

PML 9. Generative Models: A Guide Into the literature

Probabilistic Machine Learning Reading Group

David Rios Insua

March 11 2026

ICMAT-CSIC

- Generative Models: Basic Concepts
- Variational autoencoders
- Autoregressive models
- Final comments

Generative Models: Basic Concepts

- ML: Generative vs Discriminative
- GM. Joint probability distribution for $p(\mathbf{x})$ for $\mathbf{x} \in \mathcal{X}$.
- Conditional GM. Conditional on inputs $\mathbf{c} \in \mathcal{C}$, $p(\mathbf{x}|\mathbf{c})$.
1 correct y vs several correct y

Types

- Deep generative models. Use NN to learn mapping from latent z to data x
- Probabilistic graphical models. See slides Session 1. Maps interconnected latents z_1, \dots, z_N to observed x_1, \dots, x_D
- All kinds of 'oldies': STATS-101, mixture models, KDEs,...

- Generating data (Generative AI)
<https://llms-cunef-icmat-rg2024.github.io/>
 - Sample new images.
 - Synthetic data.
 - Conditional generation
 - c = text prompt, x = image, text-2-image
 - image-2-text (captioning), image-2-image (colorization), speech-2-text (ASR), sequence-2-sequence (machine translation, text generation).
- Density estimation $p(x)$
 - Outlier detection, Novelty detection, Distribution shift, mixture modeling,...

- Multiple imputation $p(x_m|x_o)$
- Structure discovery $p(z|x) \propto p(z)p(x|z)$
- Model-based reinforcement learning, when data is not abundant
- Data compression
- Decision support, Decision making

Unsupervised learning...

- Variational Autoencoders
- Autoregressive models
- Normalizing flows (sess 10)
- Energy models (sess 11)
- Diffusion models (sess 12)
- Generative adversarial networks

Murphy's chs 20, 21, 22 (and other along the slides)

- Density available. Explicit, only implicit, lower bound on density, approximation.
- Sampling. Fast(VAE) or slow (ARM).
- Training. MLE (local, global), MAP, posterior, min-max,...
- Presence of latents. Yes, No, Compressed.
- Imposed architecture restrictions (eg invertible neural nets)

Metrics capturing

- Sample quality: Samples generated by model part of data distribution?
- Sample diversity: Samples capture all modes?
- Generalization. Model generalizes beyond training data?

- KL divergence between the model q and the actual p

$$D_{KL}(p||q) = \int p(x) \log \frac{p(x)}{q(x)} dx = -H(p) + H_{ce}(p, q)$$

Approach $p(x)$ by empirical distribution, cross entropy through NLL on dataset $-\frac{1}{N} \sum_n \log q(x_n)$

- Two sample tests. Null: Both sample from same distribution. Classical tests adapted to high dimension.
- Human evaluation, e.g. in images.
-

- Murphy (2023) PML, Chap.20
- Tomczak (2022) Deep GMs
- Gallego, DRI (2022) Recent developments in NNs.

Variational autoencoders

- Definition
- Autoencoders
- Basics
- Variants

VAE. Definition

GMs with latent variables z , a prior over z , a deep NN decoder $d_\theta(z)$ for the x , and an exponential family transformation

$$z \sim p_\theta(z)$$
$$x|z \sim \text{Expfam}(x|d_\theta(z))$$

Expfam designates a distribution of the exponential family (Gaussian, product of Bernoullis,....)

Deep latent variable model (DLVM).

When prior Gaussian, Deep latent Gaussian model (DLGM)

Computing $p_\theta(z|x)$ (and $p_\theta(x)$) is computationally intractable...

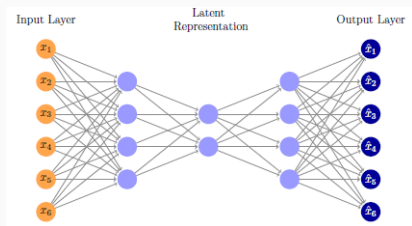
Hence VI through $q_\phi(z|x)$, recognition or inference network, to do approximate posterior inference (Sess 2, 5)

Autoencoder

A NN model that maps inputs x to a low-dimensional latent space using an encoder, $z = f_e(x)$, and then attempts to reconstruct the inputs using a decoder, $\hat{x} = f_d(z)$.

Trained to minimize $L(\theta) = \|f_d(f_e(x)) - x\|^2$.

Once trained, ignore reconstructed \hat{x} and use z as required.



A 'nonlinear PCA'

VAEs are probabilistic autoencoders

- AE objective function log likelihood, no KL term
- Encoding is deterministic. Encoder network computes $E(z|x)$ not $V(z|x)$

VAE. Basics

A basic VAE defines $p_\theta(x, z) = p_\theta(z)p_\theta(x|z)$ with $p_\theta(z)$ Gaussian; $p_\theta(x|z)$ product of ExpFam distributions with parameters defined through a NN decoder $d_\theta(z)$. For example for binary observations, use

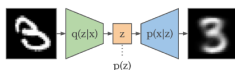
$$p_\theta(x|z) = \prod_{d=1}^D \text{Ber}(x_d | \sigma(d_\theta(z)))$$

Also fits a recognition model to perform approximate posterior inference

$$q_\phi(z, x) = q(z|e_\phi(x)) \simeq p_\theta(z|x)$$

A typical example $q_\phi(z|x) \sim N(z|\mu, \text{diag}(\exp(l)))$

$$(\mu, l) = e_\phi(x)$$



Recall Stochastic VI (session 2)

Data:

\mathcal{D} : Dataset

$q_{\phi}(z|x)$: Inference model

$p_{\theta}(x, z)$: Generative model

Result:

θ, ϕ : Learned parameters

$(\theta, \phi) \leftarrow$ Initialize parameters

while *SGD not converged* **do**

$\mathcal{M} \sim \mathcal{D}$ (Random minibatch of data)

$\epsilon \sim p(\epsilon)$ (Random noise for every datapoint in \mathcal{M})

 Compute $\tilde{\mathcal{L}}_{\theta, \phi}(\mathcal{M}, \epsilon)$ and its gradients $\nabla_{\theta, \phi} \tilde{\mathcal{L}}_{\theta, \phi}(\mathcal{M}, \epsilon)$

 Update θ and ϕ using SGD optimizer

end

Amortized Stochastic VI

Algorithm 21.1: Fitting a VAE with Bernoulli likelihood and full covariance Gaussian posterior. Based on Algorithm 2 of [KW19a].

```

1 Initialize  $\theta, \phi$ 
2 repeat
3   Sample  $x \sim p_{\mathcal{D}}$ 
4   Sample  $\epsilon \sim q_0$ 
5    $(\mu, \log \sigma, \mathbf{L}') = e_{\phi}(x)$ 
6    $\mathbf{M} = \text{np.triu}(\text{np.ones}(K), -1)$ 
7    $\mathbf{L} = \mathbf{M} \odot \mathbf{L}' + \text{diag}(\sigma)$ 
8    $z = \mathbf{L}\epsilon + \mu$ 
9    $p_{\theta} = d_{\theta}(z)$ 
10   $\mathcal{L}_{\log qz} = -\sum_{k=1}^K [\frac{1}{2}\epsilon_k^2 + \frac{1}{2}\log(2\pi) + \log \sigma_k]$  // from  $q_{\phi}(z|x)$  in Equation (10.47)
11   $\mathcal{L}_{\log pz} = -\sum_{k=1}^K [\frac{1}{2}z_k^2 + \frac{1}{2}\log(2\pi)]$  // from  $p_{\theta}(z)$  in Equation (10.48)
12   $\mathcal{L}_{\log px} = -\sum_{d=1}^D [x_d \log p_d + (1 - x_d) \log(1 - p_d)]$  // from  $p_{\theta}(x|z)$ 
13   $\mathcal{L} = \mathcal{L}_{\log px} + \mathcal{L}_{\log pz} - \mathcal{L}_{\log qz}$ 
14  Update  $\theta := \theta - \eta \nabla_{\theta} \mathcal{L}$ 
15  Update  $\phi := \phi - \eta \nabla_{\phi} \mathcal{L}$ 
16 until converged

```

- Beta-VAE

$$\mathcal{L}_\beta(\theta, \phi|x) = \underbrace{-\mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x|z)]}_{\mathcal{L}_R} + \beta \underbrace{D_{\text{KL}}(q_\phi(z|x) \parallel p_\theta(z))}_{\mathcal{L}_R}$$

- Info-VAE

$$\mathbb{I}(\theta, \phi|x) = -\lambda D_{\text{KL}}(q_\phi(z) \parallel p_\theta(z)) - \mathbb{E}_{q_\phi(z)}[D_{\text{KL}}(q_\phi(x|z) \parallel p_\theta(x|z))] + \alpha \mathbb{I}_q(x; z)$$

- Multimodal VAE. Several modes of data (e.g. text and images)
- Hierarchical VAE. With several stochastic layers
-

- Murphy (2023) PML, Chap.21
- Kingma, Welling (2013) Auto-encoding variational Bayes
- Kingma, Welling (2019) An introduction to variational autoencoders

Autoregressive models

- Basic definitions
- NNs for autoregressions
- Transformers

Autoregressive model (ARM), using chain rule

$$p(x_{1:T}) = p(x_1)p(x_2|x_1)p(x_3|x_2, x_1)\dots p(x_T|x_1, \dots, x_{T-1})$$

Conditioning on c

$$p(x_{1:T}|c) = p(x_1|c)p(x_2|x_1, c)p(x_3|x_2, x_1, c)\dots p(x_T|x_1, \dots, x_{T-1}, c)$$

Taming the beast...

- (First order) Markov assumption

$$p(x_t|x_1, \dots, x_{t-1}) = p(x_t|x_{t-1}) \text{ AR}(1)$$

- (Higher order) Markov assumptions.
- x_t depends on z_t , which compresses the past.
 - If z_t is deterministic function of x_{t-1} RNN
 - If z_t is stochastic function of z_{t-1} HMM

Use a functional form for the conditionals

$$p(y_1, y_2, \dots, y_N) = p(y_1, \dots, y_q) \prod_{t=q+1}^N p(y_t | y_{t-1}, y_{t-2}, \dots, y_{t-q})$$

Nonlinear AR(q) model

$$y_t = f(y_{t-1}, y_{t-2}, \dots, y_{t-q}) + \epsilon_t, \quad t = q+1, \dots, N$$
$$\epsilon_t \sim N(0, \sigma^2),$$

Linear AR(q)+Nonlinear AR(q)

$$f(y_{t-1}, y_{t-2}, \dots, y_{t-q}) = x_t' \lambda + \sum_{j=1}^M \beta_j \varphi(x_t' \gamma_j), \quad t = q+1, \dots, N$$

like a (shallow) NN

$$x_t = (y_{t-1}, \dots, y_{t-q}):$$

$$\hat{y}_t(x_t) = \sum_{j=1}^M \beta_j \varphi(x_t' \gamma_j) + \delta_j$$

Priors

$$\beta_j \sim N(\mu_\beta, \sigma_\beta^2), \quad \lambda \sim N(\mu_\lambda, \sigma_\lambda^2 I)$$

$$\gamma_j \sim N(\mu_\gamma, \Sigma_\gamma), \quad \sigma^2 \sim \text{InvGamma}(a_\sigma, b_\sigma).$$

$$\mu_\beta \sim N(a_{\mu_\beta}, b_{\mu_\beta}), \quad \sigma_\beta^2 \sim \text{InvGamma}(a_{\sigma_\beta}, b_{\sigma_\beta})$$

$$\mu_\lambda \sim N(a_{\mu_\lambda}, b_{\mu_\lambda}), \quad \sigma_\lambda^2 \sim \text{InvGamma}(a_{\sigma_\lambda}, b_{\sigma_\lambda})$$

$$\mu_\gamma \sim N(a_{\mu_\gamma}, b_{\mu_\gamma}), \quad \Sigma_\gamma \sim \text{InvWishart}(a_{\Sigma_\lambda}, b_{\Sigma_\lambda})$$

NNs for autoregression. Sampler

1. Given current values of χ and σ^2 (β and λ are marginalized), for each γ_j , $j = 1, \dots, M$, generate a proposal $\tilde{\gamma}_j \sim N(\gamma_j, c\Sigma_\gamma)$, calculate the acceptance probability,

$$a = \min \left[1, \frac{p(\tilde{\gamma} | \mu_\gamma, \Sigma_\gamma) p(D' | y_1, \dots, y_q, \tilde{\gamma}, \sigma^2)}{p(\gamma | \mu_\gamma, \Sigma_\gamma) p(D' | y_1, \dots, y_q, \gamma, \sigma^2)} \right]$$

where $\gamma = (\gamma_1, \dots, \gamma_{j-1}, \gamma_j, \gamma_{j+1}, \dots, \gamma_M)$, and $\tilde{\gamma} = (\gamma_1, \dots, \gamma_{j-1}, \tilde{\gamma}_j, \gamma_{j+1}, \dots, \gamma_M)$. With probability a , replace γ by $\tilde{\gamma}$ and rearrange indices if necessary to satisfy order constraint. Otherwise, leave γ_j unchanged.

2. Generate new values for parameters, drawing from their full conditional posteriors:

$\tilde{\beta} \sim p(\beta | D', \gamma, \lambda, \sigma^2, \chi)$ is a multivariate normal distribution.

$\tilde{\lambda} \sim p(\lambda | D', \gamma, \beta, \sigma^2, \chi)$ is a multivariate normal distribution.

$\tilde{\sigma}^2 \sim p(\sigma^2 | D', \gamma, \beta, \lambda)$ is an inverse gamma distribution.

3. Given current values of $(\gamma, \beta, \lambda, \sigma^2)$, generate a new value for each hyperparameter by drawing from their complete conditional posterior distributions:

$\tilde{\mu}_\beta \sim p(\mu_\beta | D', \beta, \sigma_\beta^2)$ is a normal distribution.

$\tilde{\sigma}_\beta^2 \sim p(\sigma_\beta^2 | D', \beta, \mu_\beta)$ is an inverse gamma distribution.

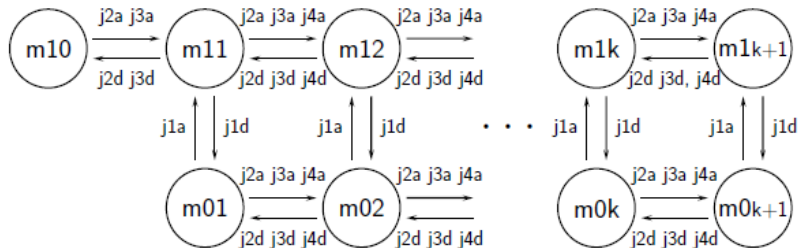
$\tilde{\mu}_\lambda \sim p(\mu_\lambda | D', \lambda, \sigma_\lambda^2)$ is a multivariate normal distribution.

$\tilde{\sigma}_\lambda^2 \sim p(\sigma_\lambda^2 | D', \lambda, \mu_\lambda)$ is an inverse gamma distribution.

$\tilde{\mu}_\gamma \sim p(\mu_\gamma | D', \gamma, \Sigma_\gamma)$ is a multivariate normal distribution.

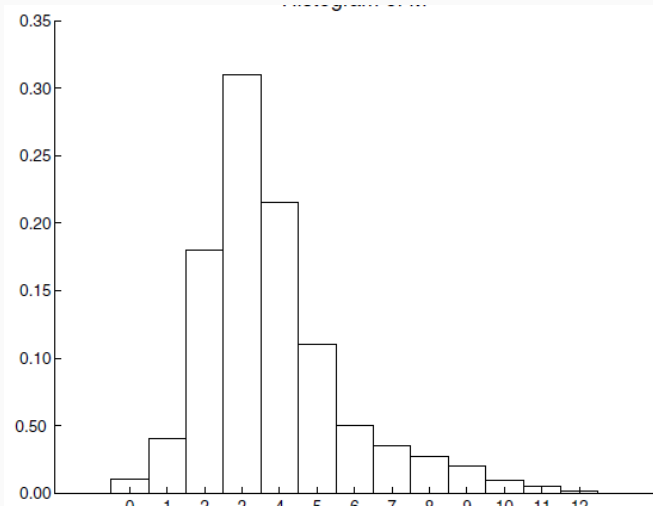
$\tilde{\Sigma}_\gamma \sim p(\Sigma_\gamma | D', \gamma, \mu_\gamma)$ is an inverse Wishart distribution.

NNs for autoregression. RJ



- Add or delete linear term
- Add or delete linearized node
- Add or delete irrelevant nodes
- Add or delete duplicate nodes

NNs for autoregression. RJ Issues



NNs for autoregression. Sampler. Inefficiencies for large problems...

- Given current values of χ and σ^2 (β and λ are marginalized), for each γ_j , $j = 1, \dots, M$, generate a proposal $\tilde{\gamma}_j \sim N(\gamma_j, c\Sigma_\gamma)$, calculate the acceptance probability,

$$a = \min \left[1, \frac{p(\tilde{\gamma} | \mu_\gamma, \Sigma_\gamma) p(D' | y_1, \dots, y_q, \tilde{\gamma}, \sigma^2)}{p(\gamma | \mu_\gamma, \Sigma_\gamma) p(D' | y_1, \dots, y_q, \gamma, \sigma^2)} \right]$$

where $\gamma = (\gamma_1, \dots, \gamma_{j-1}, \gamma_j, \gamma_{j+1}, \dots, \gamma_M)$, and $\tilde{\gamma} = (\gamma_1, \dots, \gamma_{j-1}, \tilde{\gamma}_j, \gamma_{j+1}, \dots, \gamma_M)$. With probability a , replace γ by $\tilde{\gamma}$ and rearrange indices if necessary to satisfy order constraint. Otherwise, leave γ_j unchanged.

- Generate new values for parameters, drawing from their full conditional posteriors:

$\tilde{\beta} \sim p(\beta | D', \gamma, \lambda, \sigma^2, \chi)$ is a multivariate normal distribution.

$\tilde{\lambda} \sim p(\lambda | D', \gamma, \beta, \sigma^2, \chi)$ is a multivariate normal distribution.

$\tilde{\sigma}^2 \sim p(\sigma^2 | D', \gamma, \beta, \lambda)$ is an inverse gamma distribution.

- Given current values of $(\gamma, \beta, \lambda, \sigma^2)$, generate a new value for each hyperparameter by drawing from their complete conditional posterior distributions:

$\tilde{\mu}_\beta \sim p(\mu_\beta | D', \beta, \sigma_\beta^2)$ is a normal distribution.

$\tilde{\sigma}_\beta^2 \sim p(\sigma_\beta^2 | D', \beta, \mu_\beta)$ is an inverse gamma distribution.

$\tilde{\mu}_\lambda \sim p(\mu_\lambda | D', \lambda, \sigma_\lambda^2)$ is a multivariate normal distribution.

$\tilde{\sigma}_\lambda^2 \sim p(\sigma_\lambda^2 | D', \lambda, \mu_\lambda)$ is an inverse gamma distribution.

$\tilde{\mu}_\gamma \sim p(\mu_\gamma | D', \gamma, \Sigma_\gamma)$ is a multivariate normal distribution.

$\tilde{\Sigma}_\gamma \sim p(\Sigma_\gamma | D', \mu_\gamma)$ is an inverse Wishart distribution.

- Neural autoregressive density estimators
- Causal CNNs
- Transformers

- Encoding, Decoding, Encoding-Decoding.

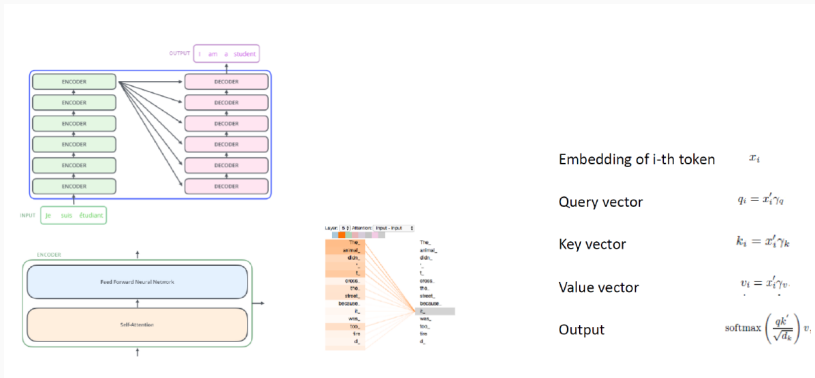
Decoder (Generator) works as follows:

For each transformer layer

- For each $y_{1:t}$ apply attention to compute attention weights $a_{1:t}$.
- Compute activation $z_t = \sum a_\tau y_\tau$
- Pass through feed-forward layer $h_t = MLP(z_t)$

Output $y_{t+1} = Cat(\text{softmax}(Wh_t))$

Transformers



Embedding of i-th token x_i

Query vector $q_i = x_i' \gamma_q$

Key vector $k_i = x_i' \gamma_k$

Value vector $v_i = x_i' \gamma_v$

Output $\text{softmax} \left(\frac{qk'}{\sqrt{d_k}} \right) v_i$

Full details with code in Session 2

<https://llms-cunef-icmat-rg2024.github.io/>

How this is used in GPTs in Session 3

<https://llms-cunef-icmat-rg2024.github.io/>

- Murphy (2023) PML, Chap.22
- Menchero, et al (2005) Bayesian Analysis of Nonlinear Autoregression Models Based on NNs, Neur. Comp.
- Prado and West (2010) Time series: modeling, computation and inference
- Vaswani Attention is all you need. Advances in NIPS 30 (2017).
- Zhao, Wayne Xin, et al. A survey of large language models. arXiv:2303.18223 1.2 (2023): 1-124.
- <https://llms-cunef-icmat-rg2024.github.io/>
- Papamarkou et al. Position: Bayesian deep learning is needed in the age of large-scale AI. arXiv:2402.00809 (2024).

Final comments

- Several of these models actually lack a PML approach
- Attacks on LLMs lagging
- Attacks on unsupervised learning lagging

- Normalizing flows
- Energy Models
- Diffusion models

Questions?