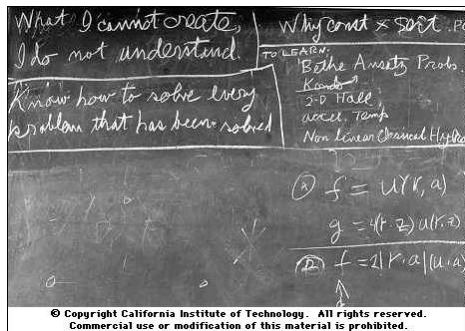


Generative vs Discriminative

2024

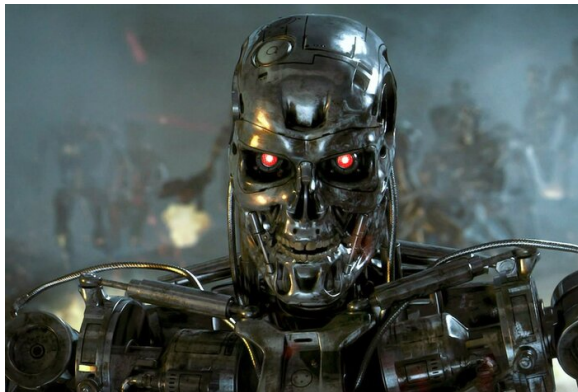
Idea of generative models



Idea of discriminative models



Plato: *"A human is featherless biped"*



Sometimes it's easier to solve a target problem (i.e. classification, regression) than describe the analyzed object nature.

Generative and discriminative models

Discriminative models

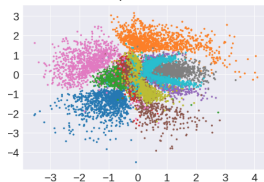
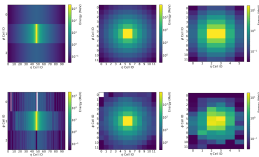
Model: $p(y|x)$.

Generative models

Model: $p(y, x)$.

Why generative models:

- When dataset generation is a target problem
- Synthetic dataset generation
- Latent properties obtaining



Does good likelihood estimation leads to good sampling?

Does good sampling estimation leads to good likelihood estimation?

High likelihood and high sampling quality can be independent (Theis et al., 2015).

- Given a noisy mixture:

$$p_w(x) = 0.01p_{\text{data}}(x) + 0.99p_{\text{noise}}(x), \log p_w(x) \geq \log p_{\text{data}}(x) - \log 100$$

- For another direction: overfitting

Generative models and unsupervised learning

Are the generative models always unsupervised?

Generative models and unsupervised learning

Are the generative models always unsupervised?

No! Linear classification is an example

Logistic regression:

$$E(\mathbf{y} | \mathbf{X}) \equiv g^{-1}(\mathbf{X}\mathbf{w}),$$

$$g^{-1}(x) \frac{e^x}{1 + e^x} \in [0, 1]$$

The decision function is a sigmoid.

Generative model:

$$p(y = 1 | x, w) = \frac{p(x | w, y = 1)p(y = 1)}{\sum_{k=0}^1 p(x | w, y = k)p(y = k)},$$

$$p(x | w, y = k) \sim \mathcal{N}(w_m^k, w_s^k).$$

The decision function is a sigmoid.

Model selection: coherent Bayesian inference

First level: find optimal parameters:

$$\mathbf{w} = \arg \max \frac{p(\mathcal{D}|\mathbf{w})p(\mathbf{w}|\mathbf{h})}{p(\mathcal{D}|\mathbf{h})},$$

Second level: find optimal model:

Evidence:

$$p(\mathcal{D}|\mathbf{h}) = \int_{\mathbf{w}} p(\mathcal{D}|\mathbf{w})p(\mathbf{w}|\mathbf{h})d\mathbf{w}.$$



What is \mathcal{D} for generative and discriminative models? Why?

Plate notation

Plate notation is an alternative visualization for graphical models.

Elements:

- White circles (random variables);
- Grey circles (observed variables);
- Small circles (deterministic values);
- Plates (batching).

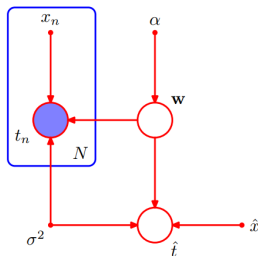
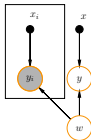


Plate notation for linear regression (Bishop)

Plate notation: discriminative and generative models

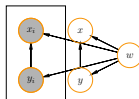
Discriminative models:

- Generate (or deterministically obtain!) x
- Generate w
- Generate $Y \sim p(y|X, w)$



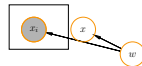
Generative model:

- Generate y
- Generate w
- Generate $x \sim p(X|y, w)$



Generative unsupervised model:

- Generate w
- Generate $x \sim p(X|w)$



Autoencoder: generative model?

(Alain, Bengio 2012): consider regularized autoencoder:

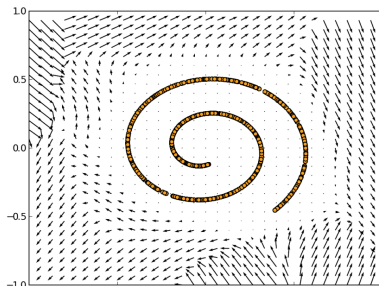
$$||\mathbf{f}(\mathbf{x}, \sigma) - \mathbf{x}||^2,$$

where σ is a noise level.

Then

$$\frac{\partial \log p(\mathbf{x})}{\partial \mathbf{x}} = \frac{||\mathbf{f}(\mathbf{x}, \sigma) - \mathbf{x}||^2}{\sigma^2} + o(1) \text{ when } \sigma \rightarrow 0.$$

Vector field induced by reconstruction error



Diffusion and score matching

Score matching idea:

$$p_{\theta}(\mathbf{x}) = \frac{1}{Z(\theta)} \exp(-E(\mathbf{x}; \theta)).$$

To approximate unknown true distribution $p(\mathbf{x})$ we try to match scores of the distribution:

$$\begin{aligned} & \mathbb{E}_{\mathbf{x} \in p(\mathbf{x})} \left\| \frac{\partial \log p_{\theta}(\mathbf{x})}{\partial \mathbf{x}} - \frac{\partial \log p(\mathbf{x})}{\partial \mathbf{x}} \right\|^2 \rightarrow \min. \\ & \approx \mathbb{E}_{\hat{\mathbf{x}}, \mathbf{x} \in q(\hat{\mathbf{x}}, \mathbf{x})} \left\| \frac{\partial \log p_{\theta}(\mathbf{x})}{\partial \mathbf{x}} \Big|_{\hat{\mathbf{x}}=\mathbf{x}} - \frac{\partial \log p(\hat{\mathbf{x}}|\mathbf{x})}{\partial \hat{\mathbf{x}}} \right\|^2 \rightarrow \min, \end{aligned}$$

where $q(\hat{\mathbf{x}}) = q(\hat{\mathbf{x}}|\mathbf{x})p(\mathbf{x})$.

Discriminative + generative

Naive approach: introduce a prior on class labels

$$p(\mathbf{x}, y | \mathbf{w}) = p(y | \mathbf{w}_y) p(\mathbf{x} | y, \mathbf{w}_x).$$

Two optimization functions:

$$L_G = p(\mathbf{w}) \prod_{\mathbf{x}, y} p(\mathbf{x}, y | \mathbf{w}),$$

$$L_D = p(\mathbf{w}) \prod_{\mathbf{x}, y} p(y | \mathbf{x}, \mathbf{w}).$$

Combine them:

$$\lambda L_G + (1 - \lambda) L_D \rightarrow \max.$$

This optimization is heuristic, it does not give us ML results, nor MAP.

Discriminative + generative

(Bishop et al., 2007): introduce two probabilistic models: “discriminative” and “generative”:

$$p(\mathbf{x}, y | \mathbf{w}_G, \mathbf{w}_D) = p(y | \mathbf{x}, \mathbf{w}_D) p(\mathbf{x} | \mathbf{w}_G) p(\mathbf{w}_G, \mathbf{w}_D), \quad p(\mathbf{x} | \mathbf{w}_G) = \sum_k p(y = k, \mathbf{x} | \mathbf{w}_G).$$

Optimization:

$$p(\mathbf{w}_G, \mathbf{w}_D) \prod_{\mathbf{x}, y} p(y | \mathbf{x}, \mathbf{w}_D) p(\mathbf{x} | \mathbf{w}_G).$$

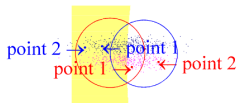
How to select $p(\mathbf{w}_G, \mathbf{w}_D)$?

- $p(\mathbf{w}_G, \mathbf{w}_D) = p(\mathbf{w}_G) p(\mathbf{w}_D)$: obtain L_D ;
- $p(\mathbf{w}_G, \mathbf{w}_D) = p(\mathbf{w}_G) \delta(\mathbf{w}_G - \mathbf{w}_D)$: obtain L_G ;
- Trade-off: $p(\mathbf{w}_G, \mathbf{w}_D) \propto p(\mathbf{w}_G) p(\mathbf{w}_D) \exp(-\frac{1}{2\sigma^2} \|\mathbf{w}_G - \mathbf{w}_D\|^2)$.

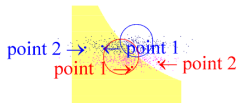
Discriminative + generative

(Bishop et al., 2007): example of different combinations of these optimizations for the synthetic dataset. The dataset contains only 2 labeled objects for each class.

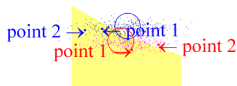
$\alpha = 0$ - generative case



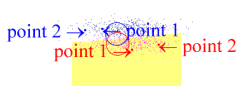
$\alpha = 0.4$



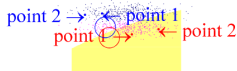
$\alpha = 0.6$



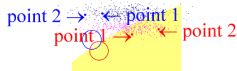
$\alpha = 0.7$



$\alpha = 0.8$



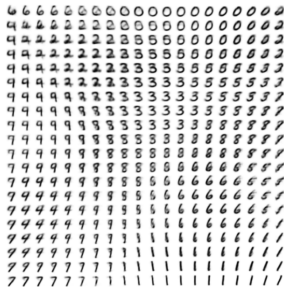
$\alpha = 1$ - discriminative case



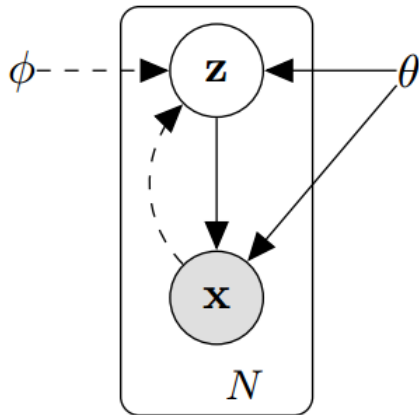
VAE: generation process



(a) Learned Frey Face manifold



(b) Learned MNIST manifold



Semi-supervised VAE (Kingma et al., 2014)

$$\text{M1: } q_{\phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_{\phi}(\mathbf{x}), \text{diag}(\boldsymbol{\sigma}_{\phi}^2(\mathbf{x}))), \quad (3)$$

$$\text{M2: } q_{\phi}(\mathbf{z}|y, \mathbf{x}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_{\phi}(y, \mathbf{x}), \text{diag}(\boldsymbol{\sigma}_{\phi}^2(\mathbf{x}))); \quad q_{\phi}(y|\mathbf{x}) = \text{Cat}(y|\boldsymbol{\pi}_{\phi}(\mathbf{x})), \quad (4)$$

For this model, we have two cases to consider. In the first case, the label corresponding to a data point is observed and the variational bound is a simple extension of equation (5):

$$\log p_{\theta}(\mathbf{x}, y) \geq \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}, y)} [\log p_{\theta}(\mathbf{x}|y, \mathbf{z}) + \log p_{\theta}(y) + \log p(\mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x}, y)] = -\mathcal{L}(\mathbf{x}, y), \quad (6)$$

For the case where the label is missing, it is treated as a latent variable over which we perform posterior inference and the resulting bound for handling data points with an unobserved label y is:

$$\begin{aligned} \log p_{\theta}(\mathbf{x}) &\geq \mathbb{E}_{q_{\phi}(y, \mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|y, \mathbf{z}) + \log p_{\theta}(y) + \log p(\mathbf{z}) - \log q_{\phi}(y, \mathbf{z}|\mathbf{x})] \\ &= \sum_y q_{\phi}(y|\mathbf{x}) (-\mathcal{L}(\mathbf{x}, y)) + \mathcal{H}(q_{\phi}(y|\mathbf{x})) = -\mathcal{U}(\mathbf{x}). \end{aligned} \quad (7)$$

The bound on the marginal likelihood for the entire dataset is now:

$$\mathcal{J} = \sum_{(\mathbf{x}, y) \sim \tilde{p}_l} \mathcal{L}(\mathbf{x}, y) + \sum_{\mathbf{x} \sim \tilde{p}_u} \mathcal{U}(\mathbf{x}) \quad (8)$$

Semi-supervised VAE (Kingma et al., 2014)

$$\mathcal{J}^\alpha = \mathcal{J} + \alpha \cdot \mathbb{E}_{\tilde{p}_l(\mathbf{x}, y)} [-\log q_\phi(y|\mathbf{x})], \quad (9)$$

Algorithm 1 Learning in model M1

```
while generativeTraining() do
   $\mathcal{D} \leftarrow \text{getRandomMiniBatch}()$ 
   $\mathbf{z}_i \sim q_\phi(\mathbf{z}_i|\mathbf{x}_i) \quad \forall \mathbf{x}_i \in \mathcal{D}$ 
   $\mathcal{J} \leftarrow \sum_n \mathcal{J}(\mathbf{x}_i)$ 
   $(\mathbf{g}_\theta, \mathbf{g}_\phi) \leftarrow (\frac{\partial \mathcal{J}}{\partial \theta}, \frac{\partial \mathcal{J}}{\partial \phi})$ 
   $(\theta, \phi) \leftarrow (\theta, \phi) + \Gamma(\mathbf{g}_\theta, \mathbf{g}_\phi)$ 
end while
while discriminativeTraining() do
   $\mathcal{D} \leftarrow \text{getLabeledRandomMiniBatch}()$ 
   $\mathbf{z}_i \sim q_\phi(\mathbf{z}_i|\mathbf{x}_i) \quad \forall \{\mathbf{x}_i, y_i\} \in \mathcal{D}$ 
   $\text{trainClassifier}(\{\mathbf{z}_i, y_i\})$ 
end while
```

Algorithm 2 Learning in model M2

```
while training() do
   $\mathcal{D} \leftarrow \text{getRandomMiniBatch}()$ 
   $y_i \sim q_\phi(y_i|\mathbf{x}_i) \quad \forall \{\mathbf{x}_i, y_i\} \notin \mathcal{O}$ 
   $\mathbf{z}_i \sim q_\phi(\mathbf{z}_i|y_i, \mathbf{x}_i)$ 
   $\mathcal{J}^\alpha \leftarrow \text{eq. (9)}$ 
   $(\mathbf{g}_\theta, \mathbf{g}_\phi) \leftarrow (\frac{\partial \mathcal{L}^\alpha}{\partial \theta}, \frac{\partial \mathcal{L}^\alpha}{\partial \phi})$ 
   $(\theta, \phi) \leftarrow (\theta, \phi) + \Gamma(\mathbf{g}_\theta, \mathbf{g}_\phi)$ 
end while
```

Semi-supervised VAE (Kingma et al., 2014)

Table 1: Benchmark results of semi-supervised classification on MNIST with few labels.

| N | NN | CNN | TSVM | CAE | MTC | AtlasRBF | M1+TSVM | M2 | M1+M2 |
|------|-------|-------|-------|-------|-------|---------------------|----------------------|----------------------|----------------------------|
| 100 | 25.81 | 22.98 | 16.81 | 13.47 | 12.03 | 8.10 (± 0.95) | 11.82 (± 0.25) | 11.97 (± 1.71) | 3.33 (± 0.14) |
| 600 | 11.44 | 7.68 | 6.16 | 6.3 | 5.13 | – | 5.72 (± 0.049) | 4.94 (± 0.13) | 2.59 (± 0.05) |
| 1000 | 10.7 | 6.45 | 5.38 | 4.77 | 3.64 | 3.68 (± 0.12) | 4.24 (± 0.07) | 3.60 (± 0.56) | 2.40 (± 0.02) |
| 3000 | 6.04 | 3.35 | 3.45 | 3.22 | 2.57 | – | 3.49 (± 0.04) | 3.92 (± 0.63) | 2.18 (± 0.04) |



(a) Handwriting styles for MNIST obtained by fixing the class label and varying the 2D latent variable z .

Generative and discriminative models hierarchy

Can we reduce generative model to discriminative?

Can we reduce discriminative model to generative?

Basic Concept of EBM¹

Energy-Based Models (EBMs) capture dependencies by associating a scalar energy to each configuration of the variables

Learning: finding an energy function that associates low energies to correct values, and higher energies to incorrect values

Inference: finding parameter that minimize the energy

Pros: no requirement for proper normalization, in comparison with probabilistic models

Probability density

$p(\mathbf{x})$ for $\mathbf{x} \in \mathbb{R}^D$: $p_{\theta}(\mathbf{x}) = \frac{\exp(-E_{\theta}(\mathbf{x}))}{Z(\theta)}$, where $E_{\theta}(\mathbf{x})$ is energy function and $Z(\theta) = \int_{\mathbf{x}} \exp(-E_{\theta}(\mathbf{x})) d\mathbf{x}$ is normalizing constant or partition function

¹See talk of Maria Kovaleva, 2022

Learning of EBM

Estimate $Z(\theta)$ can be challenging for most energy functions. Derivative of the log-likelihood for a single example \mathbf{x} with respect to θ : $\frac{\partial \log p_{\theta}(\mathbf{x})}{\partial \theta} = \mathbb{E}_{p_{\theta}(\mathbf{x}')} \left[\frac{\partial E_{\theta}(\mathbf{x}')}{\partial \theta} \right] - \frac{\partial E_{\theta}(\mathbf{x})}{\partial \theta}$

Learning

Approximate the expectation in derivative of the log-likelihood using a sampler based on Stochastic Gradient Langevin Dynamics (SGLD):

$$\mathbf{x}_0 \sim p_0(\mathbf{x}), \mathbf{x}_{i+1} = \mathbf{x}_i - \frac{\alpha}{2} \frac{\partial E_{\theta}(\mathbf{x}_i)}{\partial \mathbf{x}_i} + \epsilon, \epsilon \sim \mathcal{N}(0, \alpha),$$

where $p_0(\mathbf{x})$ is typically a Uniform distribution, α - step-size (should be decayed following a polynomial schedule).

Practically α and ϵ is often chosen separately leading to a biased sampler which allows for faster training

In classifiers

Classification: D parameters, K classes, parametric function, $f_\theta : \mathbb{R}^D \rightarrow \mathbb{R}^K$ return logits

Categorical distribution parameterized by f_θ via Softmax transfer function:

$$p_\theta(y|\mathbf{x}) = \frac{\exp(f_\theta(\mathbf{x})[y])}{\sum_{y'} \exp(f_\theta(\mathbf{x})[y'])}$$

Reinterpretation of the logits

Let's define an energy based model with $E_{\theta}(\mathbf{x}, y) = -f_{\theta}(\mathbf{x})[y]$ and unknown normalizing constant Z_{θ} :

$$p_{\theta}(\mathbf{x}, y) = \frac{\exp(f_{\theta}(\mathbf{x})[y])}{Z(\theta)}$$

Then

$$p_{\theta}(\mathbf{x}) = \sum_y p_{\theta}(\mathbf{x}, y) = \frac{\sum_y \exp(f_{\theta}(\mathbf{x})[y])}{Z(\theta)}$$

Conclusion

The $\text{LogSumExp}(\cdot)$ of the logits of any classifier can be re-used to define the energy function as $E_{\theta}(\mathbf{x}) = -\text{LogSumExp}(f_{\theta}(\mathbf{x})[y]) = -\log \sum_y \exp f_{\theta}(\mathbf{x})[y]$

A generative model hidden within every standard discriminative model!

Optimization

$$\log p_{\theta}(\mathbf{x}, y) = \log p_{\theta}(\mathbf{x}) + \log p_{\theta}(y|\mathbf{x})$$

Where:

1. $p(y|\mathbf{x})$ optimized using standard cross-entropy
2. $\log p(\mathbf{x})$ optimized with SGLD where gradients are taken with respect to $\text{LogSumExp}(f_{\theta}(\mathbf{x})[y])$ as it described for EBMs

Remarks

1. We don't know Z_{θ} , because of it we use SLGD sampler
2. The estimator of $\log p(\mathbf{x})$ will be biased when using a MCMC sampler with a finite number of steps

Classifier as generative model: results

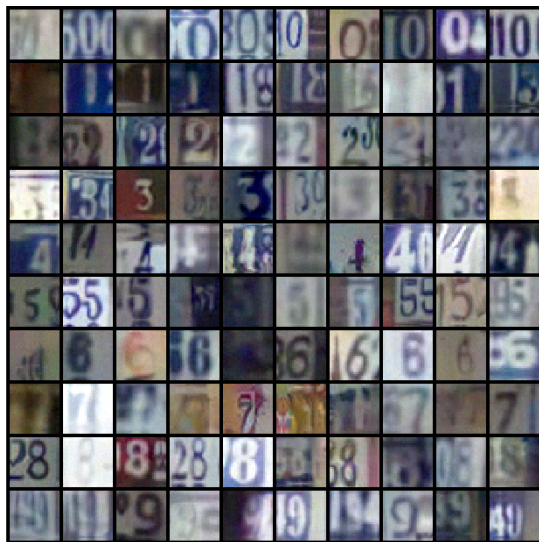
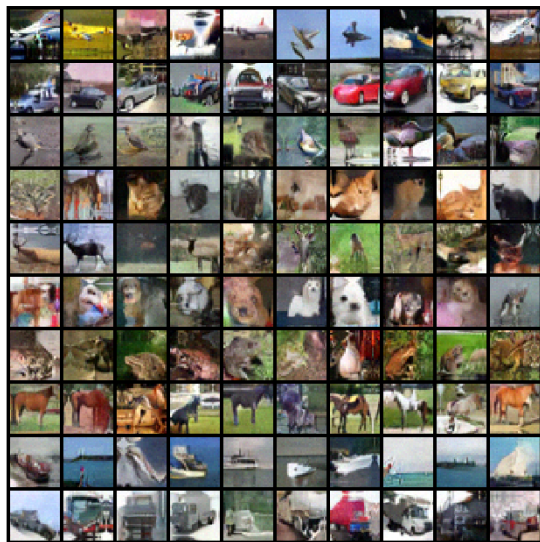


Figure 7: Class-conditional Samples. Left to right: CIFAR10, SVHN.

Dataset shift

Dataset shift is an event when distribution $p(\mathbf{X}, \mathbf{y})$ significantly differ for the training and test/inference phases.

- Covariate shift — difference in $p(\mathbf{X})$
- Prior probability shift — difference in $p(\mathbf{y})$
- Concept shift — difference in $p(\mathbf{y}|\mathbf{X})$

Dataset shift

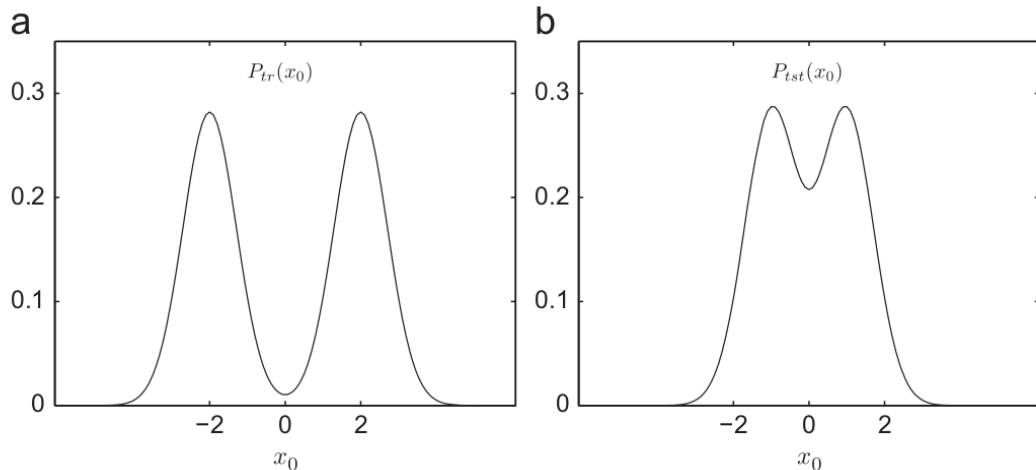
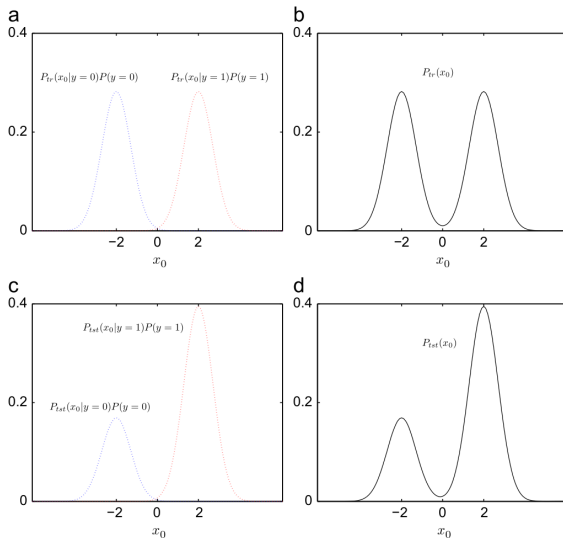
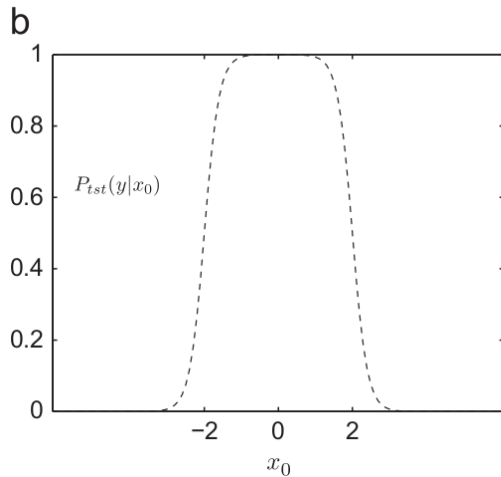
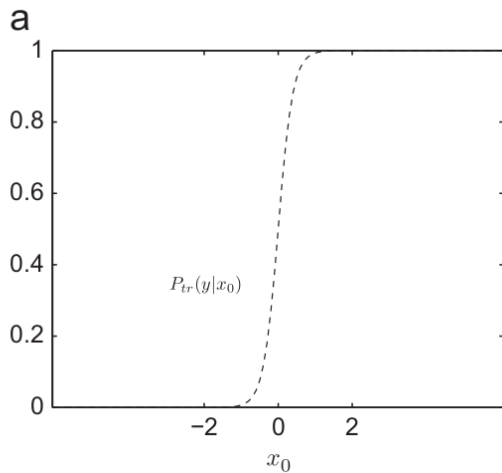


Fig. 1. Covariate shift: $P_{tst}(y|x_0) = P_{tr}(y|x_0)$ and $P_{tr}(x_0) \neq P_{tst}(x_0)$. (a) Training data and (b) test data.

Dataset shift



Dataset shift



Moreno-Torres et al., 2012

Model selection problem: recap

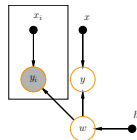
First level: find optimal parameters:

$$\mathbf{w} = \arg \max \frac{p(\mathcal{D}|\mathbf{w})p(\mathbf{w}|\mathbf{h})}{p(\mathcal{D}|\mathbf{h})},$$

Second level: find optimal model:

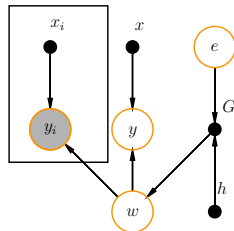
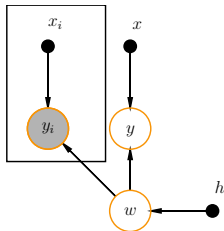
Evidence:

$$p(\mathcal{D}|\mathbf{h}) = \int_{\mathbf{w}} p(\mathcal{D}|\mathbf{w})p(\mathbf{w}|\mathbf{h})d\mathbf{w}.$$



Model selection problem: recap

Can we generate target models parameters using a generative model?



Model selection: hybrid approach

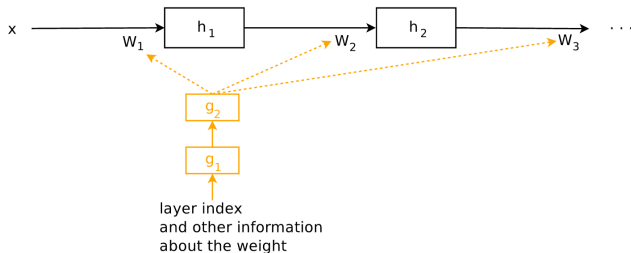
Definition

Given a set Λ .

Hypernetwork is a parametric mapping from Λ to set \mathbb{R}^n of the model \mathbf{f} parameters:

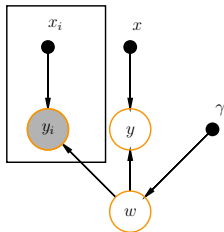
$$G : \Lambda \times \mathbb{R}^u \rightarrow \mathbb{R}^n,$$

where \mathbb{R}^u is a set of hypernetwork parameters.

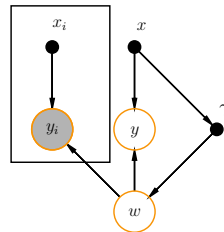


Model selection: discriminative approach

$$\mathbf{w}_{\text{MOE}} = \langle \gamma(\mathbf{x}), [\mathbf{w}_1, \dots, \mathbf{w}_n] \rangle$$



Model generation scheme



MOE optimization as a discriminative model

References

- Bishop C. M. Pattern recognition //Machine learning. – 2006. – Т. 128. – №. 9.
- Генератор котиков: <https://github.com/aleju/cat-generator>
- Paganini M., de Oliveira L., Nachman B. Accelerating science with generative adversarial networks: an application to 3D particle showers in multilayer calorimeters //Physical review letters. – 2018. – Т. 120. – №. 4. – С. 042003.
- Antoran J., Miguel A. Disentangling and learning robust representations with natural clustering //2019 18th IEEE International Conference On Machine Learning And Applications (ICMLA). – IEEE, 2019. – С. 694-699.
- LDA lectures: https://personal.utdallas.edu/~nrr150130/cs6347/2017sp/lects/Lecture_18_LDA.pdf
- Bernardo J. M. et al. Generative or discriminative? getting the best of both worlds //Bayesian statistics. – 2007. – Т. 8. – №. 3. – С. 3-24.
- Гребенькова О. С., Бахтеев О. Ю., Стрижов В. В. Вариационная оптимизация модели глубокого обучения с контролем сложности //Информатика и её применения. – 2021. – Т. 15. – №. 1. – С. 42-49.
- Ha D., Dai A., Le Q. V. Hypernetworks //arXiv preprint arXiv:1609.09106. – 2016.
- Адуенко А. А. Выбор мультимodelей в задачах классификации : дис. – Федер. исслед. центр "Информатика и управление" РАН, 2017.
- Grathwohl W. et al. Your classifier is secretly an energy based model and you should treat it like one //arXiv preprint arXiv:1912.03263. – 2019.
- Theis L., Oord A., Bethge M. A note on the evaluation of generative models //arXiv preprint arXiv:1511.01844. – 2015.
- Moreno-Torres, Jose G., et al. "A unifying view on dataset shift in classification."Pattern recognition 45.1 (2012): 521-530.

Organizational issues

- Technical meeting: moving to 26th of November, one more week, **USE IT WISELY**.
- General requirements for library and algorithms:
 - ▶ Codestyle: at least PEP-8
 - ▶ The code must be commented
 - ▶ There must be an installation file (setup.py or alternative)
 - ▶ There must be requirements file with defined versions
- Demo: must be available at the meeting (perfect case: ipynb or collab)
- Tests:
 - ▶ Coverage 75%;
 - ▶ Tests during build: not required
- Documentation: must be ready and deployed
- Blog-post: pre-final version must be ready