



UNIVERSITY  
OF AMSTERDAM

# Machine Learning for Physics and Astronomy

Juan Rojo

VU Amsterdam & Theory group, Nikhef

***Natuur- en Sterrenkunde BSc (Joint Degree), Honours Track***  
***Lecture 6, 05/10/2020***

# Today's lecture

- Adversarial Learning
- Generative Adversarial Networks
- Guest lecture by **Dr. Sascha Caron** on applications of ML to high-energy physics

# **Generative Models & Adversarial Learning**

# Generative Models

Most ML models discussed here (Supervised NNs, logistic regression, ensemble models) are **discriminative**: designed to identify **differences between groups of data**

*e.g. cats vs dogs discrimination*

these models cannot carry some tasks such as **drawing new examples** from an unknown probability distribution: for this we need **generative models**

*e.g. learn how to draw new examples  
of cat and dog images*

*e.g. generate new samples of a given phase  
of the Ising model*

generative models are Machine Learning techniques that allows to learn **how to generate new examples** similar to those found in a training dataset

Here we consider **energy-based generative models**: close connection with statistical physics

*hence the term Boltzmann Learning*

# Maximum Entropy Generative Models

basic concept is the **Shannon information-theoretic entropy**, which quantifies the statistical uncertainty one has about a random variable drawn from a probability distribution

$$S_p = \text{Tr}_x p(x) \log p(x)$$

*↑  
sum/integral over all  
possible values of variable*

*variables whose prob dist  
I would like to estimate*

assume we have a **set of models**, functions of  $x$ , whose average should coincide with some observed values. What should be their **underlying prob dist?**

$$\begin{array}{ccc} \{f_i(x)\} & \langle f_i \rangle_{\text{obs}} & i = 1, \dots, n_{\text{dat}} \\ \text{model for sample } i \text{ as function of inputs } x & \text{models} & \text{observations} \end{array}$$

**Principle of Maximum Entropy:** choose the probability distribution with the largest uncertainty (Shannon entropy) subject to the observational constraints

$$\langle f_i \rangle_{\text{model}} = \int d\mathbf{x} f_i(\mathbf{x}) p(\mathbf{x}) = \langle f_i \rangle_{\text{obs}}$$

*the selected distribution is the one that makes admits the most ignorance beyond the stated data*

# Maximum Entropy Generative Models

this condition can be expressed as a **Lagrange Multiplier problem** by minimising:

$$\mathcal{L}[p] = -S_p + \sum_i \lambda_i \left( \langle f_i \rangle_{\text{obs}} - \int d\mathbf{x} f_i(\mathbf{x}) p(\mathbf{x}) \right) + \gamma \left( 1 - \int d\mathbf{x} p(\mathbf{x}) \right)$$

*Shannon entropy*                    *observational constraints*                    *normalisation*

whose solution gives us the **Maximum Entropy distribution**

$$p(\mathbf{x}) = \frac{\exp \left( \sum_i \lambda_i f_i(\mathbf{x}) \right)}{\int d\mathbf{x} \exp \left( \sum_i \lambda_i f_i(\mathbf{x}) \right)} = \frac{\exp \left( \sum_i \lambda_i f_i(\mathbf{x}) \right)}{Z}$$

*partition function*                     $E/k_B T = - \sum_i \lambda_i f_i(\mathbf{x})$

which is nothing but the **Boltzmann distribution in statistical mechanics**, and where the parameters of the distribution are fixed by the observations

$$\partial_{\lambda_i} \log Z = \langle f_i \rangle_{\text{data}}$$

# Energy-based Generative Models

these MaxEnt models can be used to **infer the underlying probability distributions** from a finite set of observations, which subsequently can be used to **generate new instances**

training an energy-based generative model: using the data to infer the model parameters

$$E(\mathbf{x}; \boldsymbol{\theta}) = - \sum_i \theta_i f_i(\mathbf{x}) \quad \xrightarrow{\text{model parameters}}$$

as in Supervised Learning, we need to specify a **cost function**, which however in the case of generative models is much subtler: what defines a good model?

the most useful method is to **maximise the log-likelihood** of the training set

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \mathcal{L}(\{\boldsymbol{\theta}\})$$

$$\mathcal{L}(\{\boldsymbol{\theta}\}) = \langle \log p_{\boldsymbol{\theta}}(\mathbf{x}) \rangle_{\text{data}} = - \langle E(\mathbf{x}; \boldsymbol{\theta}) \rangle_{\text{data}} - \log Z(\{\boldsymbol{\theta}\})$$

where we have used that the **generative probability distribution** is of the Boltzmann form and that the partition function does not depend on the data

# Boltzmann machines

The training of **energy-based generative models** proceeds usually via SGD

MaxEnt generative models are defined by the choice of the **energy**

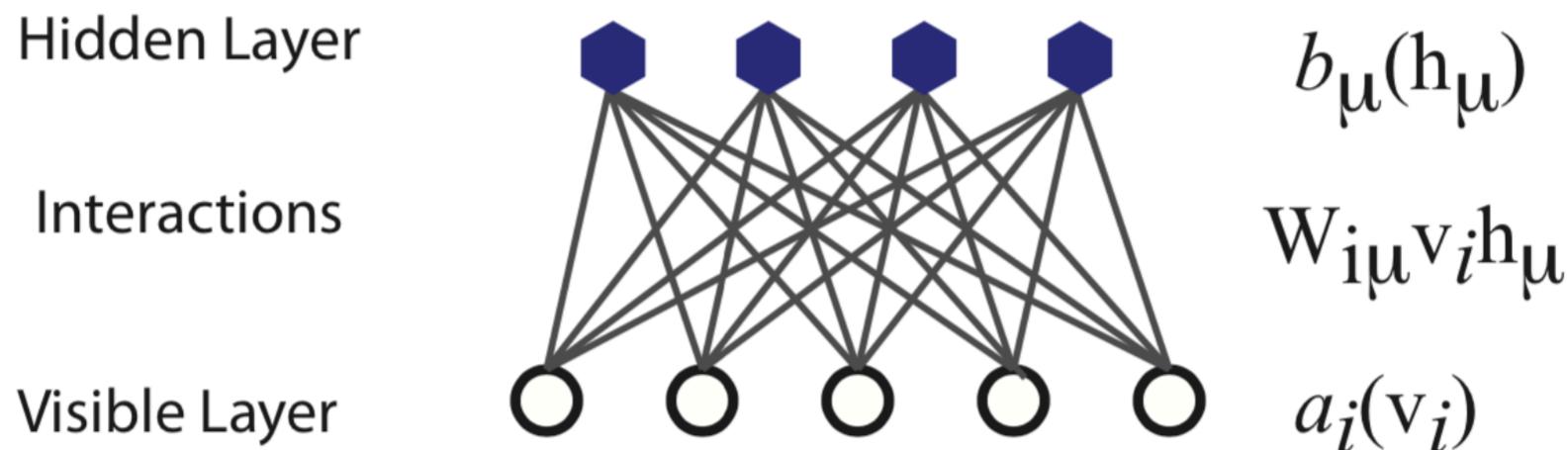
$$p_{\theta}(\mathbf{x}) = \frac{\exp \left( \sum_i \theta_i f_i(\mathbf{x}) \right)}{\int d\mathbf{x} \exp \left( \sum_i \theta_i f_i(\mathbf{x}) \right)} = \frac{\exp \left( \sum_i \theta_i f_i(\mathbf{x}) \right)}{Z} = \frac{\exp (-E(\mathbf{x}; \boldsymbol{\theta}))}{Z}$$

*model parameters* 

one can construct various **other generative models** with different choices of the energy

e.g. *Restricted Boltzmann machines*

$$E(\mathbf{x}; \boldsymbol{\theta}) = - \sum_i a_i(v_i) - \sum_{\mu} b_{\mu}(h_{\mu}) - \sum_{i\mu} W_{i\mu} v_i h_{\mu}$$



# The Kullback-Leibler divergence

The KL divergence, a measure of the similarity between two probability distributions  $p(\mathbf{x})$  and  $q(\mathbf{x})$  plays an important role in machine learning applications

$$D_{KL}(p \parallel q) = \int d\mathbf{x} p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} \quad D_{KL}(q \parallel p) = \int d\mathbf{x} q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})}$$

which can be symmetrised to construct a **squared metric** (distance)

$$D_{JS}(p \parallel q) = \frac{1}{2} \left( D_{KL}\left(p \middle\| \frac{p+q}{2}\right) + D_{KL}\left(q \middle\| \frac{p+q}{2}\right) \right)$$

*Jensen-Shannon divergence*

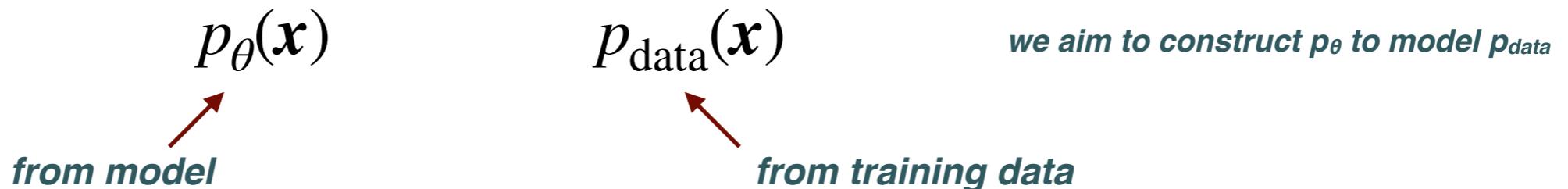
The KL-divergence is **positive-definite**, and only vanishes when  $p(\mathbf{x})=q(\mathbf{x})$

$$D_{KL}(p \parallel q) \geq 0$$

in general the integral cannot be computed and one needs to sample the two probability distributions by means of a suitable binning

# The Kullback-Leibler divergence

In **generative models** one deals with two probability distributions (data and model), which we would like to have as similar as possible



however subtleties about how we define **similarity** have large implications for the model training

maximising the **log-likelihood of the data under the model** is the same as **minimising the KL divergence** between the data distribution and the model distribution

$$\begin{aligned} D_{KL}(p_{\text{data}} || p_\theta) &= \int d\mathbf{x} p_{\text{data}}(\mathbf{x}) \log \frac{p_{\text{data}}(\mathbf{x})}{p_\theta(\mathbf{x})} \\ &= \int d\mathbf{x} p_{\text{data}}(\mathbf{x}) \log p_{\text{data}}(\mathbf{x}) - \int d\mathbf{x} p_{\text{data}}(\mathbf{x}) \log p_\theta(\mathbf{x}) \\ &= S_p[p_{\text{data}}] - \langle \log p_\theta \rangle_{\text{data}} \end{aligned}$$

# The Kullback-Leibler divergence

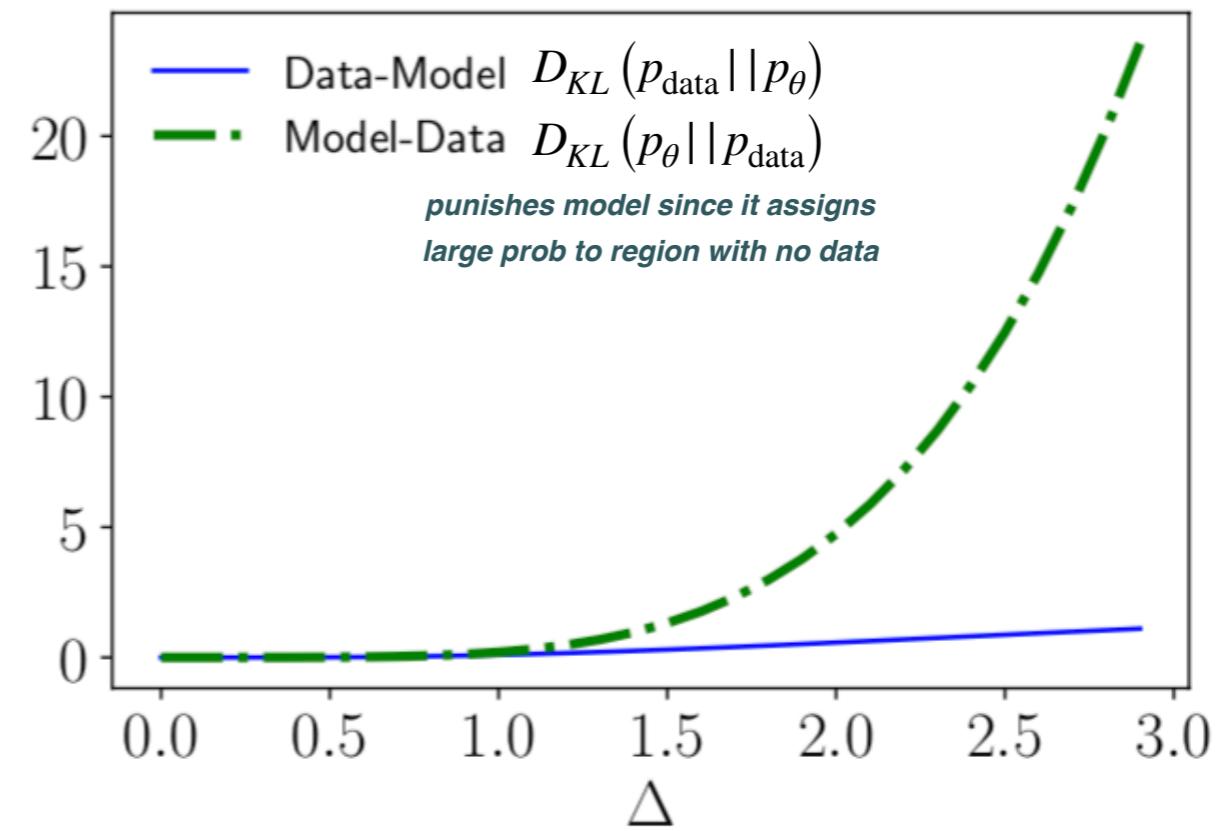
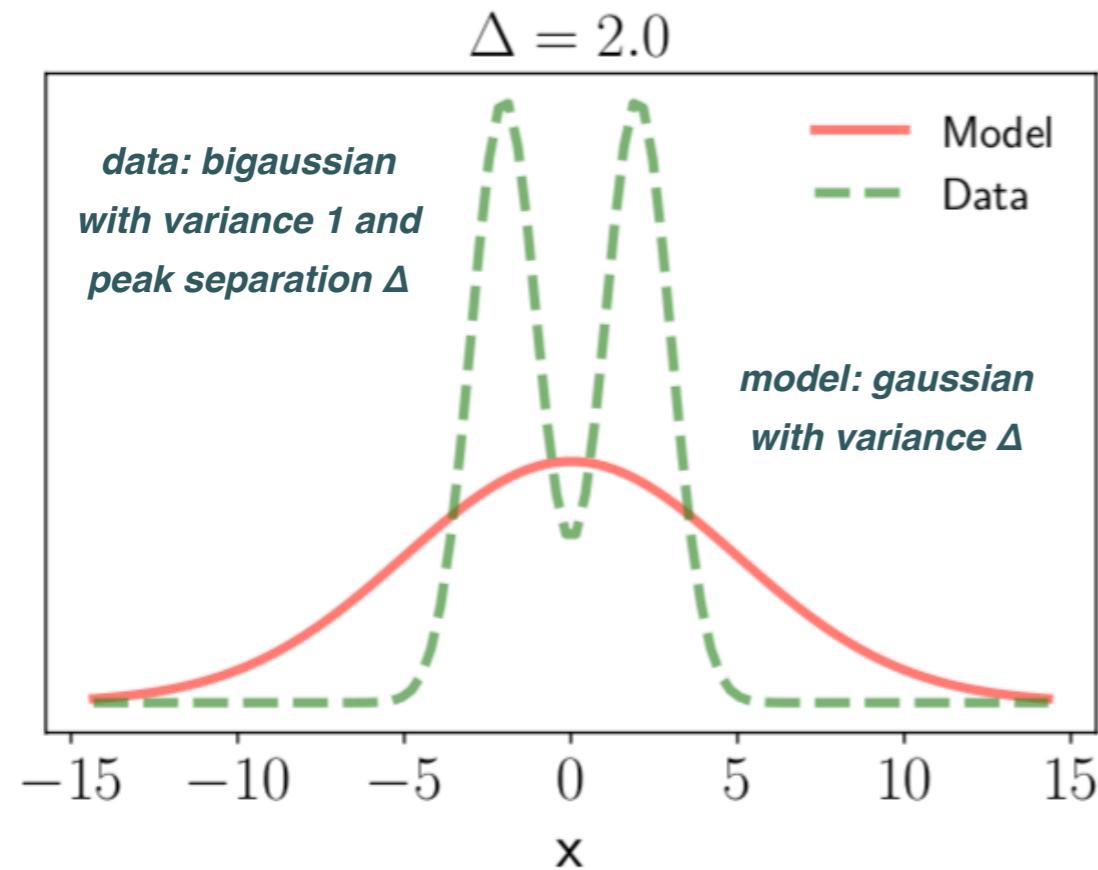
maximising the **log-likelihood** of the data under the model is the same as **minimising the KL divergence** between the data distribution and the model distribution

$$\langle \log p_\theta(x) \rangle_{\text{data}} = S_p[p_{\text{data}}] - D_{KL}(p_{\text{data}} \parallel p_\theta)$$

*↑*                           *↑*                           *↑*  
*Log-likelihood of data under model*      *entropy of data:  
independent of model parameters*      *KL-divergence*

# Similarity

**Similarity** between probability distributions is a subtle concept



$$D_{KL}(p_{\text{data}} \parallel p_{\theta}) \longrightarrow$$

misses important information when comparing the data and theory probability distributions

# Adversarial Learning

$$(1) \ D_{KL} (p_{\text{data}} || p_{\theta})$$

*Calculable using sampling*

*(draw values of  $\theta$  accordingly to  $p_{\theta}$  and compare with data)*

$$(2) \ D_{KL} (p_{\theta} || p_{\text{data}})$$

*Large when model over-weights low-density regions near real peaks*

*but not calculable since  $p_{\text{data}}$  unknown ....*

In **Adversarial Learning** we achieve a similar goal as that of minimising **(2)** by training a **discriminator** to distinguish between real data points and samples from the model

By punishing the model for generating points that can be easily discriminated from the data, Adversarial Learning decreases the **weight of regions in the model space that are far away from data points**, regions that inevitably arise when maximising the likelihood

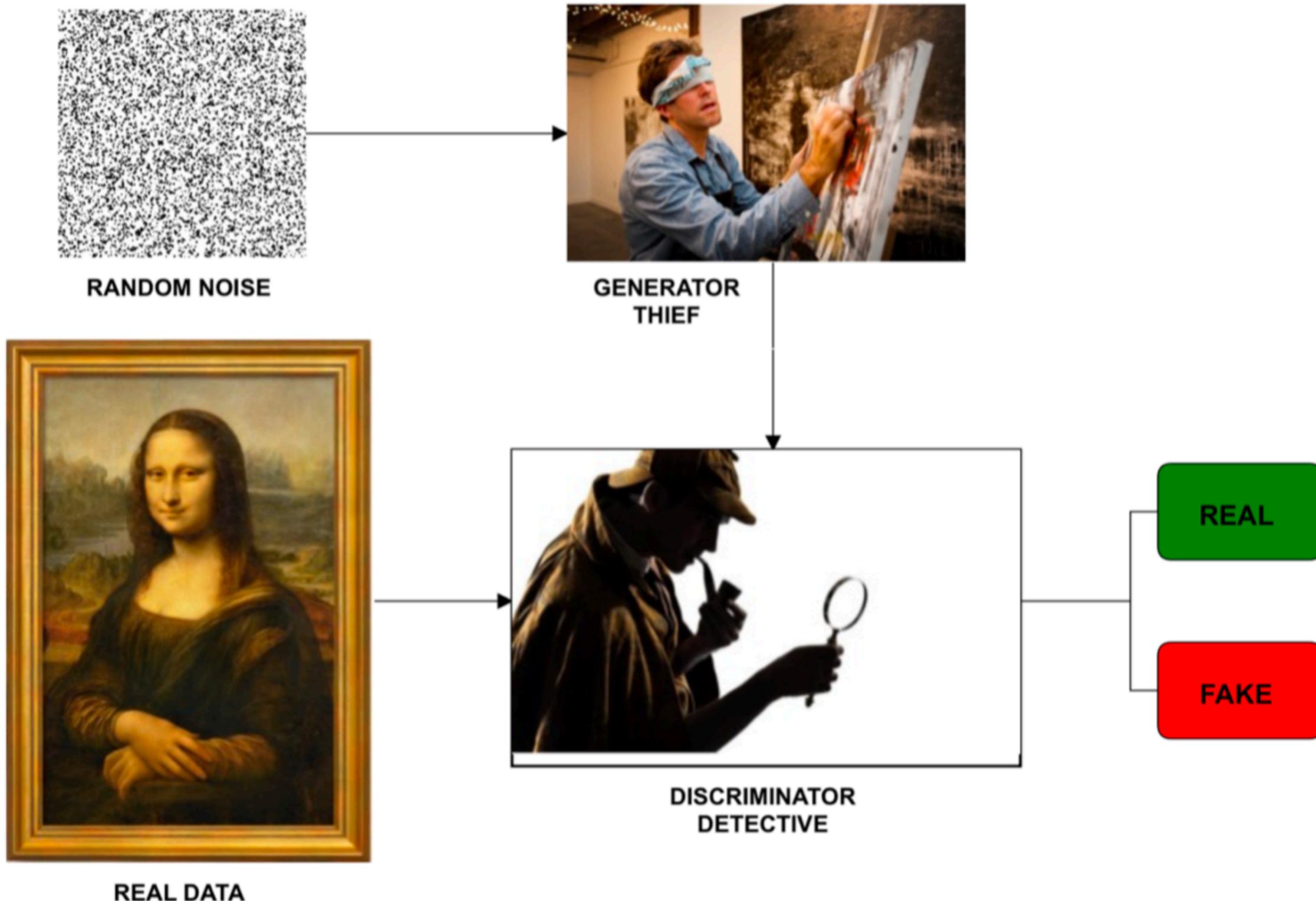
# Generative adversarial networks

Generative Adversarial Networks (GANs) are deep neural network architectures, composed by two independent NNs which **compete against each other**

- ✿ (1) A **generator G** NN that creates (samples) pseudo-data by inferring the probability distribution associated to the training dataset
- ✿ (2) A **discriminator D** NN which determines the probability of a given sample arises from the actual training data rather than having been produced by **G**

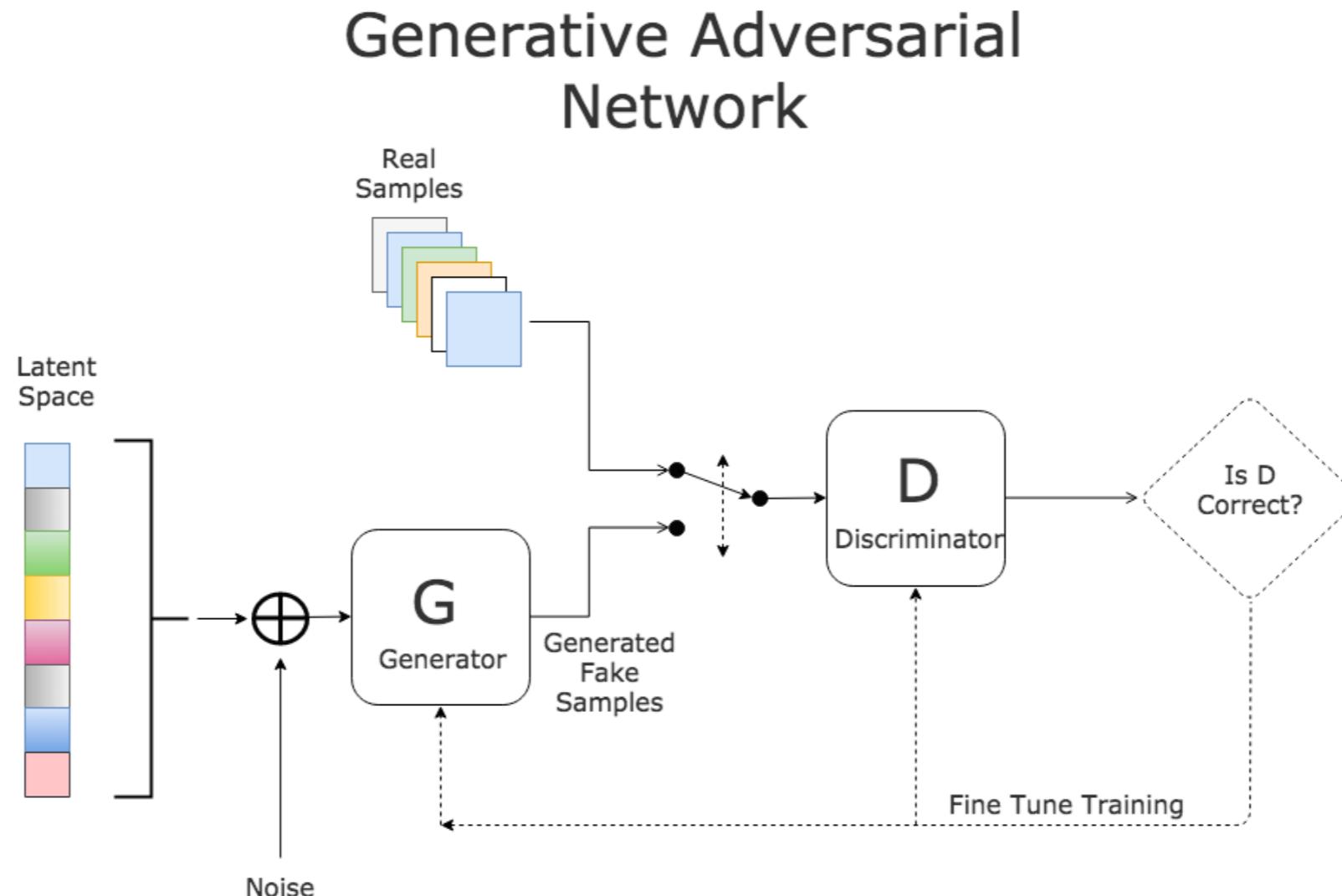
the generator network **G** should be trained to maximise the probability that the discriminator network **D** makes a mistake: that is, **G** should generate pseudo-data samples that are virtually **indistinguishable** from the actual data

# Generative adversarial networks



# Generative Adversarial Networks

- New architecture for an **unsupervised neural network training** (unlabelled samples)
- Based on two **independent nets** that work separately and act as **adversaries**:
  - the **Discriminator (D)** undergoes training and plays the role of classifier
  - the **Generator (G)** and is tasked to generate random samples that **resemble real samples** with a twist rendering them as fake samples.



# Generative Adversarial Networks

# GAN training

As with other NN architectures one uses GD to train GANs, but now one has to **update sequentially** the model parameters of both **G** and **D**

- Take a sample of  $N$  data points from the training set

$$\{\mathbf{x}_n\}_{n=1}^N \quad \mathbf{x}_n = (x_{n,1}, x_{n,2}, \dots, x_{n,p}) \quad p = \text{number of features per sample}$$

- Produce a sample of  $N$  pseudo-data points from generator **G** (at  $\text{ite}_0$  this is random noise)

$$\{\mathbf{z}_n\}_{n=1}^N \quad \mathbf{z}_n = (z_{n,1}, z_{n,2}, \dots, z_{n,p})$$

- Evaluate the cost function: since we are dealing with binary classification (true/false) the appropriate cost function is the **cross-entropy**

$$C(\theta_D, \theta_G) = \frac{1}{N} \sum_{n=1}^N (\log D(\mathbf{x}_i) + \log(1 - D(G(\mathbf{z}_i))))$$

NN params of G      output of D when input  
NN params of D      a real data sample      a "fake" data sample produced by G

- Train **D** using GD to maximise its discrimination capability

# GAN training

- Evaluate the cost function: since we are dealing with binary classification (true/false) the appropriate cost function is the **cross-entropy**

$$C(\theta_D, \theta_G) = \frac{1}{N} \sum_{n=1}^N (\log D(x_i) + \log(1 - D(G(z_i))))$$

- Train **D** using GD to maximise its discrimination capability

$$\mathbf{v}_t = \eta_t \nabla_{\theta_D} C(\theta_{D,t}, \theta_{G,t}), \quad \theta_{D,t+1} = \theta_{D,t} - \mathbf{v}_t$$

- At this point **D** can tell apart data from pseudo-data pretty well, so we need to train **G** to generate better (closer to the training set) pseudo-data samples
- Produce a sample of  $N$  pseudo-data points from the generator **G**  $\{z_n\}_{n=1}^N$

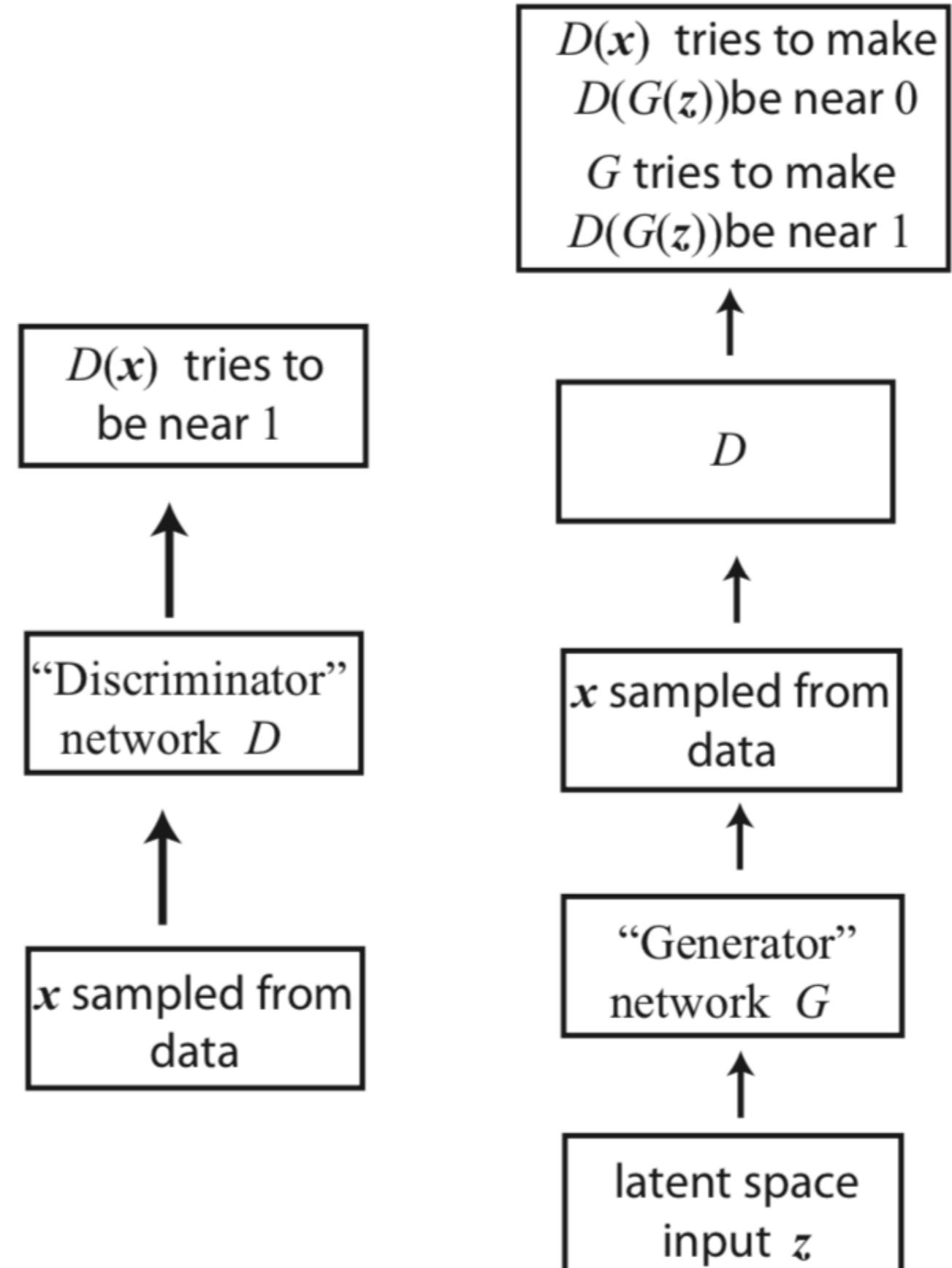
$$C(\theta_D, \theta_G) = \frac{1}{N} \sum_{n=1}^N \log(1 - D(G(z_i)))$$

*output of D (now with its parameters fixed)*

$$\mathbf{v}_t = \eta_t \nabla_{\theta_G} C(\theta_{D,t}, \theta_G), \quad \theta_{G,t+1} = \theta_{G,t} - \mathbf{v}_t$$

# GAN training

the generator and discriminator are sequentially trained and iterated until convergence is achieved, at this point **D** cannot tell apart the pseudo-data from **G** from the real data



# Image generation with GANs



*<https://thispersondoesnotexist.com/>*

# Image generation with GANs

**<https://thispersondoesnotexist.com/>**

# Generative Adversarial Networks