# Semi-supervised Learning with Deep Generative Models
## Kingma et. al. (2014)

Tyler Brown

CS 7180

# Motivating Question

How can we model data of increasing size when obtaining label information is difficult?

# High-level Answer

We can estimate missing label information by using a probabilistic model.

# Specifying the Probabilistic Model for Missing Labels

- Data appears as pairs $(\mathbf{X}, \mathbf{Y}) = \{(\mathbf{x}_1, y_1), ..., (\mathbf{x}_N, y_N)\}$ with the $i$-th observation $x_i \in \mathbb{R}^D$ and a corresponding class label $y_i \in \{1, ..., L\}$
    - Each pair of observations $(x_i, y_i)$ has a corresponding latent variable $z_i$
    - Empirical distribution over the labelled and unabelled subsets is referred to as $\tilde{p}_l(\mathbf{x}, y)$ and $\tilde{p}_u(\mathbf{x})$
- We can estimate $y_i$ for $x_i$ in distribution $\tilde{p}_u(\mathbf{x})$ by finding the maximum probability of $p(y_i)$ by using a set of features related to $z_i$ and a predictive model
    1. **Latent-feature discriminative model (M1)**
    2. **Generative semi-supervised model (M2)**
    3. **Stacked generative semi-supervised model (M1+M2)**

# Bayes Rule is used when specifying M1 & M2

$$p(x, y) = p(x)p(y|x)$$
$$= p(y)p(x|y)$$
$$p(x|y) = \frac{p(x)p(y|x)}{p(y)}$$

for models M1 [1], $p(z|x)$, and M2 [2]; $p(y|x)$

---

[1]Kingma et. al. (2014) equation (1)
[2]Kingma et. al. (2014) equation (2)

# (M1) Latent-feature discriminative model

$$y \Leftarrow p(z|x) = \frac{p(z)p(x|z)}{p(x)}$$

where

$p(z) = \mathcal{N}(z|0, I)$      Gaussian distribution of $z$ given a missing label $y$

$p(x|z) = f(x; z, \theta)$      likelihood function, parameters $\theta$ of a set of $z$

$p(x) = \tilde{p}_u(x)$      unlabelled subset of $x_i \in \mathbb{R}^D$

Kingma et. al. (2014) eq. (1)

# (M1) Predicting Class Labels $y$

Approximate samples from the posterior distribution over the latent variables $p(z|x)$ are used as features to train a classifier that predicts class labels $y$

- ▶ (transductive) SVM
- ▶ multinomial regression

**TODO:** Add pictures or simulation here

# (M2) Generative semi-supervised model

$$p(y|\mathbf{x}) = \frac{p(\mathbf{x}|y)p(y)}{p(\mathbf{x})} \approx \frac{p_\theta(\mathbf{x}|y, \mathbf{z})p(y)}{p(\mathbf{x})}$$

where

$$p(y) = Cat(y|\pi) \quad \text{multinomial distribution, } y \text{ can be latent}$$
$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I}) \quad \text{Gaussian distribution of } z \text{ when missing } y$$
$$p_\theta(\mathbf{x}|y, z) = f(\mathbf{x}; y, \mathbf{z}, \theta) \quad \text{likelihood function, nonlinear parameters}$$
$$p(x) \quad \text{all } x \text{ in dist. of real numbers; } x \in \mathbb{R}^D$$

# Stacked generative semi-supervised model (M1 + M2)

Combine M1 and M2

1. Learn a new latent representation $z_1$ from M1
2. Use embeddings from $z_1$ instead of raw data $x$, to create a generative semi-supervised model M2

**TODO:** Add a picture or something here

# Scaling Up: Lower Bound Objective

Lower Bound Objective[3]: computation of the exact posterior distribution is intractable for models M1 and M2

$$\text{M1: } q_\phi(z|x) = \mathcal{N}(z|u_\phi(x), \text{diag}(\sigma_\phi^2(\mathbf{x}))) \tag{3}$$

$$\text{M2: } q_\phi(\mathbf{z}|y, \mathbf{x}) = \mathcal{N}(\mathbf{z}|\mu_\phi(y, \mathbf{x}), \text{diag}(\sigma_\phi^2(\mathbf{x})));$$

$$q_\phi(y|\mathbf{z}) = \text{Cat}(y|\pi_\phi(x)), \tag{4}$$

where

$$\sigma_\phi(x) \quad \text{vector of standard deviations}$$
$$\pi_\phi(x) \quad \text{probability vector}$$
$$\mu_\phi(x), \sigma_\phi(x), \pi_\sigma(x) \quad \text{Maximum likelihood Priors (MLPs)}$$

---

[3]Kingma et. al. (2014) equations (3), (4)

# Scaling Up: M1 Model Objective

$$\log p_\theta(x) \geq \mathbb{E}$$

# Scaling Up: M2 Model Objective