

Introduction

- We consider two classes of models over continuous data: energy based models (EBMs) and feed-forward autoencoders.
- EBMs define a full probabilistic model, while autoencoders admit fast inference and tractable learning.
- Score matching is a statistical estimator designed for fully visible EBMs with intractable partition functions, but applying it to different EBMs often yields learning objectives that are difficult to interpret.
- We generalize score matching to the class of EBMs with latent variables. We further show that these lead to novel autoencoder architectures with comparable empirical performance on a recently proposed EBM.

Models

A latent variable energy based model defines a probability distribution $P_{\theta}(\mathbf{v}, \mathbf{h}) = \frac{\exp(-E_{\theta}(\mathbf{v}, \mathbf{h}))}{Z(\theta)}$ over real data vectors \mathbf{v} with latent variables \mathbf{h} , parameterized by θ . $E_{\theta}(\mathbf{v}, \mathbf{h})$ is called the *energy function* and the distribution is normalized by the *partition function* $Z(\theta)$. We also define the marginal distribution $P_{\theta}(\mathbf{v}) = \frac{\exp(-F_{\theta}(\mathbf{v}))}{Z(\theta)}$ in terms of the *free energy* $F_{\theta}(\mathbf{v}) = -\log(\int_{\mathbf{h}} \exp(-E_{\theta}(\mathbf{v}, \mathbf{h})) d\mathbf{h})$. Some examples include:

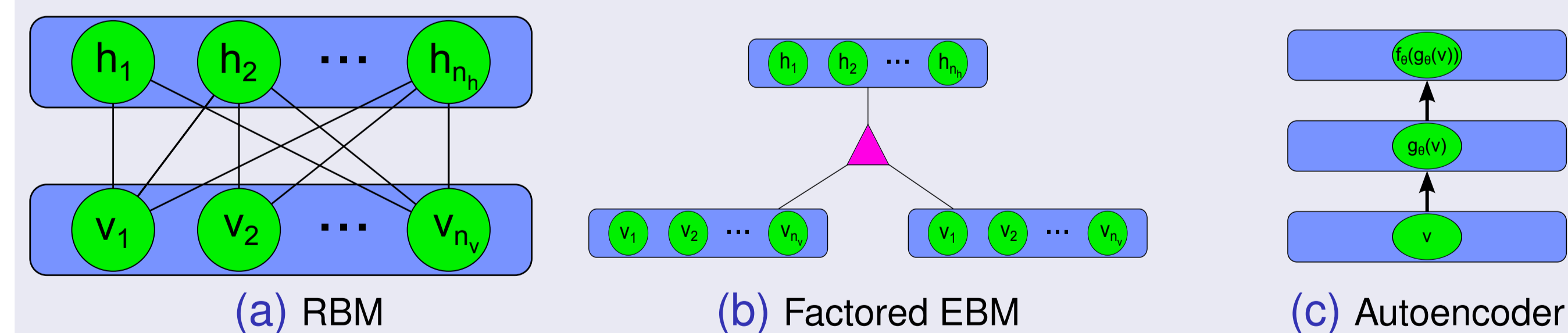
- Gaussian-binary RBMs: $\theta = (\mathbf{W}, \sigma, \mathbf{b}, \mathbf{c})$ and $h_j \in \{0, 1\}$,

$$E_{\theta}(\mathbf{v}, \mathbf{h}) = -\sum_{i=1}^{n_v} \sum_{j=1}^{n_h} \frac{v_i}{\sigma_i} \mathbf{W}_{ij} h_j - \sum_{j=1}^{n_h} \mathbf{b}_j h_j + \frac{1}{2} \sum_{i=1}^{n_v} \frac{(c_i - v_i)^2}{\sigma_i^2}$$

- mPoT: $\theta = (\gamma, \mathbf{W}, \mathbf{C}, \mathbf{b}^v, \mathbf{b}^m)$, $h_j^m \in \{0, 1\}$ and $h_k^c \in \mathbb{R} > 0$,

$$E_{\theta}(\mathbf{v}, \mathbf{h}) = \sum_{k=1}^{n_{hc}} \left[h_k^c \left(1 + \frac{1}{2} \sum_{i=1}^{n_v} \mathbf{C}_{ik} v_i \right) + (1 - \gamma) \log(h_k^c) \right] + \frac{1}{2} \sum_{i=1}^{n_v} v_i^2 - \sum_i \mathbf{b}_i^v v_i - \sum_{i=1}^{n_v} \sum_{j=1}^{n_{hm}} h_j^m \mathbf{W}_{ij} v_i - \sum_{j=1}^{n_{hm}} \mathbf{b}_j^m h_j^m$$

An autoencoder learns to encode its input using an *activation function* $g_{\theta}(\mathbf{v})$ and produces reconstructions by a decoder $f_{\theta}(g_{\theta}(\mathbf{v}))$. Learning consists of setting θ to minimize the reconstruction error $\|f_{\theta}(g_{\theta}(\mathbf{v})) - \mathbf{v}\|^2$.



Score Matching

We assume the data is generated from some unknown distribution $\tilde{p}(\mathbf{v})$. The score matching (SM) estimator is defined as:

$$\mathbf{J}(\theta) = \mathbb{E}_{\tilde{p}(\mathbf{v})} \left[\sum_{i=1}^{n_v} \frac{1}{2} (\psi_i(\mathbf{p}_{\theta}(\mathbf{v})))^2 + \frac{\partial \psi_i(\mathbf{p}_{\theta}(\mathbf{v}))}{\partial v_i} \right]$$

where $\psi_i(\mathbf{p}(\mathbf{v})) = \frac{\partial \log p(\mathbf{v})}{\partial v_i}$ is called the *score function*. Since $\frac{\partial Z(\theta)}{\partial v_i} = \mathbf{0}$, we avoid having to compute the partition function.

Applying and Generalizing Score Matching

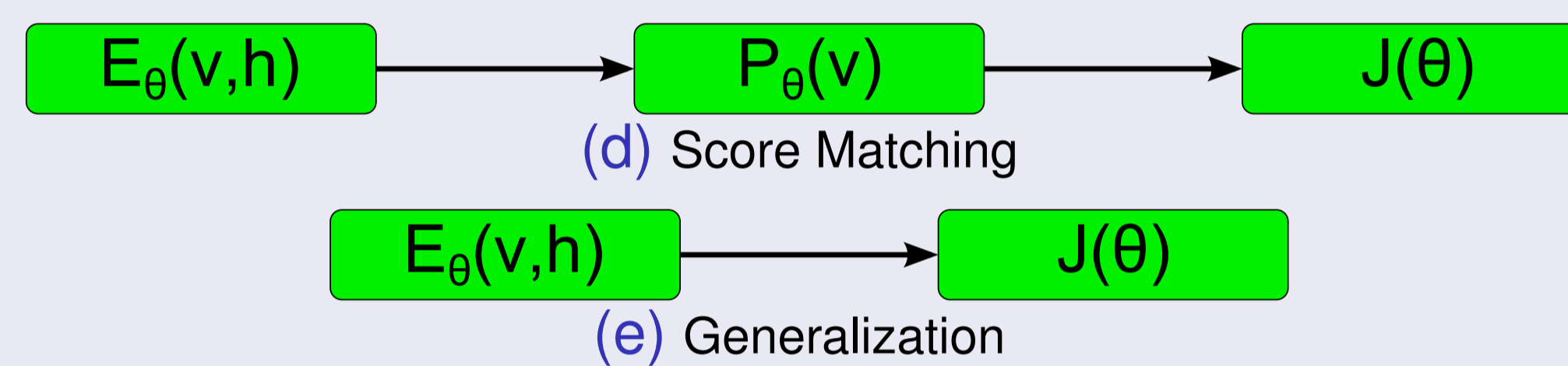
- Score Matching for Gaussian-binary RBMs

$$F_{\theta}(\mathbf{v}) = \frac{1}{2} \sum_{i=1}^{n_v} \frac{(c_i - v_i)^2}{\sigma_i^2} - \sum_{j=1}^{n_h} \log \left(1 + \exp \left(\sum_{i=1}^{n_v} \frac{v_i}{\sigma_i} \mathbf{W}_{ij} + \mathbf{b}_j \right) \right)$$

$$\mathbf{J}(\theta) = \mathbb{E}_{\tilde{p}(\mathbf{v})} \left[\sum_{i=1}^{n_v} \frac{1}{2} \left(\frac{v_i}{\sigma_i^2} - \frac{c_i}{\sigma_i^2} - \sum_{j=1}^{n_h} \frac{\mathbf{W}_{ij}}{\sigma_i} \hat{h}_j \right)^2 - \frac{1}{\sigma_i^2} + \sum_{j=1}^{n_h} \frac{\mathbf{W}_{ij}^2}{\sigma_i^2} \hat{h}_j (1 - \hat{h}_j) \right]$$

where $\hat{h}_j = \text{sigm} \left(\sum_{i=1}^{n_v} \frac{v_i}{\sigma_i} \mathbf{W}_{ij} + \mathbf{b}_j \right)$ and $\text{sigm}(\mathbf{x}) := \frac{1}{1 + \exp(-\mathbf{x})}$

- When $\sigma = 1$, this objective is equivalent to a regularized autoencoder.
- We can also apply score matching to factored models like mPoT.
- When applied to the free energy of a general EBM, the resulting objectives are often difficult to interpret. Instead, we apply score matching without directly marginalizing over the hidden variables.



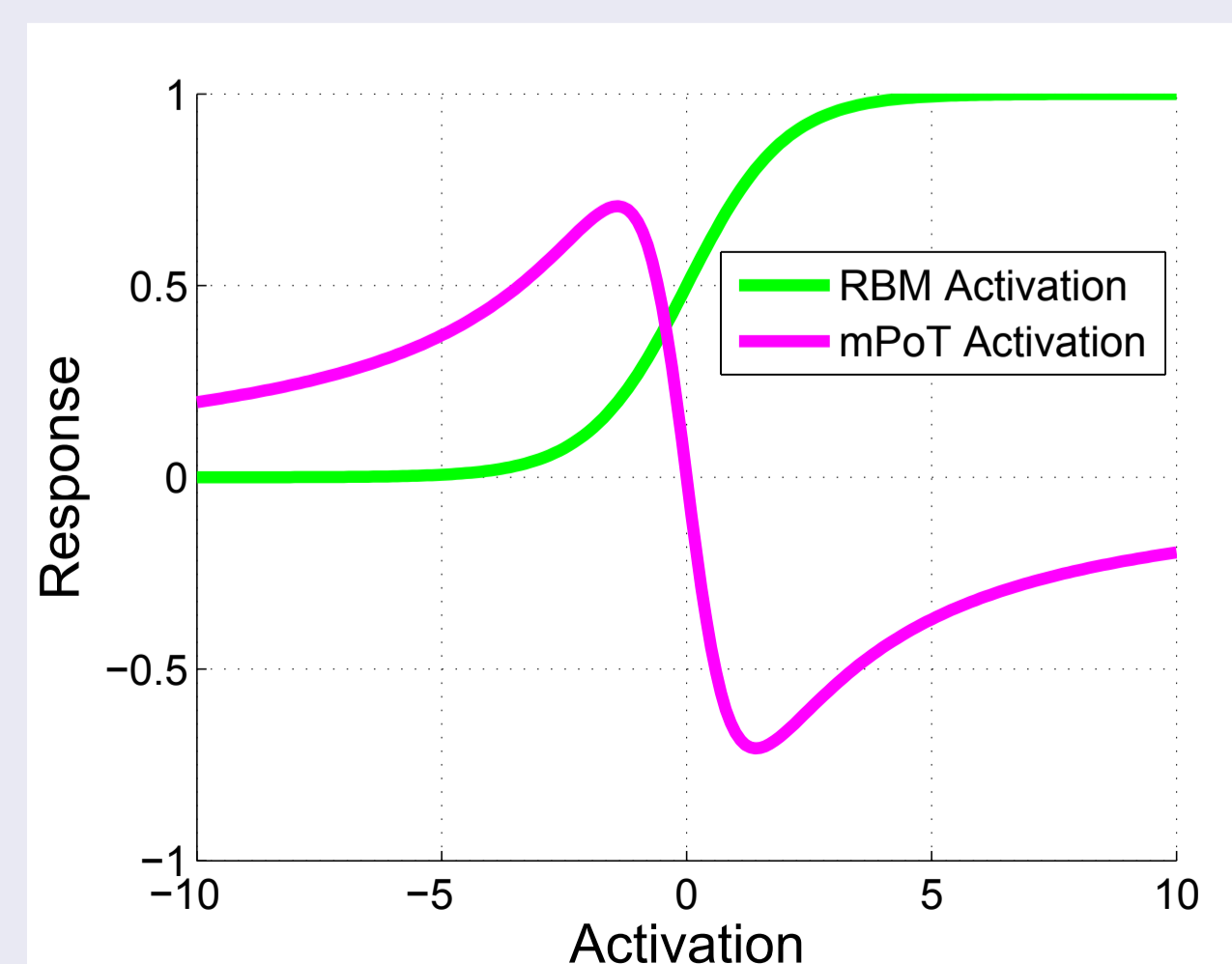
Theorem

The score matching objective for a latent EBM can be expressed succinctly in terms of expectations of the energy with respect to the conditional distribution $P_{\theta}(\mathbf{h}|\mathbf{v})$.

$$\mathbf{J}(\theta) = \mathbb{E}_{\tilde{p}(\mathbf{v})} \left[\sum_{i=1}^{n_v} \frac{1}{2} \left(-\frac{\partial F_{\theta}(\mathbf{v})}{\partial v_i} \right)^2 - \frac{\partial^2 F_{\theta}(\mathbf{v})}{\partial v_i^2} \right]$$

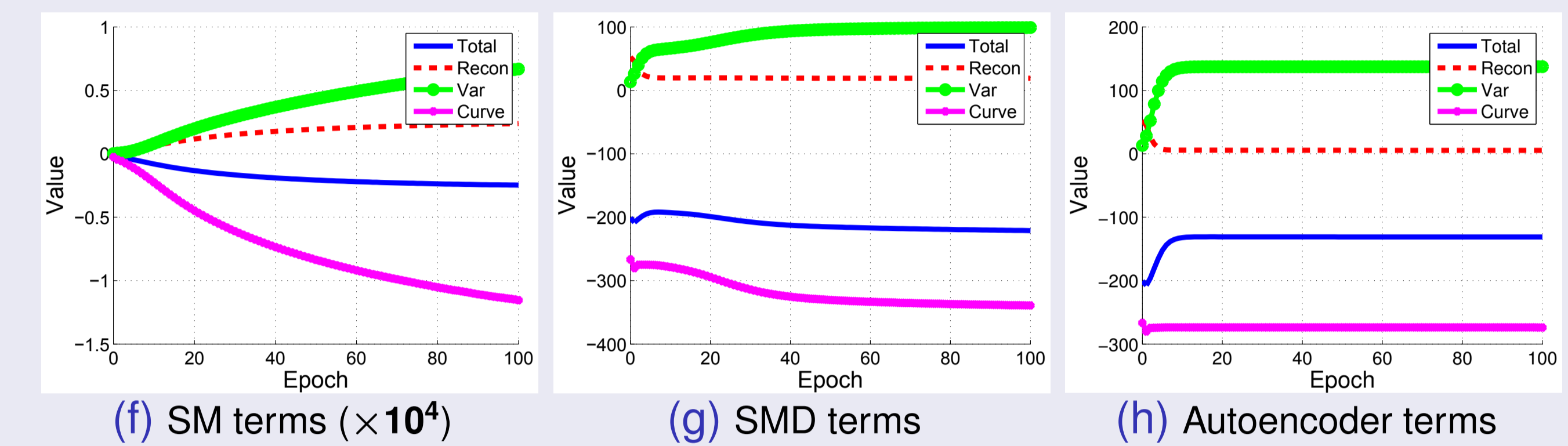
$$= \mathbb{E}_{\tilde{p}(\mathbf{v})} \left[\sum_{i=1}^{n_v} \frac{1}{2} \left(\mathbb{E}_{P_{\theta}(\mathbf{h}|\mathbf{v})} \left[\frac{\partial E_{\theta}(\mathbf{v}, \mathbf{h})}{\partial v_i} \right] \right)^2 + \text{var}_{P_{\theta}(\mathbf{h}|\mathbf{v})} \left[\frac{\partial E_{\theta}(\mathbf{v}, \mathbf{h})}{\partial v_i} \right] - \mathbb{E}_{P_{\theta}(\mathbf{h}|\mathbf{v})} \left[\frac{\partial^2 E_{\theta}(\mathbf{v}, \mathbf{h})}{\partial v_i^2} \right] \right]$$

- Score matching will always minimize expected error and variance, and maximize expected energy curvature.
- When the conditional distribution $P_{\theta}(\mathbf{v}|\mathbf{h})$ over the visible units is Gaussian, the error will be a scaled reconstruction term.
- Score matching can be used to derive new autoencoder models from EBMs that account for covariance structure in \mathbf{v} .
- Applying score matching to factored models leads to autoencoders with peculiar activation functions.

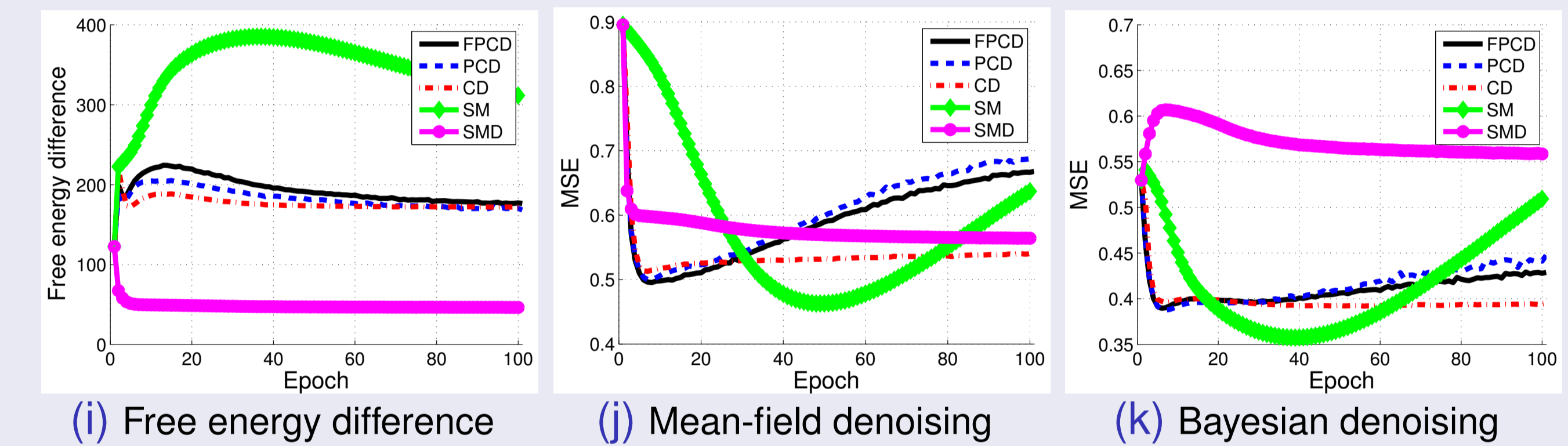


Experiments using mPoT

- We analyze the empirical properties of score matching and compare against several other training methods.
 - CD: Contrastive Divergence
 - PCD: Persistent Contrastive Divergence
 - FPCD: Fast Persistent Contrastive Divergence
 - SM: Score Matching
 - SMD: Denoising Score Matching (Vincent, 2011)
 - AE: mPoT-based autoencoder without regularization
- Objective function analysis



- Denoising



- Learned features



- Classification accuracy on CIFAR 10

CD	PCD	FPCD	SM	SMD	AE
64.6%	64.7%	65.5%	65.0%	64.7%	57.6%

Conclusion

- Applying score matching to a general EBM in terms of the energy function leads to an interpretable objective.
- Score matching for an RBM yields a standard autoencoder, while factored models specify new architectures.
- Regularization, particularly on the curvature of the energy seems to be important for learning good features with autoencoders.
- Methods based on score matching are comparable in performance to stochastic estimators and give tractable objectives.

References

Vincent, P. A connection between score matching and denoising autoencoders. *Neural Computation*, 2011.