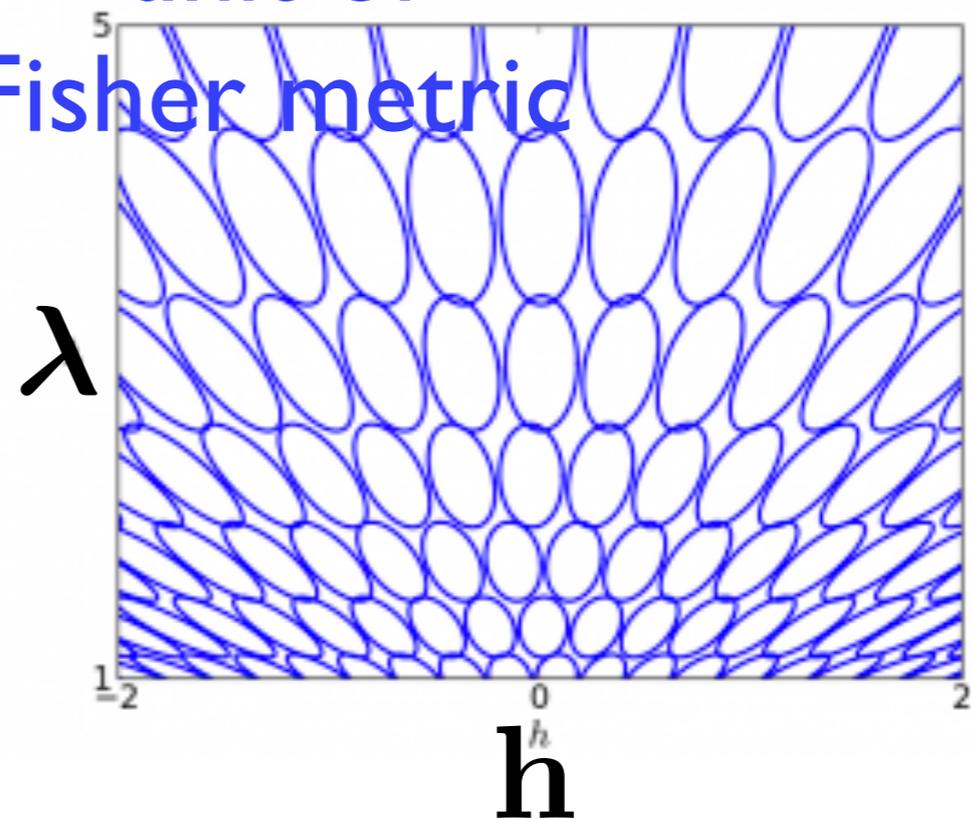
mean and variance



$$p(x) \propto \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

information form  
unit of

Fisher metric



$$p(x) \propto \exp\left(hx - \frac{\lambda}{2}x^2\right)$$

# Natural gradient

(Amari, 1998)

## gradient ascent

steepest ascent in Euclidean norm

depends on parameterization

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha \nabla_{\boldsymbol{\theta}} \ell$$

## Recall:

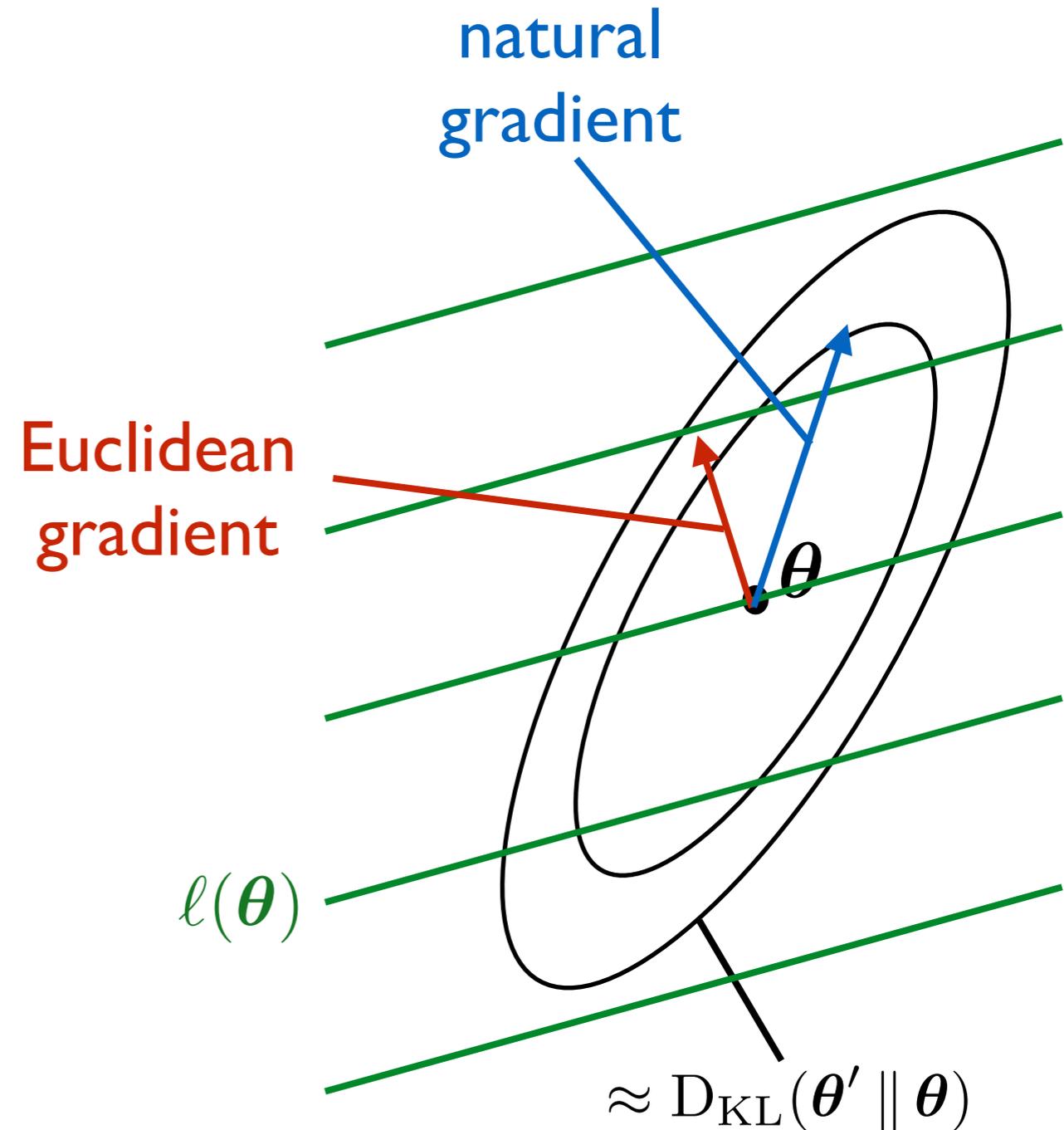
$$D_{\text{KL}}(\boldsymbol{\theta}' \parallel \boldsymbol{\theta}) \approx \frac{1}{2} (\boldsymbol{\theta}' - \boldsymbol{\theta})^T \mathbf{G}_{\boldsymbol{\theta}} (\boldsymbol{\theta}' - \boldsymbol{\theta})$$

## natural gradient ascent

steepest ascent in Fisher norm

independent of parameterization

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha \mathbf{G}^{-1} \nabla_{\boldsymbol{\theta}} \ell$$



# Natural gradient: RBMs

gradient ascent:  $\boldsymbol{\eta} \leftarrow \boldsymbol{\eta} + \alpha (\mathbf{s}_{\text{data}} - \mathbf{s}_{\text{model}})$

natural gradient ascent:  $\boldsymbol{\eta} \leftarrow \boldsymbol{\eta} + \alpha \mathbf{G}^{-1} (\mathbf{s}_{\text{data}} - \mathbf{s}_{\text{model}})$

to first order:  $\mathbf{s} \leftarrow \mathbf{s} + \alpha (\mathbf{s}_{\text{data}} - \mathbf{s}_{\text{model}})$   $(d\mathbf{s} = \mathbf{G}d\boldsymbol{\eta})$

Unfortunately...

$\mathbf{G}$  has  $(N_v + N_h + N_v N_h)^2$  entries

= 155 billion for an MNIST RBM!

same problem as for other second order methods (Newton, etc.)

# Natural gradient: RBMs

- Several approximations by analogy with Quasi-Newton
- Approximating  $G$ 
  - diagonal (e.g. Adagrad; Duchi et al., 2011)
  - block diagonal (e.g. Le Roux et al., 2008)
- Iterative methods
  - Park et al. (2000)
  - metric-free natural gradient (Desjardins et al., 2013)

# Factorized Natural Gradient (FaNG)

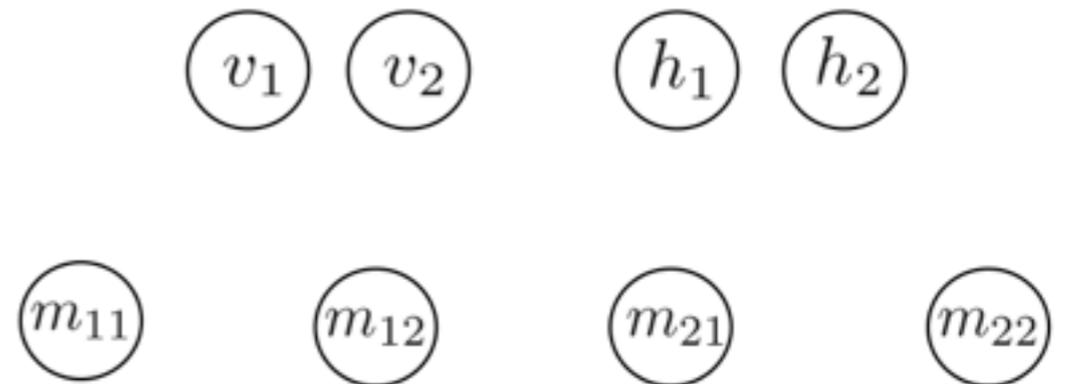
Fisher information is the covariance of the sufficient statistics:

$$G = \text{Cov}_{\mathbf{x}}(\mathbf{g}(\mathbf{x}))$$

We can approximate a covariance matrix with a Gaussian graphical model.

A diagonal approximation corresponds to a fully disconnected graph.

To come up with a better structure, let's analyze the Fisher information of an RBM.



# Fisher information of RBMs

computing  $G$  for a small RBM (20 hidden units)

$$\mathbb{E}[\mathbf{g}] = \sum_{\mathbf{h}} p(\mathbf{h}) \mathbb{E}[\mathbf{g} | \mathbf{h}]$$

$$\mathbb{E}[\mathbf{g}\mathbf{g}^T] = \sum_{\mathbf{h}} p(\mathbf{h}) (\mathbb{E}[\mathbf{g} | \mathbf{h}]\mathbb{E}[\mathbf{g} | \mathbf{h}]^T + \text{Cov}[\mathbf{g} | \mathbf{h}]) .$$

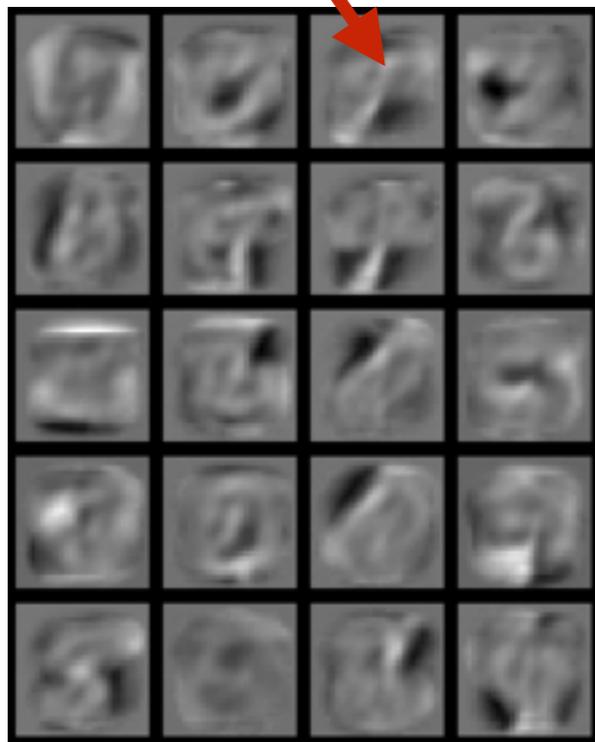
$$\mathbf{G} = \mathbb{E}[\mathbf{g}\mathbf{g}^T] - \mathbb{E}[\mathbf{g}]\mathbb{E}[\mathbf{g}]^T$$

# Fisher information of RBMs

$$ds = \underbrace{G}_{\text{Fisher information}} d\underbrace{\eta}_{\text{change in natural parameters}}$$

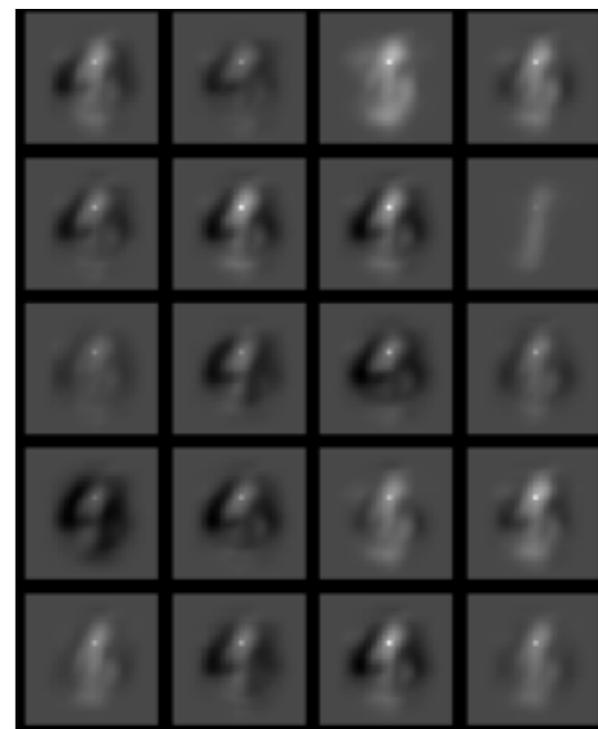
change in moments

Changing this weight...



$W$

... does this to the moments



$dE[\mathbf{v}\mathbf{h}^T]$

# Fisher information of RBMs

$$ds = G d\eta$$

change in moments (blue arrow pointing to  $ds$ )

Fisher information (green arrow pointing to  $G$ )

change in natural parameters (red arrow pointing to  $d\eta$ )

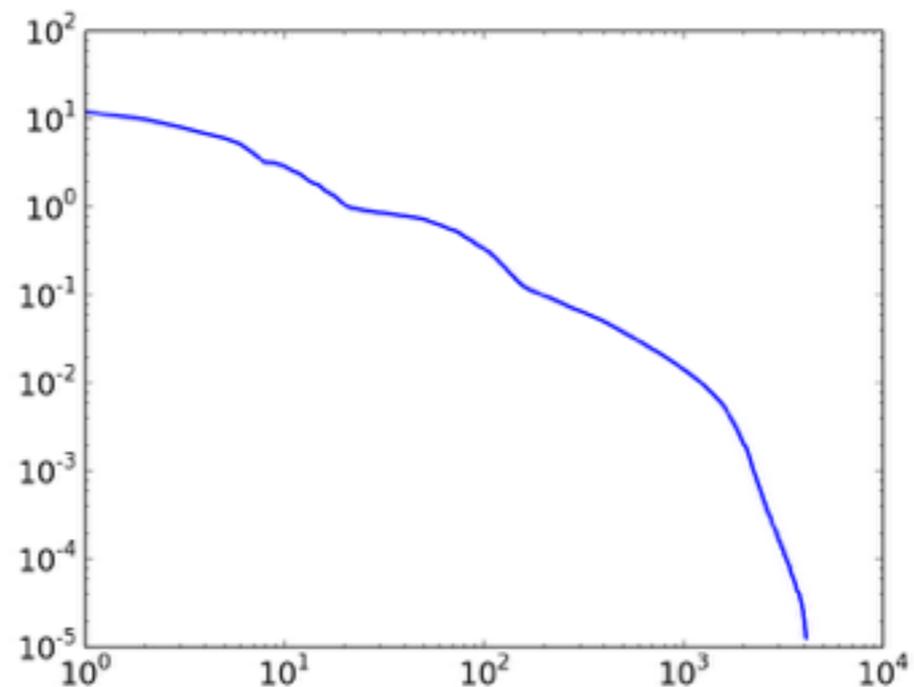
Most of the sufficient statistics are second moments, so let's think about how these vary with the parameters.

$$d\mathbb{E}[\mathbf{v}\mathbf{h}^T] = \underbrace{d\mathbb{E}[\mathbf{v}]\mathbb{E}[\mathbf{h}]^T + \mathbb{E}[\mathbf{v}]d\mathbb{E}[\mathbf{h}]^T}_{\text{change in second moments}} + \underbrace{d\text{Cov}(\mathbf{v}, \mathbf{h})}_{\text{change in covariances (interesting)}}.$$

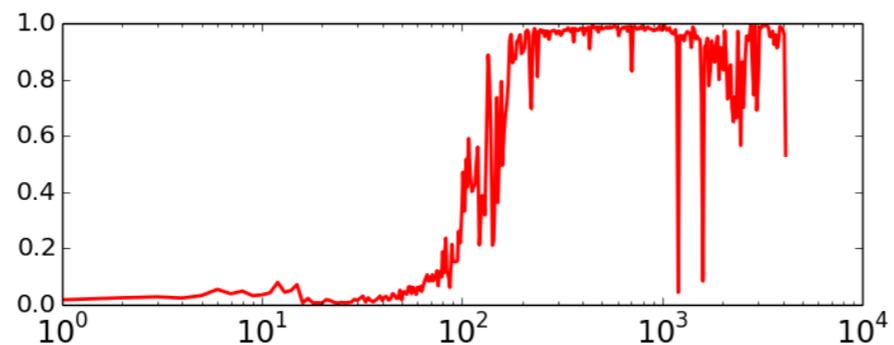
change in second moments = change in means (boring) + change in covariances (interesting)

# Fisher information of RBMs

eigenspectrum of G

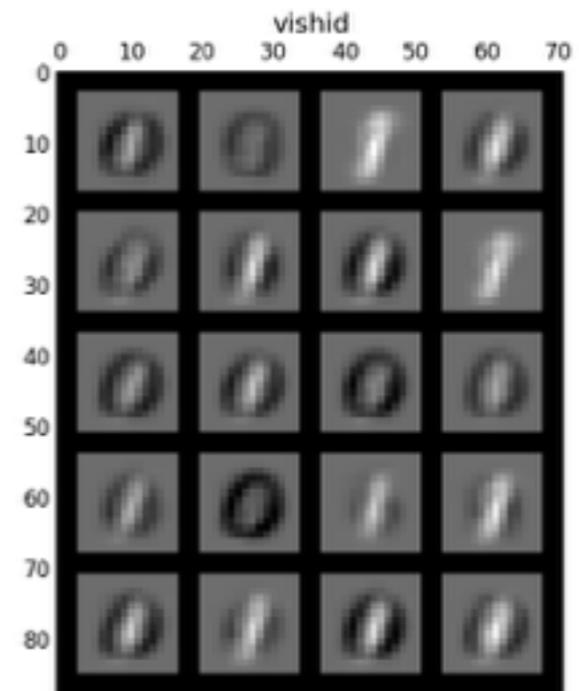


$$\frac{\|d \text{Cov}(\mathbf{v}, \mathbf{h})\|^2}{\|d\mathbf{s}\|^2}$$



mostly  
1st order  
stats

mostly  
covariances

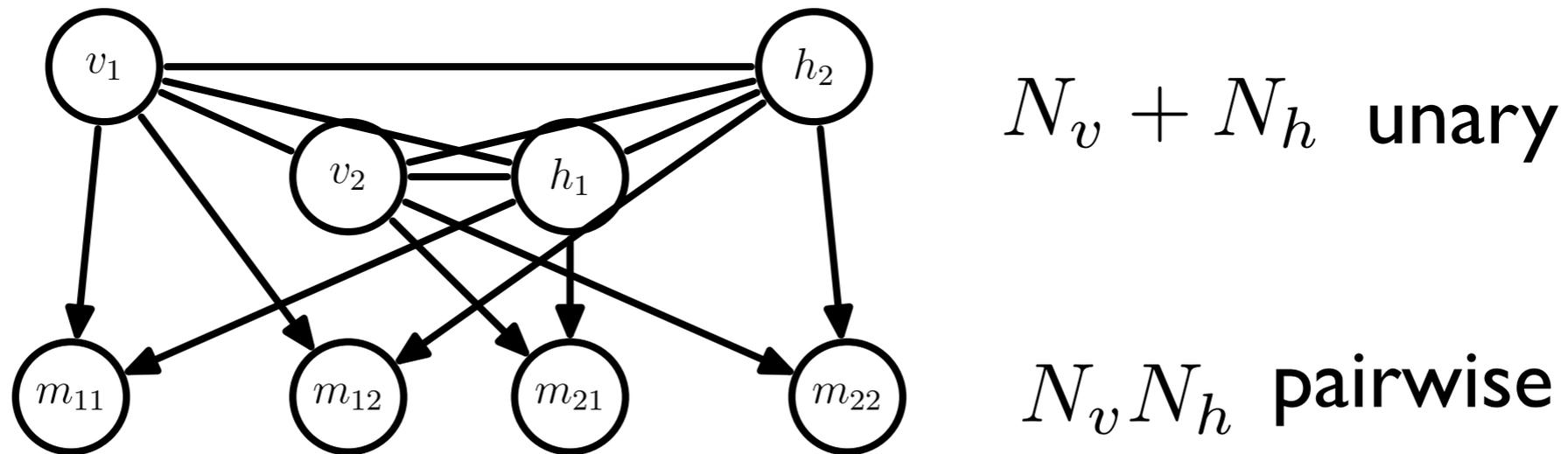


top eigenvector  
(direction of highest curvature)

# Factorized Natural Gradient (FaNG)

We just saw that the first-order statistics are driving everything.

We propose a graphical model structure which models dependencies between sufficient statistics.

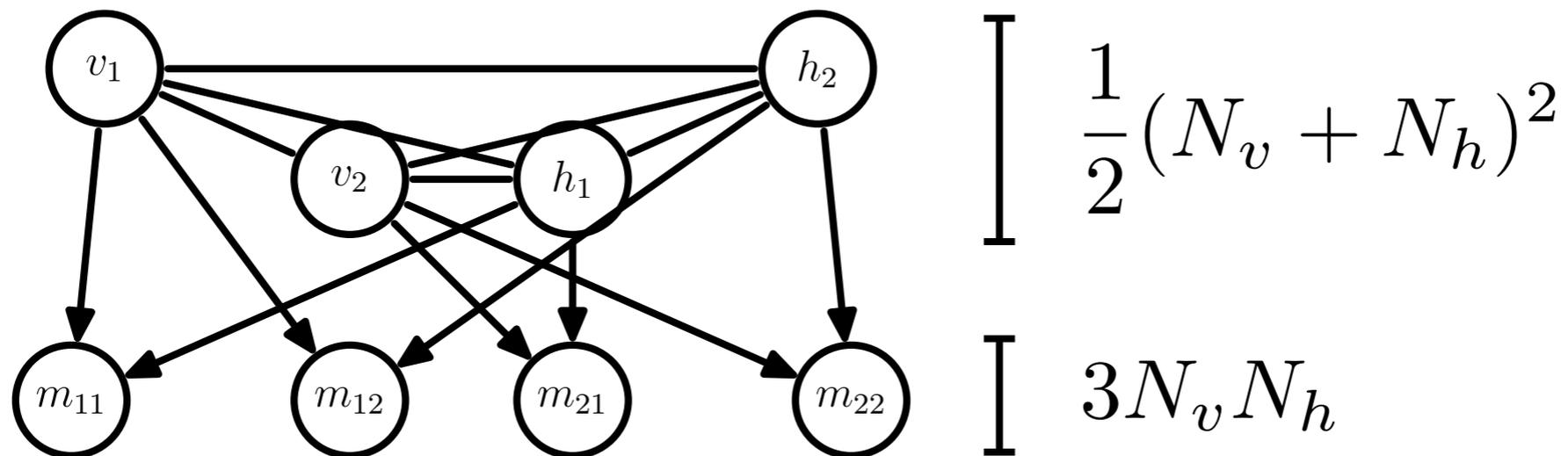


If the variables were binary, this would represent the joint distribution exactly.

The approximation is that we're modeling them as jointly Gaussian.

# Factorized Natural Gradient (FaNG)

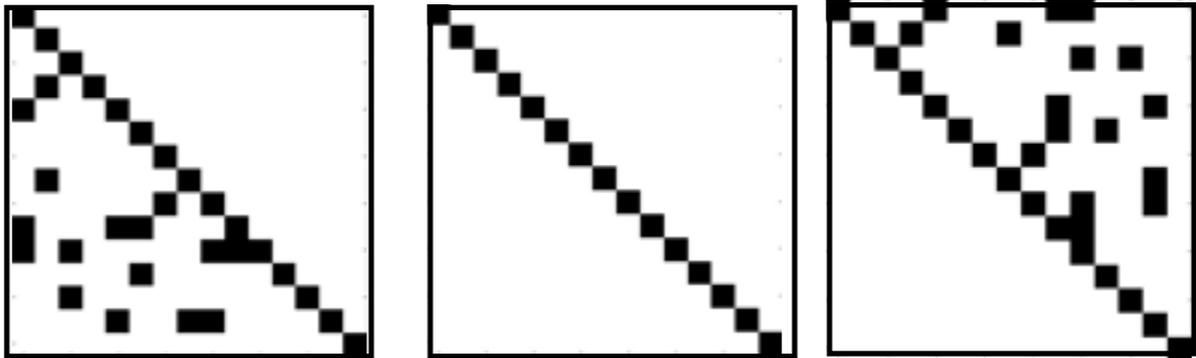
Number of parameters



Compared with  $(N_v + N_h + N_v N_h)^2$  for full G

# Factorized Natural Gradient (FaNG)

Directed model corresponds to sparse Cholesky factorization

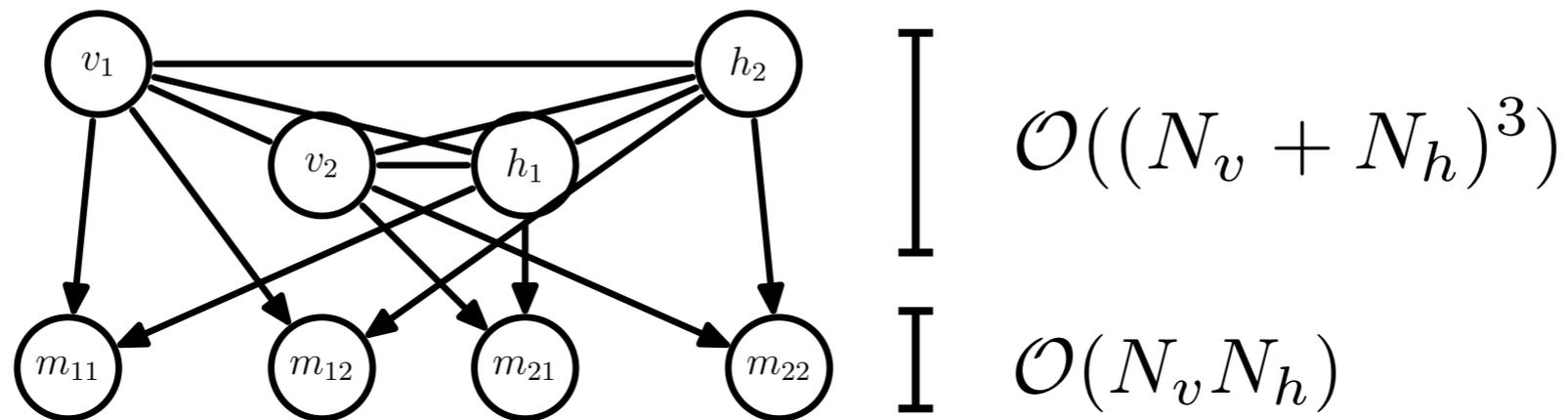
$$G^{-1} \approx L D L^T$$


Natural gradient updates are sparse matrix-vector products

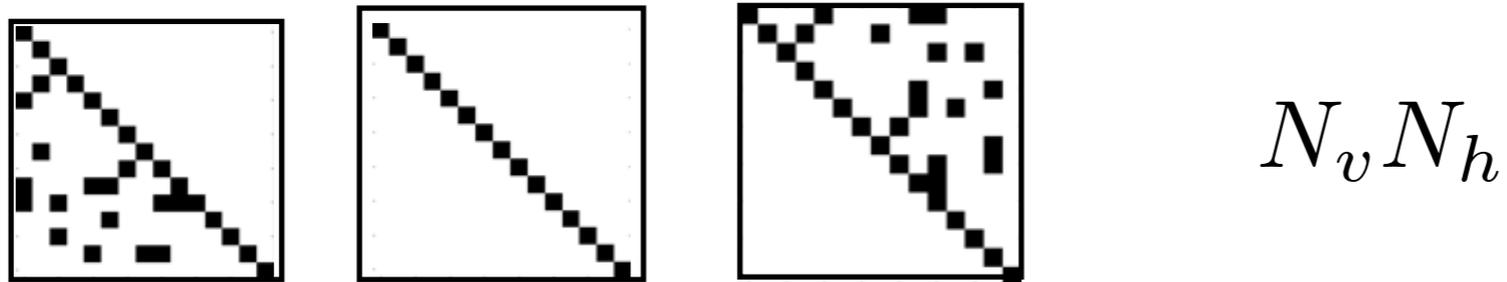
Fit maximum likelihood parameters using linear regression

# Factorized Natural Gradient (FaNG)

Computational cost of fitting parameters (every  $\sim 100$  iterations)



Computational cost of updates



# Explaining the “centering trick”

- Enhanced gradient (Cho et al., 2011); centering trick (Montavon and Muller, 2012)

- Reparameterize the RBM:

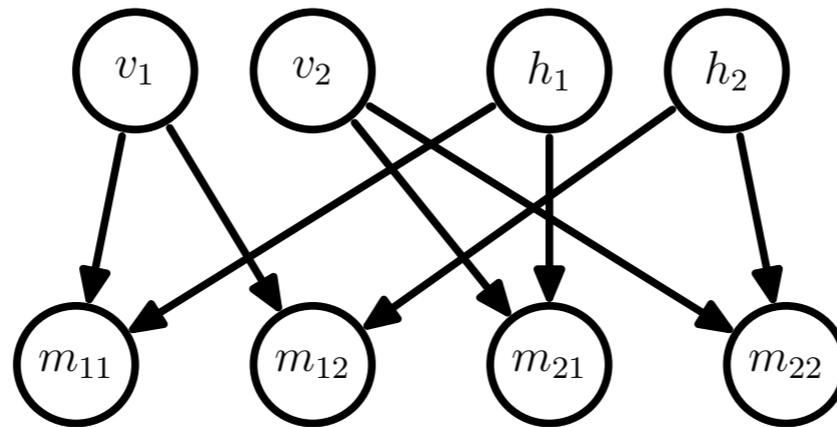
$$\tilde{\mathbf{v}} = \mathbf{v} - \mathbb{E}[\mathbf{v}]$$

$$\tilde{\mathbf{h}} = \mathbf{h} - \mathbb{E}[\mathbf{h}]$$

- Invariant to flipping on/off
- Train deep Boltzmann machines without pre-training (Montavon and Muller, 2012)
- Need this even in addition to metric-free natural gradient (Desjardins et al., 2013)
- But why should this make such a difference?

# Explaining the “centering trick”

- Actually an approximation to natural gradient!

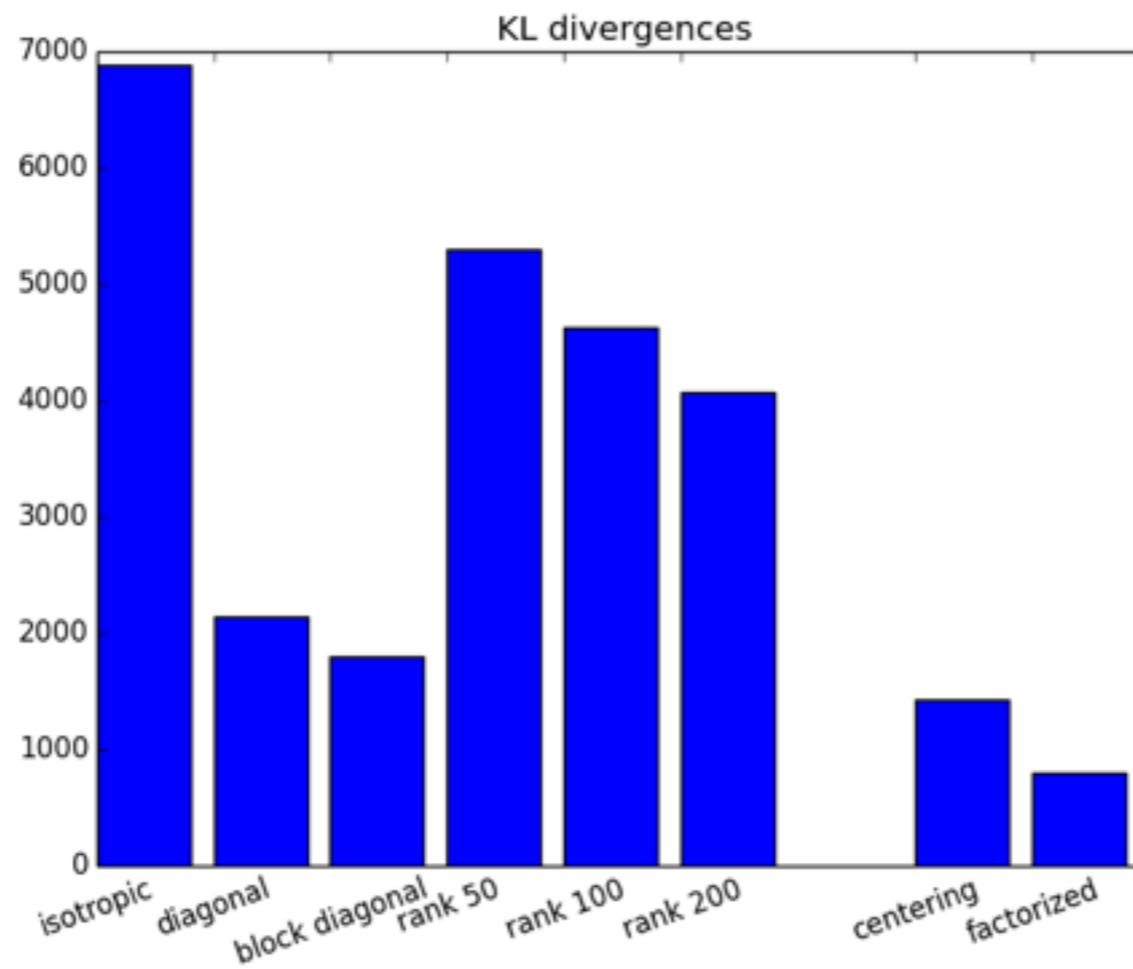


$$m_{ij} | v_i, h_j \sim \mathcal{N} \left( \mathbb{E}[m_{ij}] + \mathbb{E}[v_j](h_j - \mathbb{E}[h_j]) + \mathbb{E}[h_j](v_i - \mathbb{E}[v_i]), \sigma^2 \mathbf{I} \right)$$

- Not the optimal graphical model parameters, but a fairly good guess

# Explaining the “centering trick”

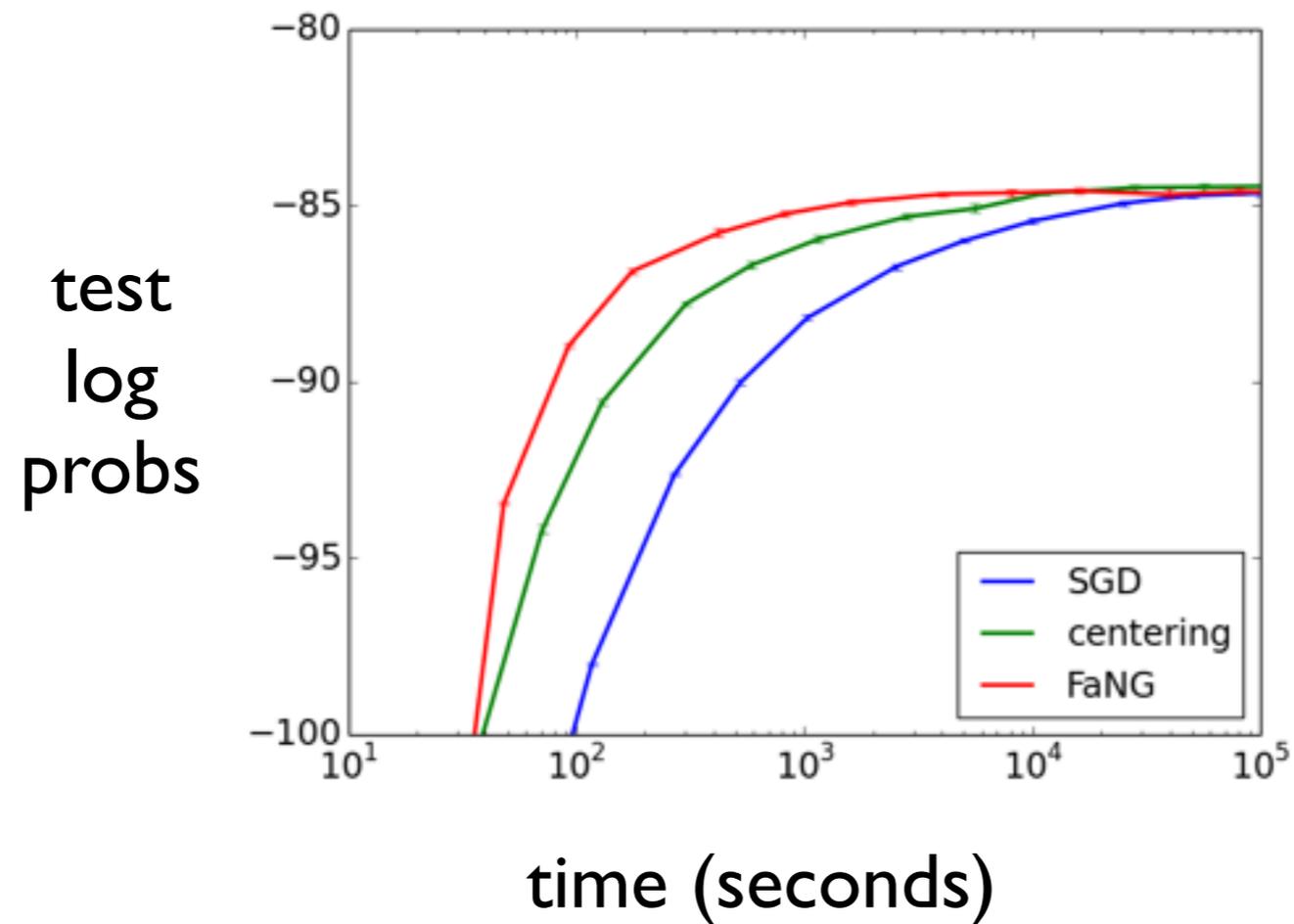
KL divergence to Gaussian with covariance  $G$



# RBM training experiments

## MNIST

(784 × 500)



# Aside: estimating RBM log-likelihoods

- Markov random fields (e.g. RBMs)

$$p(\mathbf{x}) = \frac{f(\mathbf{x})}{\mathcal{Z}}$$
$$\mathcal{Z} = \sum_{\mathbf{x}} f(\mathbf{x})$$

- Evaluating the likelihood requires estimating the intractable  $\mathcal{Z}$

# Aside: estimating RBM log-likelihoods

- Annealed importance sampling (AIS) gives an unbiased estimate of

$$\mathbb{E}[\hat{\mathcal{Z}}_b] = \mathcal{Z}_b$$

- But it gives a biased estimate of

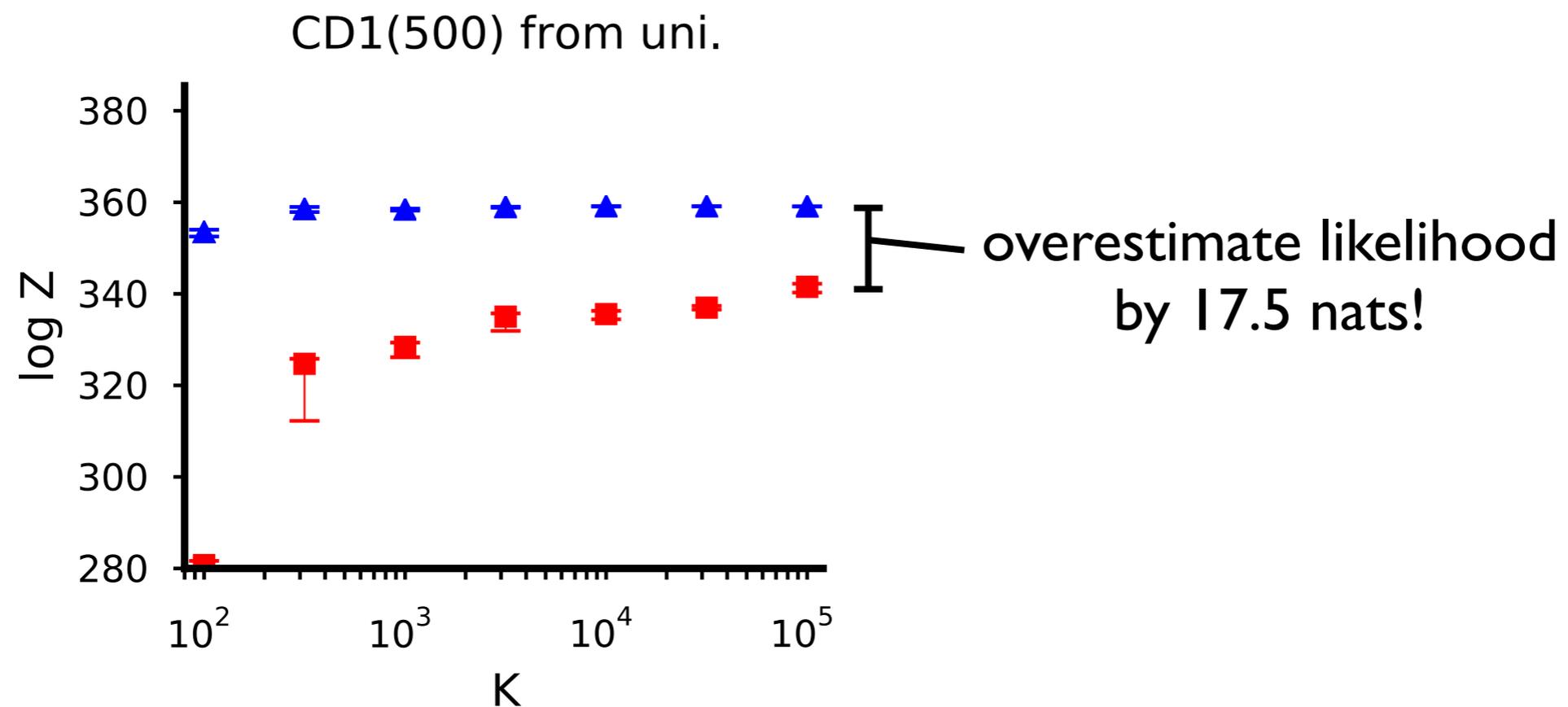
$$\mathbb{E}[\log \hat{\mathcal{Z}}_b] \leq \log \mathcal{Z}_b$$

- Do you have a good model or a bad partition function estimator?

overestimate  $\rightarrow p(\mathbf{x}) = f(\mathbf{x})/\hat{\mathcal{Z}} \leftarrow$  underestimate

# Aside: estimating RBM log-likelihoods

- Is this a problem in practice?



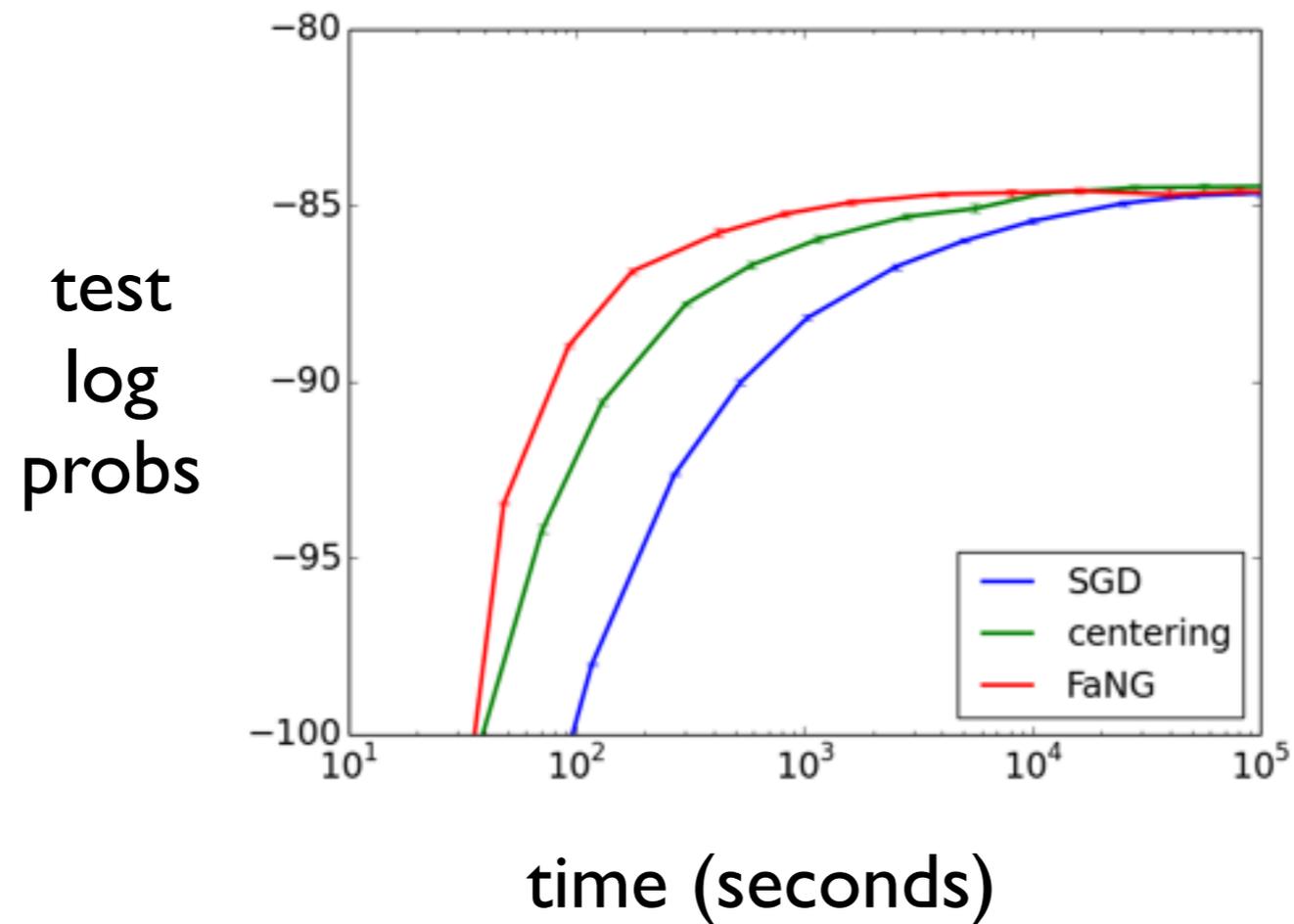
# Aside: estimating RBM log-likelihoods

- Used AIS with geometric averages
  - careful choice of initial distribution: set weights to 0, biases to log odds of data conditional distribution
- Computed confidence intervals using the bootstrap (rather than Gaussian approximation, as is typical)
- Checked that final results were consistent with other estimation methods
  - AIS with moment averaging (Grosse et al., 2013)
  - RAISE (Burda et al., 2014)

# RBM training experiments

## MNIST

(784 × 500)



# RBM training experiments

## Omniglot

(784 x 500)

|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 𐄂 | 𐄃 | 𐄄 | 𐄅 | 𐄆 | 𐄇 | 𐄈 | 𐄉 | 𐄊 | 𐄋 | 𐄌 | 𐄍 | 𐄎 | 𐄏 | 𐄐 |
| 𐄑 | 𐄒 | 𐄓 | 𐄔 | 𐄕 | 𐄖 | 𐄗 | 𐄘 | 𐄙 | 𐄚 | 𐄛 | 𐄜 | 𐄝 | 𐄞 | 𐄟 |
| 𐄠 | 𐄡 | 𐄢 | 𐄣 | 𐄤 | 𐄥 | 𐄦 | 𐄧 | 𐄨 | 𐄩 | 𐄪 | 𐄫 | 𐄬 | 𐄭 | 𐄮 |
| 𐄯 | 𐄰 | 𐄱 | 𐄲 | 𐄳 | 𐄴 | 𐄵 | 𐄶 | 𐄷 | 𐄸 | 𐄹 | 𐄺 | 𐄻 | 𐄼 | 𐄽 |

# RBM training experiments

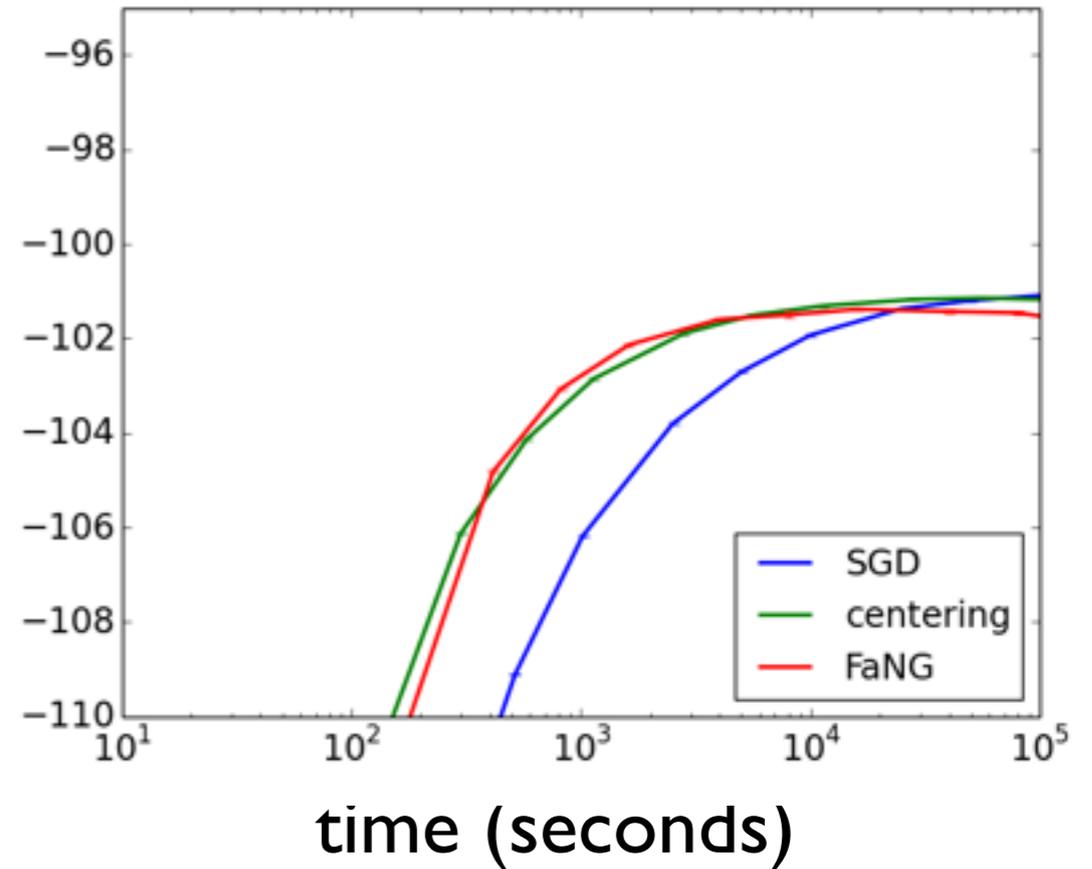
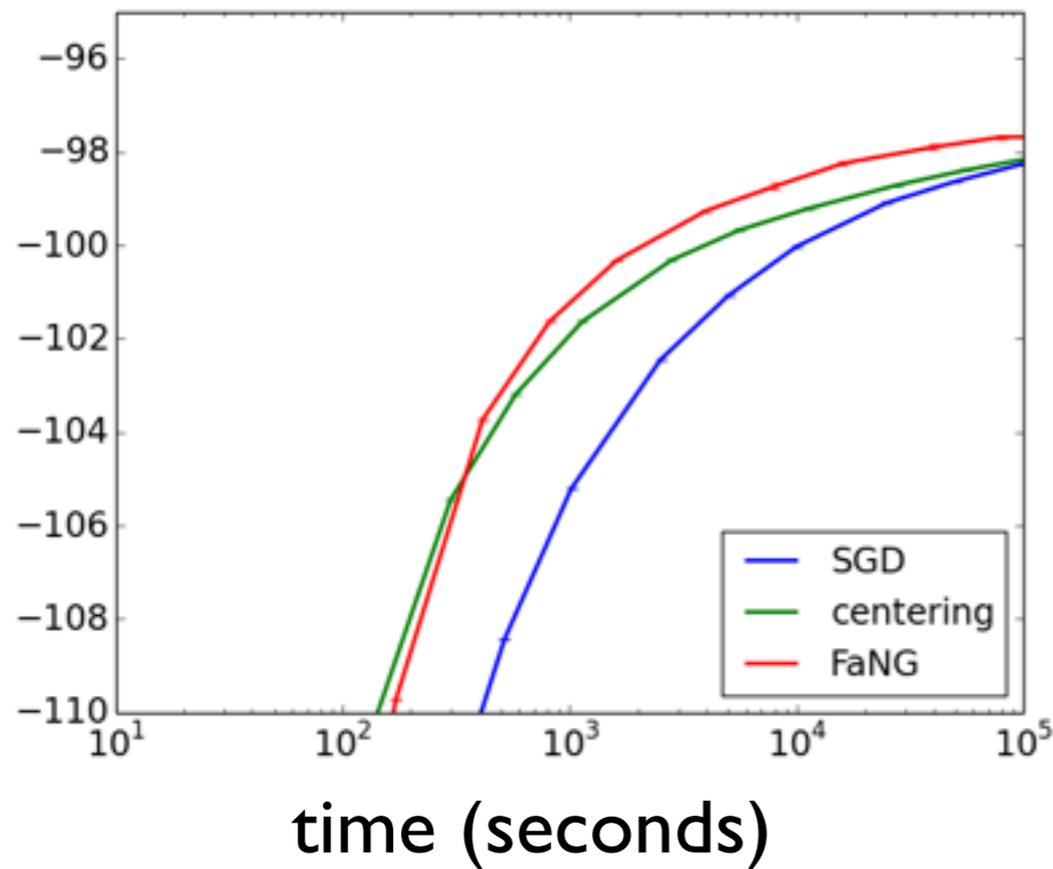
## Omniglot

(784 × 500)

training

test

log  
probs



# Conclusions

- RBM training suffers from nasty curvature
- important covariance information hidden in smaller eigenvalues
- Factorized Natural Gradient: non-iterative approximation
  - helps explain the impressive performance of the centering trick
- more generally, looking at the structure of the model can help us improve optimization