

A collection of approximately 20 squares in three shades of blue and grey, scattered across the top half of the slide.

Generación de datos con DL

Data Mining

Ester Vidaña Vila

Generación de datos

Hasta ahora hemos visto cómo clasificar o reconstruir datos a partir de técnicas de *machine learning* o *deep learning*.

También podemos **generar** datos nuevos.

Los modelos generativos son una técnica no supervisada de *machine learning* que descubren y aprenden los patrones de los datos de entrada de forma que el modelo pueda usarlos para generar nuevos datos.

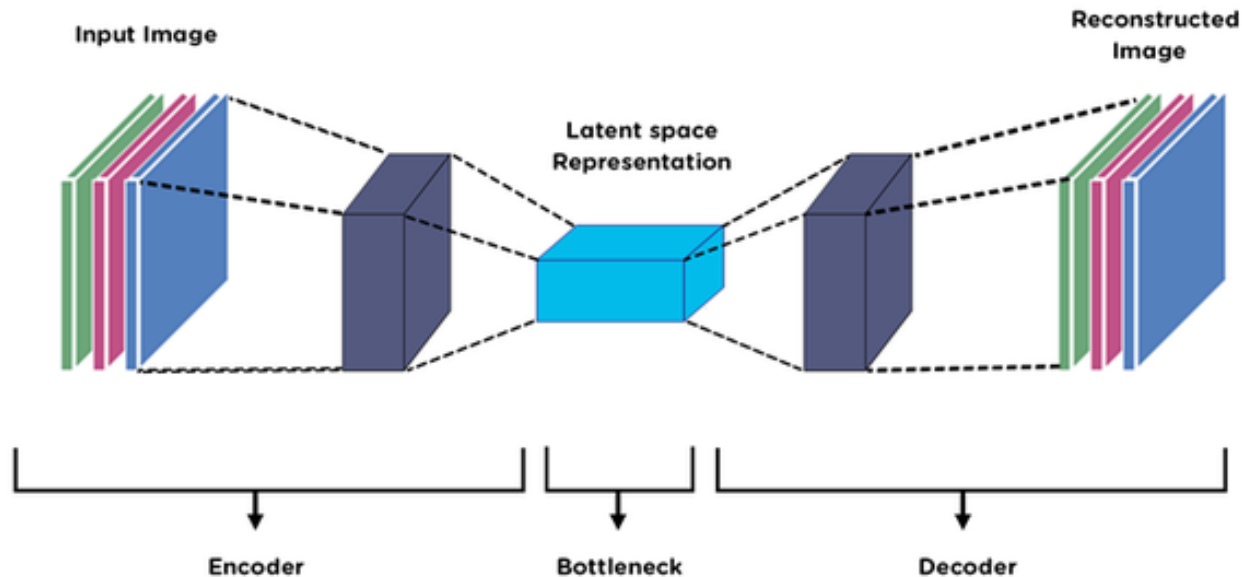
Vamos a ver algunas técnicas para generar nuevos datos:

1. Variational AutoEncoders (VAE).
2. Generative Adversarial Networks (GAN).
3. Long Short Term Memory (LSTM).

Variational AutoEncoder (VAE)

Variational AutoEncoder (VAE)

Recordatorio: un **autoencoder** es un modelo que tiene un modulo de *encoder* y otro de *decoder* que han sido entrenados para **minimizar el error de reconstrucción** entre los datos decodificados y los datos de entrada.

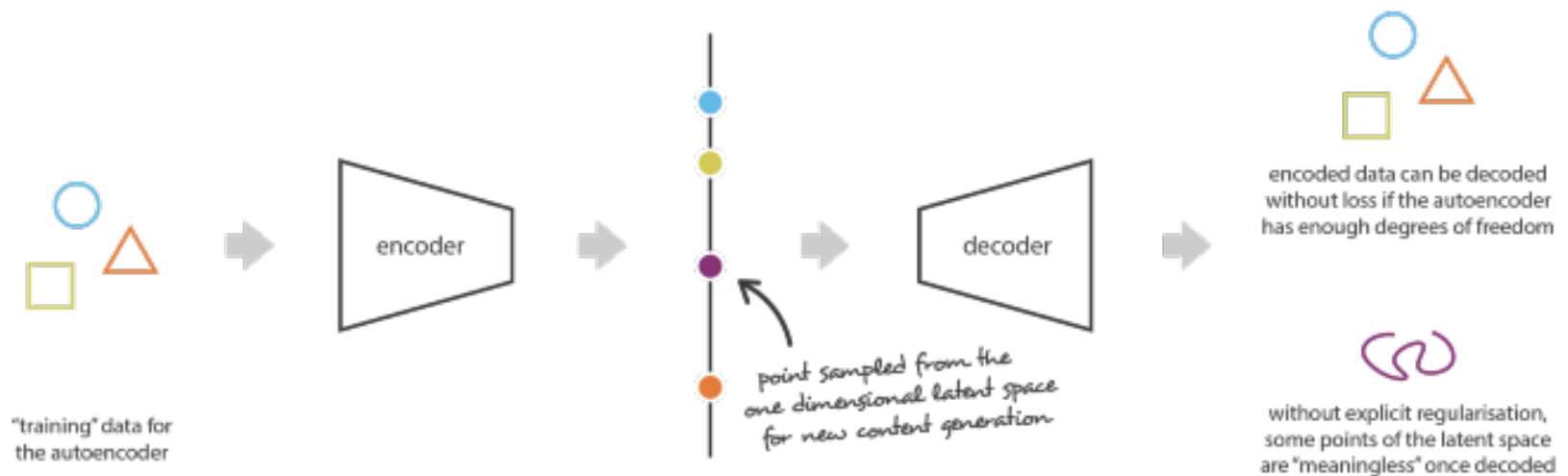


Fuente: <https://medium.com/@birla.deepak26/autoencoders-76bb49ae6a8f>

Variational AutoEncoder (VAE)

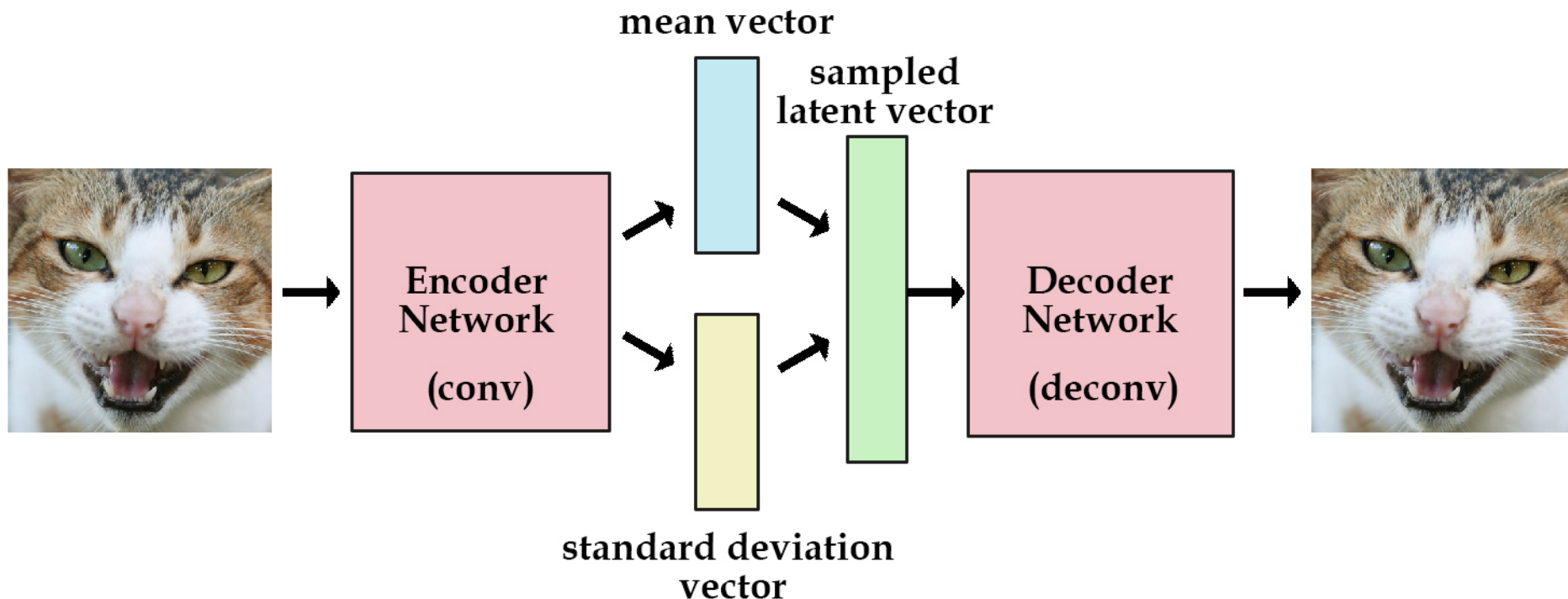
¿Podemos utilizar un autoencoder como generador?

Si lo usáramos directamente, probablemente nos generaría cosas extrañas.



Variational AutoEncoder (VAE)

En un VAE, en lugar de codificar una entrada como un punto del *latent space*, lo codificamos como una **distribución** sobre el *latent space*.



Variational AutoEncoder (VAE)

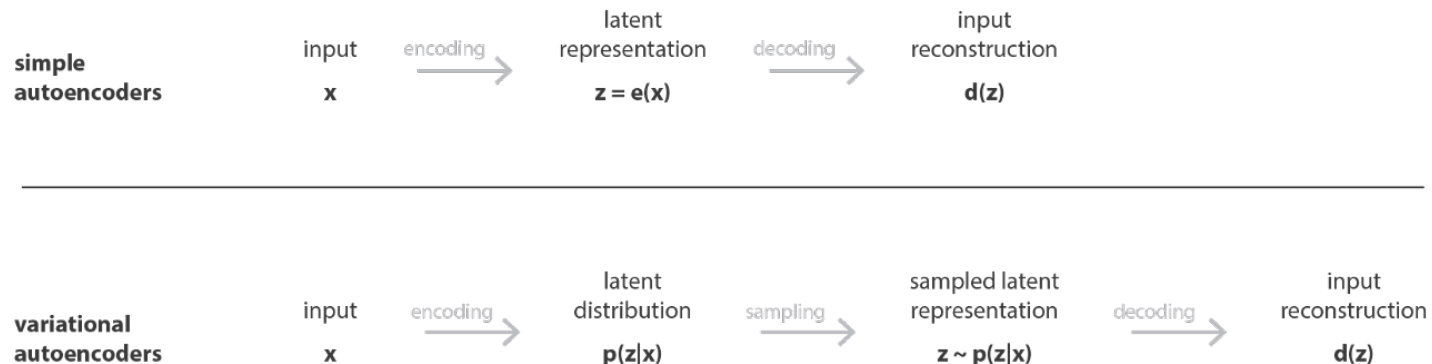
Así, el proceso de entrenamiento sería:

1. La entrada se codifica como una distribución sobre el *latent space*.
2. Se escoge un punto del *latent space* a partir de la distribución anterior.
3. Se decodifica la muestra del punto anterior para poder calcular el error de reconstrucción.
4. Se aplica *backpropagation* a la red.

A términos prácticos, se fuerza que la distribución sea normal.

Para conseguirlo, cuando se calcula la *loss*, se tiene en cuenta tanto el error de reconstrucción como la **KLDivergence**.

Básicamente, la KLDivergence evalúa cómo son de diferentes dos distribuciones de probabilidad. Si tenemos dos distribuciones completamente iguales, la KLDivergence será 0.



Fuente: <https://towardsdatascience.com/understanding-variational-autoencoders-vaes-f70510919f73>

Variational AutoEncoder (VAE)

Ejemplo visual:

<https://www.youtube.com/watch?v=sV2FOdGqIX0>



Fuente: <https://towardsdatascience.com/understanding-variational-autoencoders-vaes-f70510919f73>

Generative Adversarial Network (GAN)

Generative Adversial Network (GAN)

Dos redes neuronales se enfrentan entre ellas:

1. **Red generadora:** trata de generar imágenes que sean lo más realistas posible.

La red se alimenta de datos, por ejemplo, imágenes de personas al azar y, a partir de la información recibida, crea su propia imagen completamente nueva.

2. **Red discriminadora** (detective): trata de discernir qué imágenes son reales i qué imágenes son falsas.

La imagen no solo se declara falsa si se desvía demasiado de los datos básicos, sino también si es una imitación demasiado perfecta, ya que entonces los datos nuevos no se verían naturales.

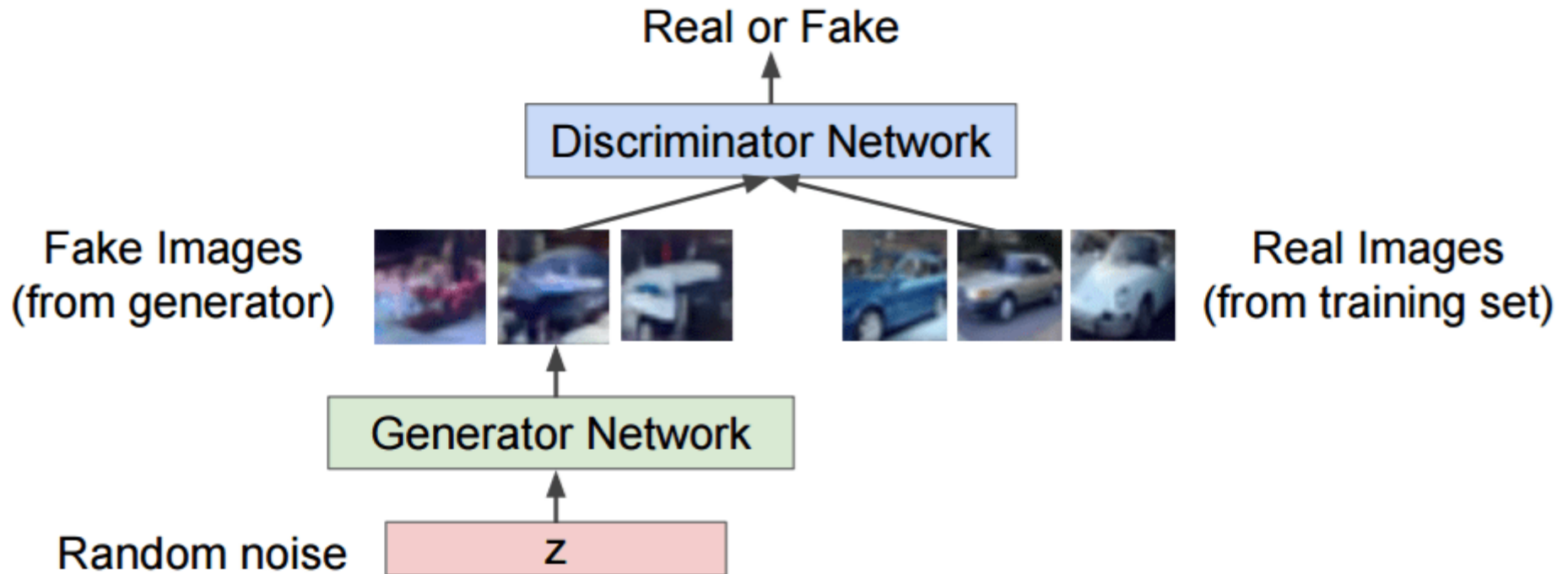
Generative Adversarial Network (GAN)

- **Primer entrenamiento de la red discriminadora:** la red discriminadora analiza los datos verdaderos, para aprender a comprenderlos.
- **Primer entrenamiento de la red generadora:** la red generadora comienza a generar datos falsos.
- **Segundo entrenamiento de la red discriminadora:** la red discriminadora se alimenta con los nuevos datos falsificados de la generadora y debe decidir qué datos considera verdaderos y qué datos considera falsos.
- **Segundo entrenamiento de la red generadora:** la red generadora mejora aún más con el resultado del segundo entrenamiento de la red discriminadora. Aquí la red generadora aprende sobre los puntos débiles de la discriminadora e intenta aprovecharlos, generando registros de datos falsos de apariencia aún más real.

Generative Adversarial Network (GAN)

- **Primer entrenamiento de la red discriminadora:** la red discriminadora analiza los datos verdaderos, para aprender a comprenderlos.
- **Primer entrenamiento de la red generadora:** la red generadora comienza a generar datos falsos.
- **Segundo entrenamiento de la red discriminadora:** la red discriminadora se alimenta con los nuevos datos falsificados de la generadora y debe decidir qué datos considera verdaderos y qué datos considera falsos.
- **Segundo entrenamiento de la red generadora:** la red generadora mejora aún más con el resultado del segundo entrenamiento de la red discriminadora. Aquí la red generadora aprende sobre los puntos débiles de la discriminadora e intenta aprovecharlos, generando registros de datos falsos de apariencia aún más real.

Generative Adversarial Network (GAN)



Fuente: <https://puentesdigitales.com/2019/04/05/todo-lo-que-necesitas-saber-sobre-las-gan-redes-generativas-antagonicas/>

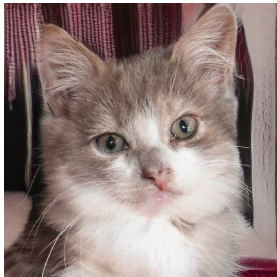
Ejemplo: generación de rostros de personas

<https://www.youtube.com/watch?v=kSLJriaOumA&t=283s>

Ejemplo 2: gatos (o personas) que no existen

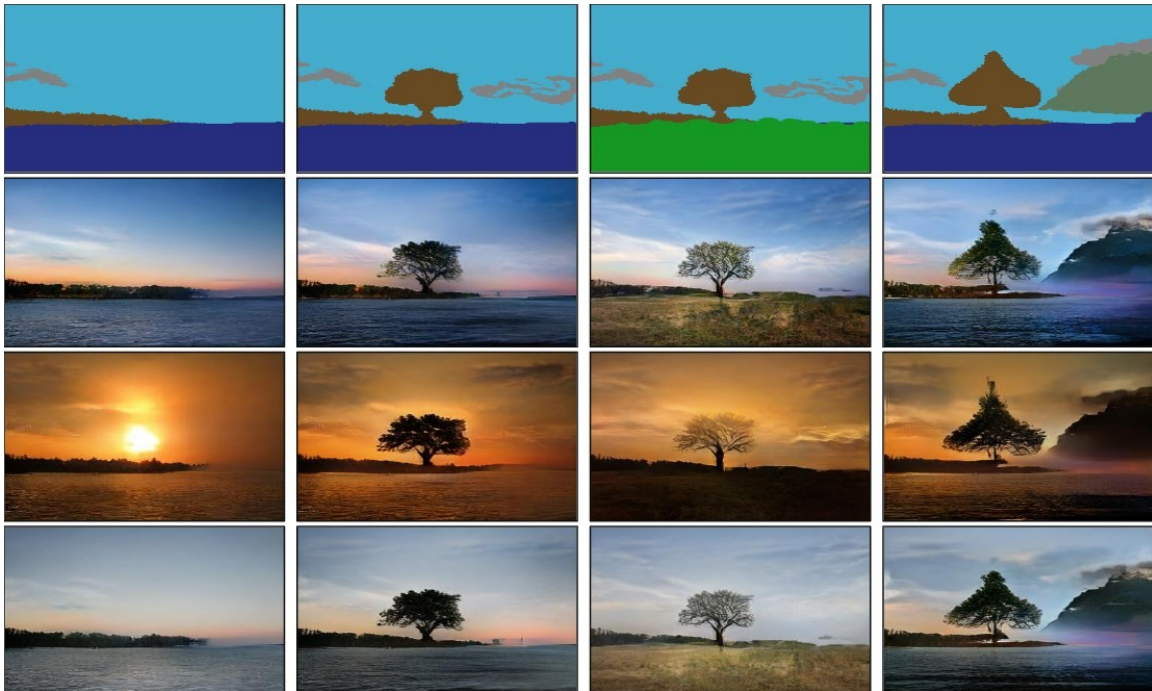
<https://thiscatdoesnotexist.com/>

<https://thispersondoesnotexist.com/>



Ejemplo 3: de dibujos a paisajes

<https://www.youtube.com/watch?v=p5U4NgVGAwg>



<http://gaugan.org/gaugan2/>

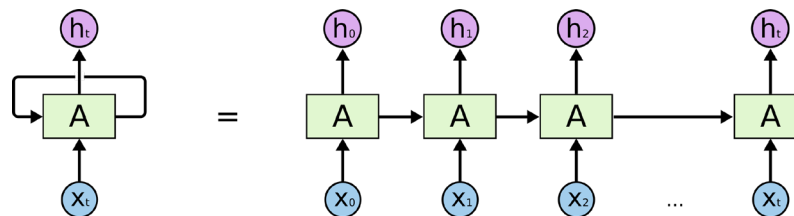
Ejemplo interactivo

<https://poloclub.github.io/ganlab/>

Long Short Term Memory (LSTM)

LSTM

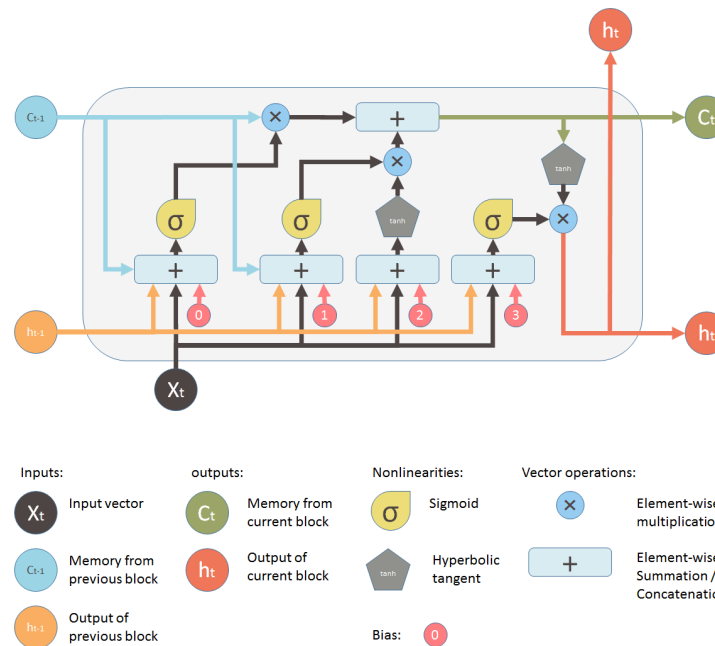
- Long Short Term Memory:
- Los humanos no empezamos a aprender desde cero cada día, ¿verdad?
- Entendemos las palabras basándonos en el entendimiento de las palabras que hemos entendido antes.
- Tenemos **persistencia**. Así pues, hay redes neuronales que también tienen persistencia. Estas se llaman **Recurrent Neural Networks**:
 - Son redes que tienen bucles.
 - Podemos pensar en ellas como múltiples copias de una misma red, que va pasando un mensaje nodo a nodo como si fuera su sucesor.



Fuente: <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

LSTM

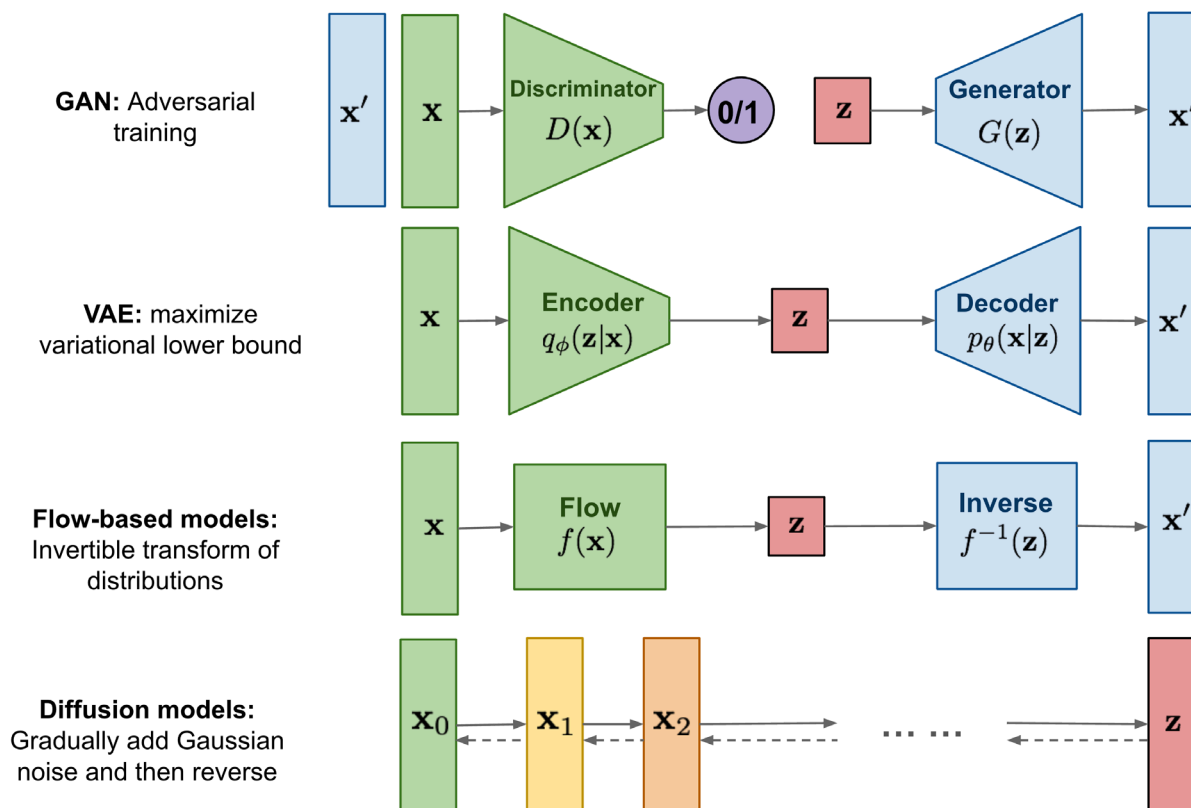
- Las RNN tienen problemas de memorización a largo plazo, para ello nacieron las LSTM, que solventan este inconveniente.



Fuente: <https://www.analyticsvidhya.com/blog/2022/02/explaining-text-generation-with-lstm/>

Diffussion models

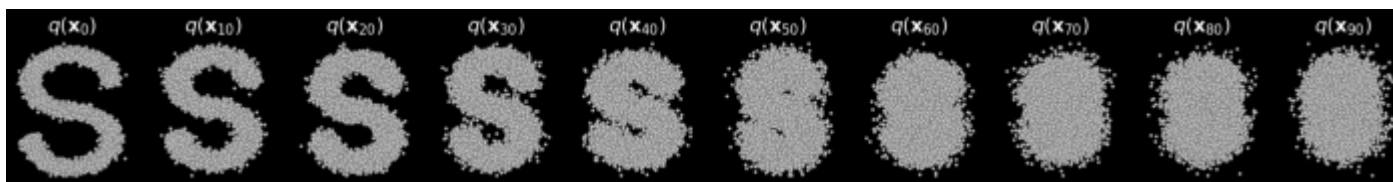
Diffusion models



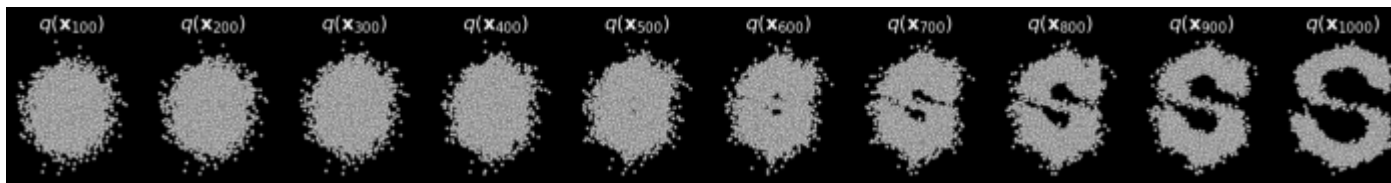
Fuente: <https://lilianweng.github.io/posts/2021-07-11-diffusion-models/>

Diffusion models

- Forward pass: vamos añadiendo ruido a los datos.



- Reconstruction: miramos de reconstruir los datos.



- Nuestro modelo debe aprender a reconstruir los datos a partir de ruido.

Fuente: <https://towardsdatascience.com/diffusion-models-made-easy-8414298ce4da>