

Inverting VAEs for Improved Generative Accuracy

Ian Gemp^{*1}, Mario Parente², Sridhar Mahadevan¹

¹ College of Information and Computer Sciences, University of Massachusetts at Amherst

² Electrical and Computer Engineering, University of Massachusetts at Amherst

*corresponding author: imgemp@cs.umass.edu



Problem: When labeled data (x,y) are scarce, semi-supervised learning can improve model performance by leveraging large amounts of unlabeled data (x) . Under this setting, the deep semi-supervised VAE (M2) learns a generative latent variable model of the data. However, there are cases when this model does not recognize data that *it* generated! How do we encourage the learned model to be internally consistent with desired semantics?

MNIST Inconsistency

Desired semantics (y) : $[0, 1, 2, 3, 4, 8, 9, 5, 7, 6]$

Generate image of 8 with random style: $y \leftarrow [0, 0, 0, 0, 0, 1, 0, 0, 0, 0]$
 $z \sim p(z)$



$x \sim p(x|y, z)$

mismatch!

$y' \sim q(y|x)$

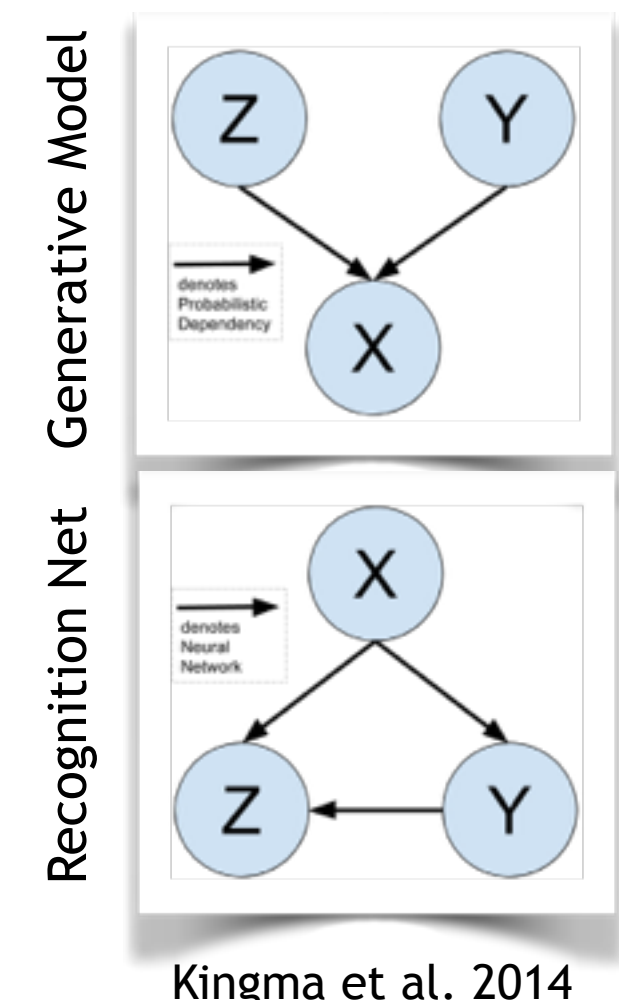
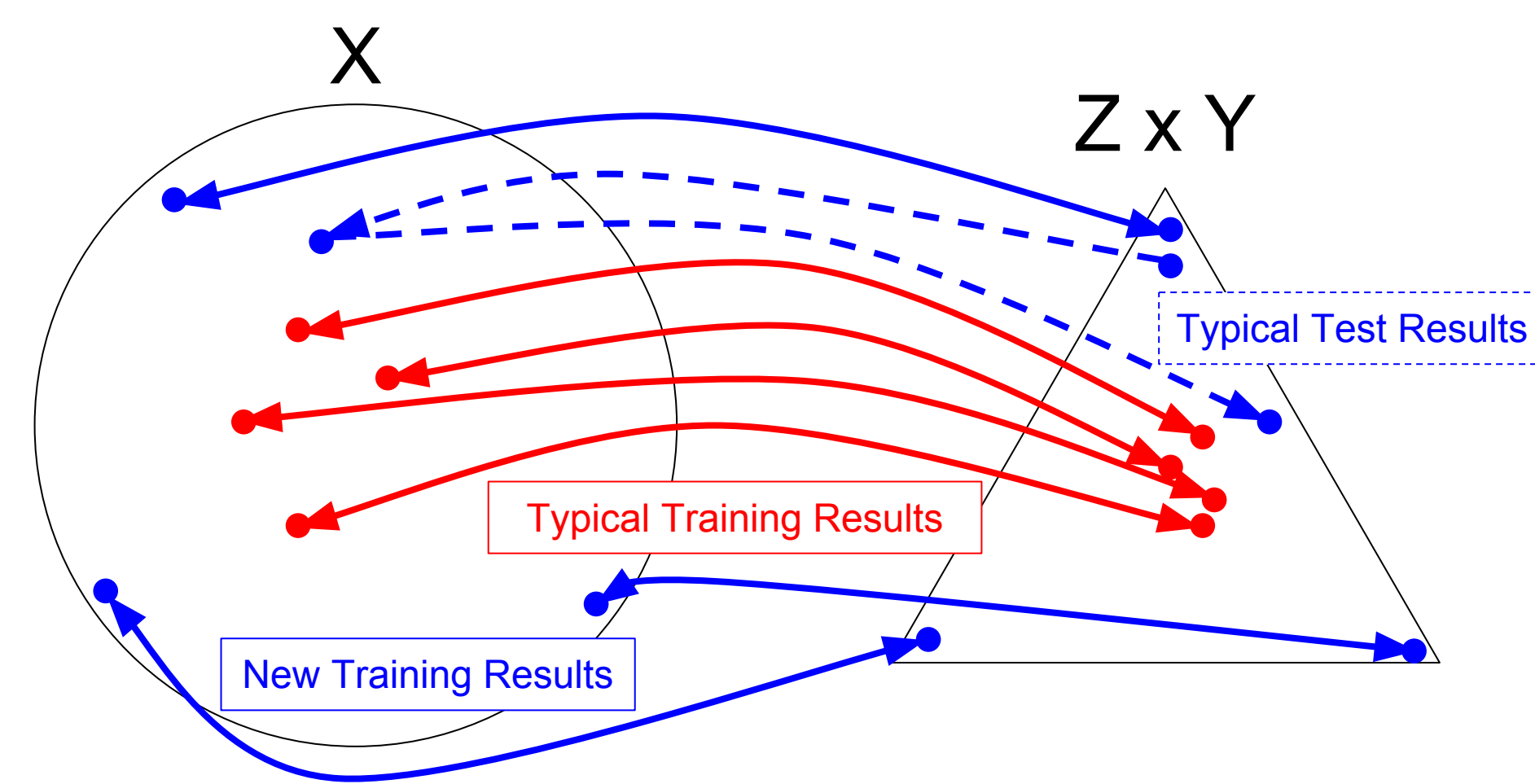
$[0, 0, 0, 0, 0, 0, 1, 0, 0, 0]$

M2
0 : 2 3 4 8 8 5 7 6

0 0 0 0
1 1 1 1
2 2 2 2
3 3 3 3
4 4 4 4
8 8 8 8
8 8 8 8
5 5 5 5
7 7 7 7
6 6 6 6

8 is duplicated!

General Problem



Kingma et al. 2014

Model: To encourage internal consistency, we “invert” the semi-supervised VAE and train under the adversarial scenario exhibited above. We derive the lower bound for maximizing the marginal likelihood of *unfeatured* labels (y) and add this to the standard M2 objective.

$$\log q_\phi(y) \geq \mathbb{E}_{p(z)p_\theta(x|y,z)} \left[\log q_\phi(y, z|x) - \log p_\theta(x|y, z) + \log q(x) - \log p(z) \right] = ELBO_r^u$$

$$ELBO_r^u = \log q_\phi(y) - KL(p_\theta(x, z|y) | q_\phi(x, z|y))$$

Bound becomes tight as standard and “inverted” conditionals become similar

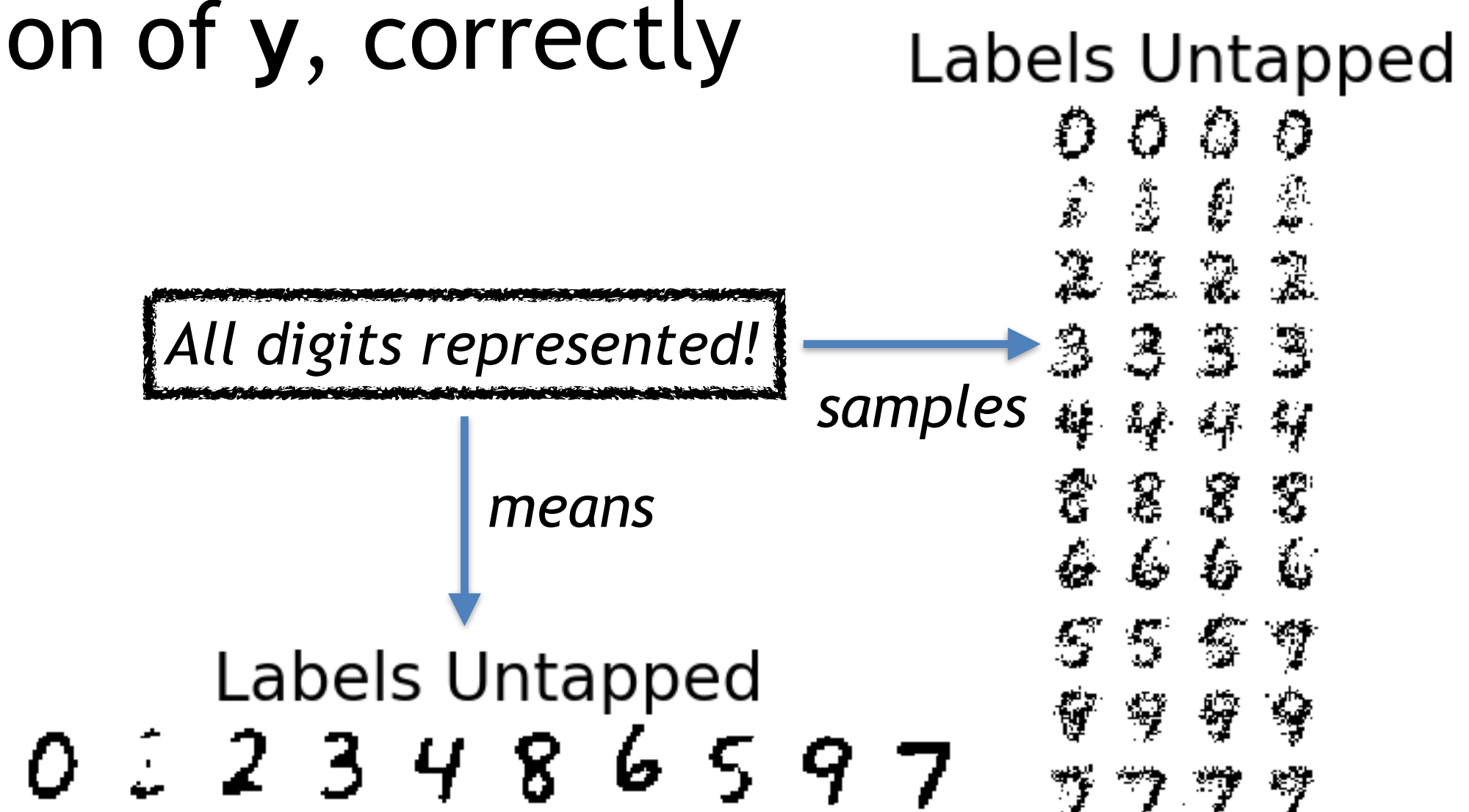
Experiment: We compare the proposed model, *Untapped*, with the original semi-supervised VAE model, M2, by Kingma et al. on the following MNIST task. For both models, labeled data (x,y) is provided for digits 0-4 and unlabeled data (x) for digits 0-9. In addition to these, *Untapped*, leverages *unfeatured* data (y) for digits 0-9 by maximizing the above evidence lower bound. Our hypothesis is that *Untapped* is better at disentangling the semantics of y .

Results: *Untapped* assigns a unique digit to each dimension of y , correctly disentangling semantics.

Improving the generative model results in improved discriminative performance as well.

Closest semantic permutation for digits learned with pretrained logistic regression classifier. Models tested on classifying held-out digits 5-9.

	M2	Untapped
Cross-Entropy	2.30	2.07



Conclusion:

- Improved disentangling of semantics using *Untapped* over M2 VAE model.
- Future work: improving prior $q(x)$ tightens lower bound - improves model performance.