

# Aprendizaje profundo

## AUTOCODIFICADORES VARIACIONALES

---

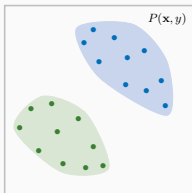
Gibran Fuentes-Pineda

Noviembre 2023

# Introducción

- Entrenamiento discriminativo busca modelar directamente  $P(y|\mathbf{x})$ , por lo que generar muestras de la misma distribución puede ser bastante difícil.
- Redes generativas aprenden  $P(\mathbf{x}, y)$  y pueden generar muestras de forma directa.

Modelo generativo



Modelo discriminativo

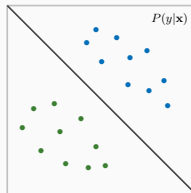


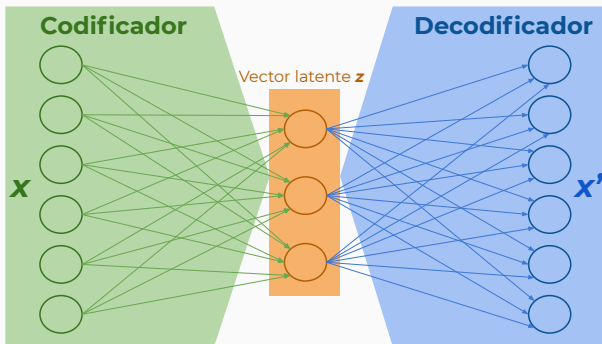
Imagen cortesía de Ricardo Montalvo Lezama

# Modelos generativos profundos

- Modelos autoregresivos
- Flujos normalizadores
- Modelos basados energía (por ej. Máquinas de Boltzmann)
- Autocodificadores variacionales
- Redes generativas antagónicas
- Modelos de difusión

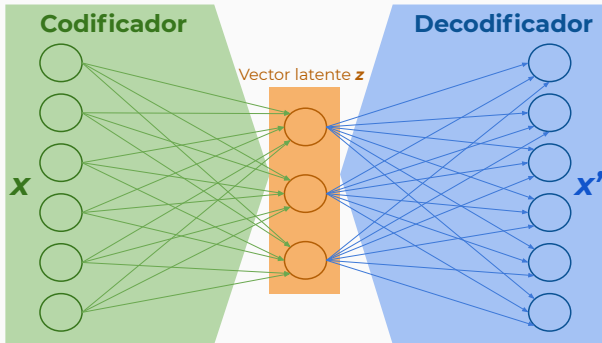
# Esquema codificador-decodificador

- Red entrenada generando sus mismas entradas
  - Codificador:  $\mathbf{z} = F_{\mathbf{W}_c, \mathbf{b}_c}(\mathbf{x})$
  - Decodificador:  $\mathbf{x}' = G_{\mathbf{W}_d, \mathbf{b}_d}(\mathbf{z})$
  - $\mathbf{W}_c, \mathbf{b}_c, \mathbf{W}_d, \mathbf{b}_d = \arg \min_{\mathbf{W}_c, \mathbf{b}_c, \mathbf{W}_d, \mathbf{b}_d} \|\mathbf{x} - \mathbf{x}'\|_2^2$



# Autocodificadores contractivos

- Hace  $z$  de menor dimensionalidad que  $x$



# Autocodificadores quita ruido

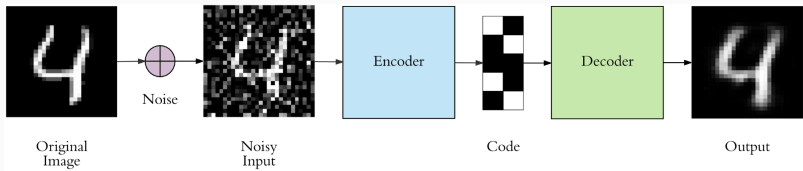


Imagen tomada de <https://towardsdatascience.com/applied-deep-learning-part-3-autoencoders-1c083af4d798>

# Autocodificadores dispersos

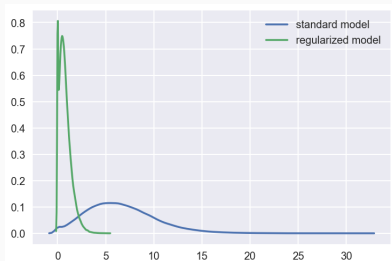


Figura tomada de <https://medium.com/towards-data-science/applied-deep-learning-part-3-autoencoders-1c083af4d798>

# Autocodificadores como estrategia de inicialización de pesos por capa (Restricted Boltzman Machines)

- Usado como estrategia de inicialización de pesos por capa (Restricted Boltzman Machines)

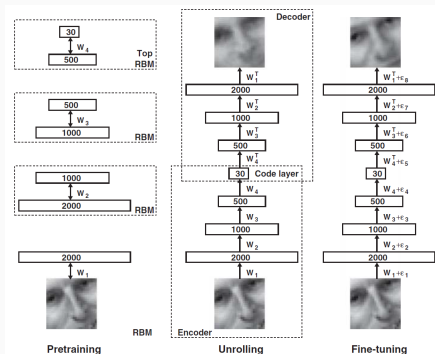


Imagen tomada de Hinton and Salakhutdinov. Reducing the Dimensionality of Data with Neural Networks, 2006.



# Modelos de variables latentes profundas

- Sea  $\mathbf{x}$  una observación muestreada aleatoriamente de una distribución no conocida, se presupone que

$$\mathbf{x} \sim P_{\theta}(\mathbf{x}) \approx P_{real}(\mathbf{x}).$$

- Los modelos de variables latentes representan la distribución  $P_{\theta}(\mathbf{x})$  usando variables latentes  $\mathbf{z}$

$$P_{\theta}(\mathbf{x}) = \int P_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z}.$$

- La distribución conjunta comúnmente se factoriza como

$$P_{\theta}(\mathbf{x}, \mathbf{z}) = P_{\theta}(\mathbf{x}|\mathbf{z}) \cdot P_{\theta}(\mathbf{z}).$$

- En los modelos profundos de variables latentes (DLVM), estas distribuciones están parametrizadas por redes neuronales.

# Proceso generativo de modelos de variables latentes

- Proceso generativo:
  1. Se obtiene una muestra  $\mathbf{z}$  de la distribución a priori  $\mathbf{z} \sim P_{\theta}(\mathbf{z})$ .
  2. Se obtiene una muestra  $\mathbf{x}$  de la distribución condicional  $\mathbf{x} \sim P_{\theta}(\mathbf{x}|\mathbf{z})$ .

- La probabilidad a posteriori está dada por

$$P_{\theta}(\mathbf{z}|\mathbf{x}) = \frac{P_{\theta}(\mathbf{x}, \mathbf{z})}{P_{\theta}(\mathbf{x})}.$$

- Tanto  $P_{\theta}(\mathbf{x})$  como  $P_{\theta}(\mathbf{z}|\mathbf{x})$  son intratables en DLVMs.

- Un autocodificador variacional es un DLVM en el que se presupone que

$$P_{\theta}(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$P_{\theta}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\sigma}^2)$$

- Para realizar aprendizaje de forma eficiente, se aproxima  $P_{\theta}(\mathbf{z}|\mathbf{x})$  con  $Q_{\phi}(\mathbf{z}|\mathbf{x})$ :

$$\boldsymbol{\mu}, \log \boldsymbol{\sigma} \leftarrow \text{RedNeuronalCodificadora}_{\phi}(\mathbf{x})$$

$$Q_{\phi}(\mathbf{z}|\mathbf{x}) \leftarrow \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}))$$

- A esto se le conoce como inferencia amortizada.

# Evidence Lower BOund (ELBO)

- En el entrenamiento se maximiza la cota inferior de evidencia:

$$\begin{aligned}\log P_{\theta}(\mathbf{x}) &= \mathbb{E}_{Q_{\phi}(\mathbf{z}|\mathbf{x})} [\log P_{\theta}(\mathbf{x})] \\&= \mathbb{E}_{Q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \log \left( \frac{P_{\theta}(\mathbf{x}, \mathbf{z})}{P_{\theta}(\mathbf{z}|\mathbf{x})} \right) \right] = \mathbb{E}_{Q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \log \left( \frac{P_{\theta}(\mathbf{x}, \mathbf{z})}{Q_{\phi}(\mathbf{z}|\mathbf{x})} \frac{Q_{\phi}(\mathbf{z}|\mathbf{x})}{P_{\theta}(\mathbf{z}|\mathbf{x})} \right) \right] \\&= \underbrace{\mathbb{E}_{Q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \log \left( \frac{P_{\theta}(\mathbf{x}, \mathbf{z})}{Q_{\phi}(\mathbf{z}|\mathbf{x})} \right) \right]}_{\mathcal{L}_{\theta, \phi}(\mathbf{x}) \text{ (ELBO)}} + \underbrace{\mathbb{E}_{Q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \log \left( \frac{Q_{\phi}(\mathbf{z}|\mathbf{x})}{P_{\theta}(\mathbf{z}|\mathbf{x})} \right) \right]}_{D_{\text{KL}}[Q_{\phi}(\mathbf{z}|\mathbf{x}) \| P_{\theta}(\mathbf{z}|\mathbf{x})] \geq 0}\end{aligned}$$

- Debido a que la divergencia KL es no negativa

$$\mathcal{L}_{\theta, \phi}(\mathbf{x}) \leq \log P_{\theta}(\mathbf{x})$$

# Entrenamiento de codificadores variacionales (1)

- Cota inferior de  $\log P_{\theta}(\mathbf{x})$  (aproxima máxima verosimilitud):  
 $\phi$  y  $\theta$  optimizadas conjuntamente
  - Equivalente a minimizar la divergencia de Kullback Leibler (KL) dada por

$$\begin{aligned} D_{KL} [Q_{\phi}(\mathbf{z}|\mathbf{x})||P_{\theta}(\mathbf{z}|\mathbf{x})] &= \mathbb{E}_{Q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \log \frac{Q_{\phi}(\mathbf{z}|\mathbf{x})}{P_{\theta}(\mathbf{z}|\mathbf{x})} \right] \\ &= \sum_{\mathbf{x}} Q_{\phi}(\mathbf{z}|\mathbf{x}) \log \frac{Q_{\phi}(\mathbf{z}|\mathbf{x})}{P_{\theta}(\mathbf{z}|\mathbf{x})} \\ &= \sum_{\mathbf{x}} Q_{\phi}(\mathbf{z}|\mathbf{x}) [\log Q_{\phi}(\mathbf{z}|\mathbf{x}) - \log P_{\theta}(\mathbf{z}|\mathbf{x})] \\ &= \mathbb{E}_{Q_{\phi}(\mathbf{z}|\mathbf{x})} [\log Q_{\phi}(\mathbf{z}|\mathbf{x}) - \log P_{\theta}(\mathbf{z}|\mathbf{x})] \end{aligned}$$

- Por regla de bayes

$$\begin{aligned} D_{KL} [Q_{\phi}(\mathbf{z}|\mathbf{x})||P_{\theta}(\mathbf{z}|\mathbf{x})] &= \mathbb{E}_{Q_{\phi}(\mathbf{z}|\mathbf{x})} [\log Q_{\phi}(\mathbf{z}|\mathbf{x}) - \log P_{\theta}(\mathbf{x}|\mathbf{z}) - \log P_{\phi}(\mathbf{z}) \\ &\quad + \log P_{\theta}(\mathbf{x})] \end{aligned}$$

- Despejando tenemos

$$\log P_{\theta}(\mathbf{x}) - D_{KL} [Q_{\phi}(\mathbf{z}|\mathbf{x})||P_{\theta}(\mathbf{z}|\mathbf{x})] = \mathbb{E}_{Q_{\phi}(\mathbf{z}|\mathbf{x})} [\log P_{\theta}(\mathbf{x}|\mathbf{z})] \\ - D_{KL} [Q_{\phi}(\mathbf{z}|\mathbf{x})||P_{\theta}(\mathbf{z})]$$

- Se busca maximizar  $\mathbb{E}_{Q_{\phi}(\mathbf{z}|\mathbf{x})} [\log P_{\theta}(\mathbf{x}|\mathbf{z})]$  y minimizar  $D_{KL} [Q_{\phi}(\mathbf{z}|\mathbf{x})||P_{\theta}(\mathbf{z})]$  mediante descenso por gradiente estocástico o variantes

- **Red codificadora** genera media  $\mu_z$  y covarianza diagonal  $\Sigma_z$  de  $Q_\phi(\mathbf{z}|\mathbf{x})$
- **Red dedificadora** genera media  $\mu_x$  y covarianza diagonal  $\Sigma_x$  de  $P_\theta(\mathbf{x}|\mathbf{z})$
- $\mathbf{x}'^{(i)}$  se obtiene
  1. Muestreando  $\mathbf{z}^{(i)}$  de  $Q_\phi(\mathbf{z}|\mathbf{x})$
  2. Muestreando  $\mathbf{x}'^{(i)}$  de  $P_\theta(\mathbf{x}|\mathbf{z})$

# Esquema general del modelo de autocodificador variacional

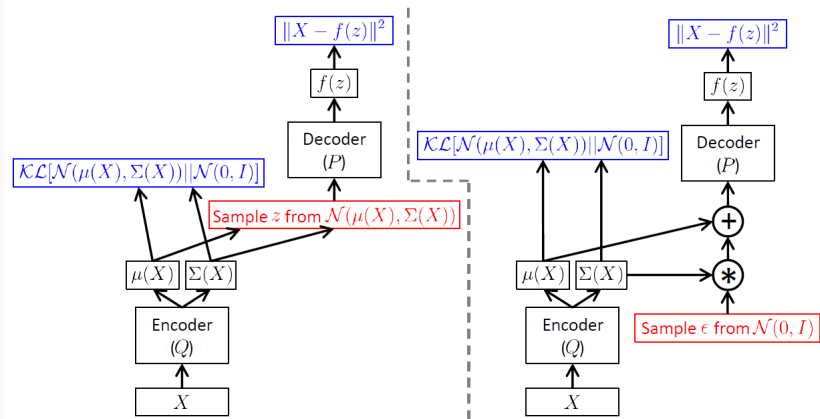


Imagen tomada de Carl Doersch. *Tutorial on Variational Autoencoders*, arXiv:1606.05908, 2016



# Generación de datos

- Para generar nuevos datos se muestrea  $z$  de  $\mathcal{N}(0, \mathbb{I})$

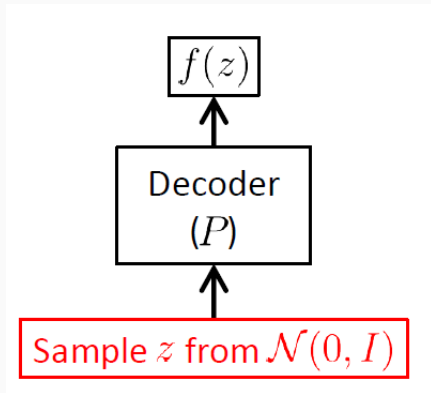


Imagen tomada de Carl Doersch. *Tutorial on Variational Autoencoders*, arXiv:1606.05908, 2016

# Aplicaciones: generación de dígitos

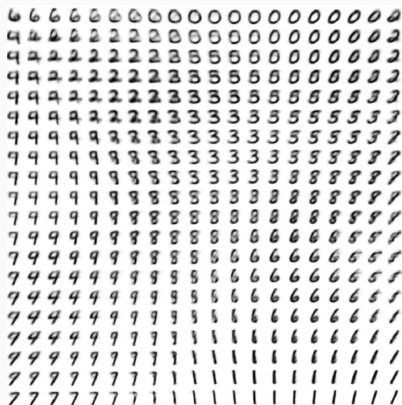


Imagen tomada de Kingma and Welling. Auto-Encoding Variational Bayes, 2013

# Aplicaciones: generación de imágenes



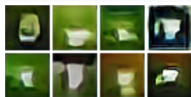
# Aplicaciones: generación de imágenes a partir de texto



A stop sign is flying in blue skies.



A herd of elephants flying in the blue skies.



A toilet seat sits open in the grass field.



A person skiing on sand clad vast desert.

Imagen tomada de Mansimov et al. Generating Images from Captions with Attention, 2015

# Aplicaciones: generación de dibujos (1)

- sketch-rnn

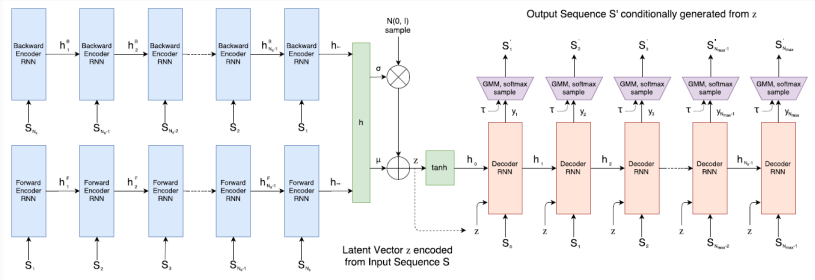


Imagen tomada de Ha and Eck. A Neural Representation of Sketch Drawings, 2017.

# Aplicaciones: generación de dibujos (2)

- sketch-rnn

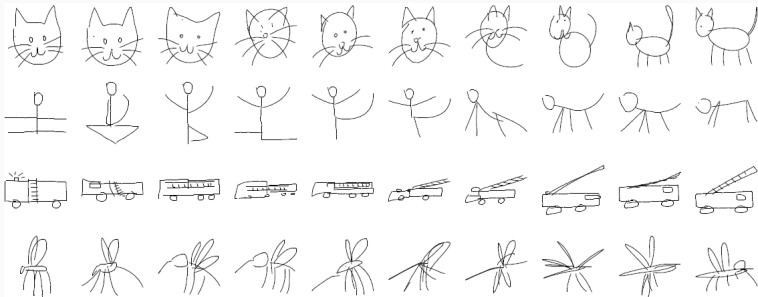


Imagen tomada de Ha and Eck. A Neural Representation of Sketch Drawings, 2017.

# Colapso del aposteriori

- Ocurre cuando la distribución variacional se aproxima a la distribución a priori no informativa para un conjunto de variables latentes
  - La generación es independiente del vector latente
- Estrategias de mitigación
  - Escalar pérdida KL por una tasa  $\beta \in (0, 1)$  y programarla respecto a épocas
  - Usar apriori y aposteriori que no sean gaussianas (por ej. von Mises-Fisher)
  - Agregar conexiones de salto
  - Mejorar el entrenamiento de la red codificadora
  - Reemplazar ELBO (por ej. con autocodificadores antagónicos)