



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 7, 19/10/2021***

# Today's lecture

- ✿ Information Theory for Machine Learning
- ✿ Generative Models
- ✿ Adversarial Learning and Generative Adversarial Networks
- ✿ Kernel Methods and Support Vector Machines

# **Information Theory for Machine Learning**

# Information Theory

the main ideas of **information theory** are ubiquitous in machine learning. In order to discuss **generative models** and **adversarial learning** it is useful to review a few basic concepts

consider a discrete random variable  $x$ . How much information is brought in if we measure an specific value of this variable (its **degree of surprise**)?

*the greater the surprise, the bigger the amount of info  
e.g. ``the sun did not rise this morning'' contains more surprise than ``the sun rose this morning''!*

clearly a measure of ``*how much information a new measurement brings*'' depends on the probability distribution of  $x$ ,  $p(x)$ . Thus we look for an information-measuring function  $h[p(x)]$

*Q: have you encountered before a function that we can use to measure information*

# Information Theory

the main ideas of **information theory** are ubiquitous in machine learning. In order to discuss generative models and adversarial learning it is useful to review a few basic concepts here

consider a discrete random variable  $x$ . How much information is brought in if we measure an specific value of this variable (its **degree of surprise**)?

*the greater the surprise, the bigger the amount of info  
e.g. ``the sun did not rise this morning'' contains more surprise than ``the sun rose this morning''!*

clearly a measure of ``*how much information a new measurement brings*'' depends on the probability distribution of  $x$ ,  $p(x)$ . Thus we look for an information-measuring function  $h[p(x)]$

*if two events  $x$  and  $y$  are independent, the information we gain from observing both of them is the sum of the information we gain from each of them separately (additive)*

$$h(x, y) = h(x) + h(y)$$

*But these two events are statistically independent,  
thus their probabilities obey*

$$p(x, y) = p(x)p(y)$$

# Information Theory

the main ideas of **information theory** are ubiquitous in machine learning. In order to discuss generative models and adversarial learning it is useful to review a few basic concepts

consider a discrete random variable  $x$ . How much information is brought in if we measure an specific value of this variable (its **degree of surprise**)?

*the greater the surprise, the bigger the amount of info  
e.g. ``the sun did not rise this morning'' contains more surprise than ``the sun rose this morning''!*

clearly a measure of ``*how much information a new measurement brings*'' depends on the probability distribution of  $x$ ,  $p(x)$ . Thus we look for an information-measuring function  $h[p(x)]$

*if two events  $x$  and  $y$  are independent, the information we gain from observing both of them is the sum of the information we gain from each of them separately (additive)*

$$h(x, y) = h(x) + h(y)$$

*But these two events are statistically independent, thus their probabilities obey*

$$p(x, y) = p(x)p(y)$$

$$h(x) = -\ln p(x)$$

*information is positive!*

$$0 \leq p(x) \leq 1$$

# Information Theory

what is the **total amount of information** contained in a set of instances of  $x$ ?

*It will be the expectation value of  $h(x)$  wrt the probability distribution  $p(x)$*

$$H[x] = \sum_x p(x)h(x) = - \sum_x p(x)\ln p(x)$$

which is known as the **entropy** associated to the random variable  $x$

*Q1: which kind of events dominate the calculation of the total entropy?*

*Q2: is entropy larger in a ``smooth'' or in a ``spiky'' probability distribution?*

# Information Theory

what is the **total amount of information** contained in a set of instances of  $x$ ?

*It will be the expectation value of  $h(x)$  wrt the probability distribution  $p(x)$*

$$H[x] = \sum_x p(x)h(x) = - \sum_x p(x)\ln p(x)$$

which is known as the **entropy** associated to the random variable  $x$

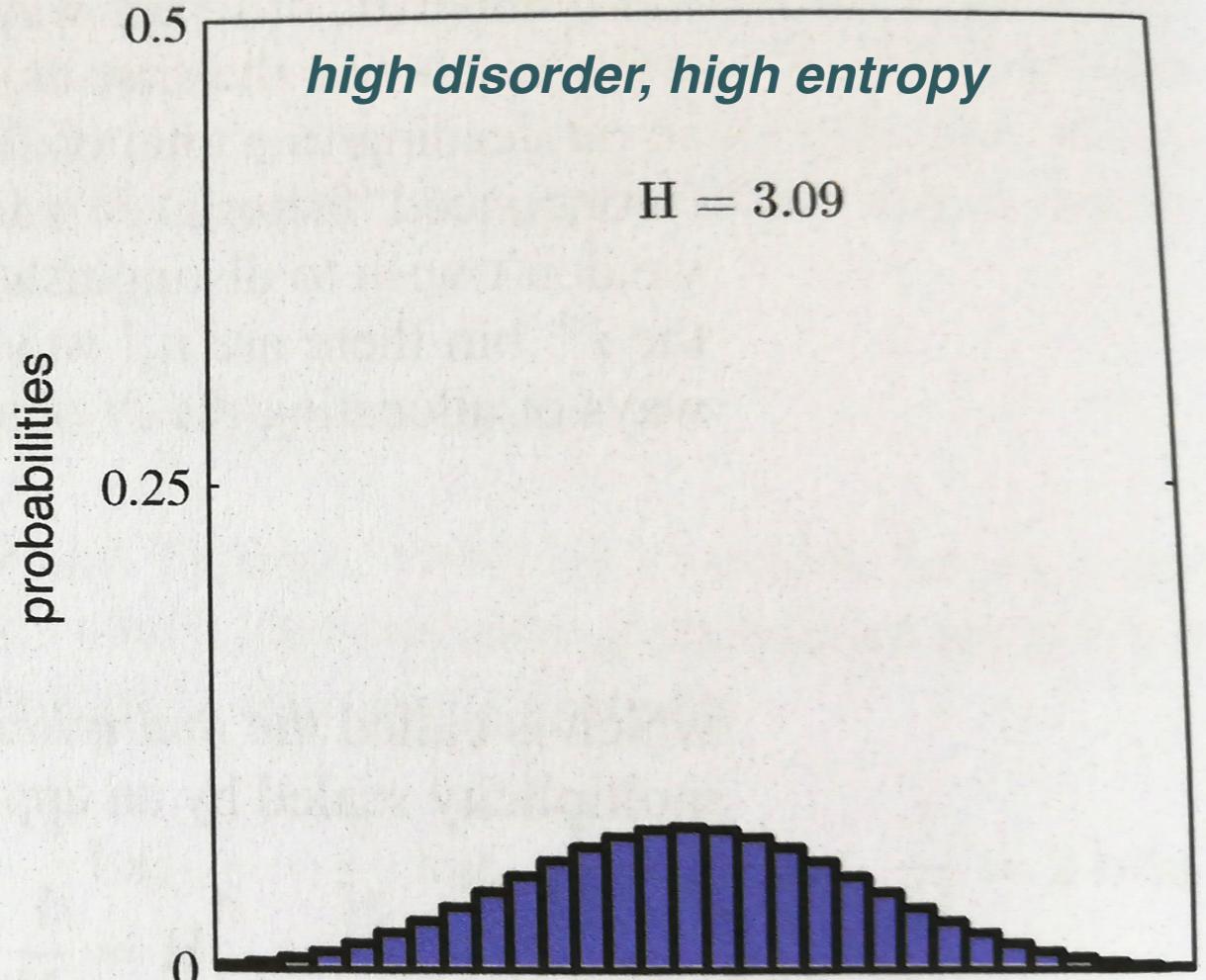
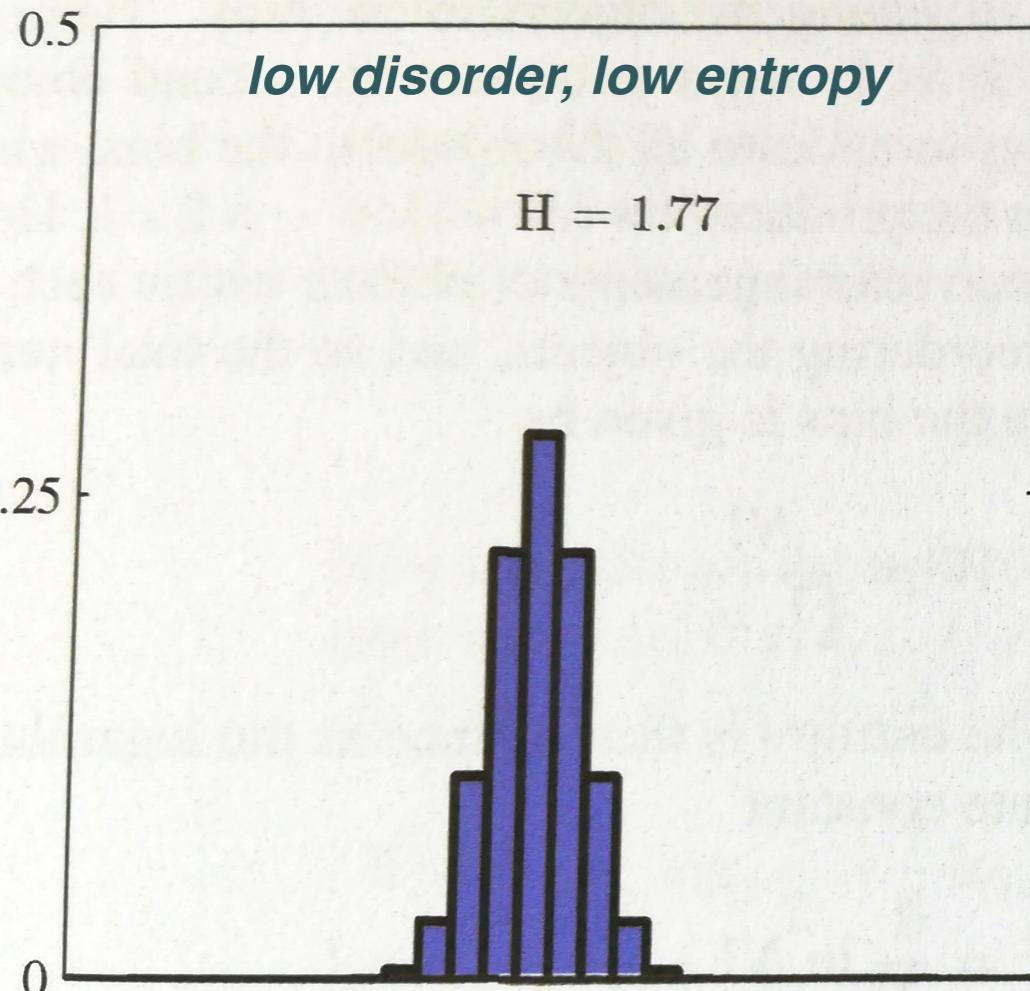
Note some **important properties** about the **surprise function**  $h(x)$  and the **entropy**  $H[x]$ :

- Low-probability events have associated a **high information content**
- However low-probability events carry a **small weight** in the total entropy
- One can show that steeply-varying probability distributions lead to smaller entropies than smooth distributions (largest entropy is for uniform distributions)
- Entropy can be understood to as the **average amount of information** needed to specify the state of a random variable

*entropy is of course closely connected  
to disorder: the higher the disorder, the largest the associated entropy*

# Information Theory

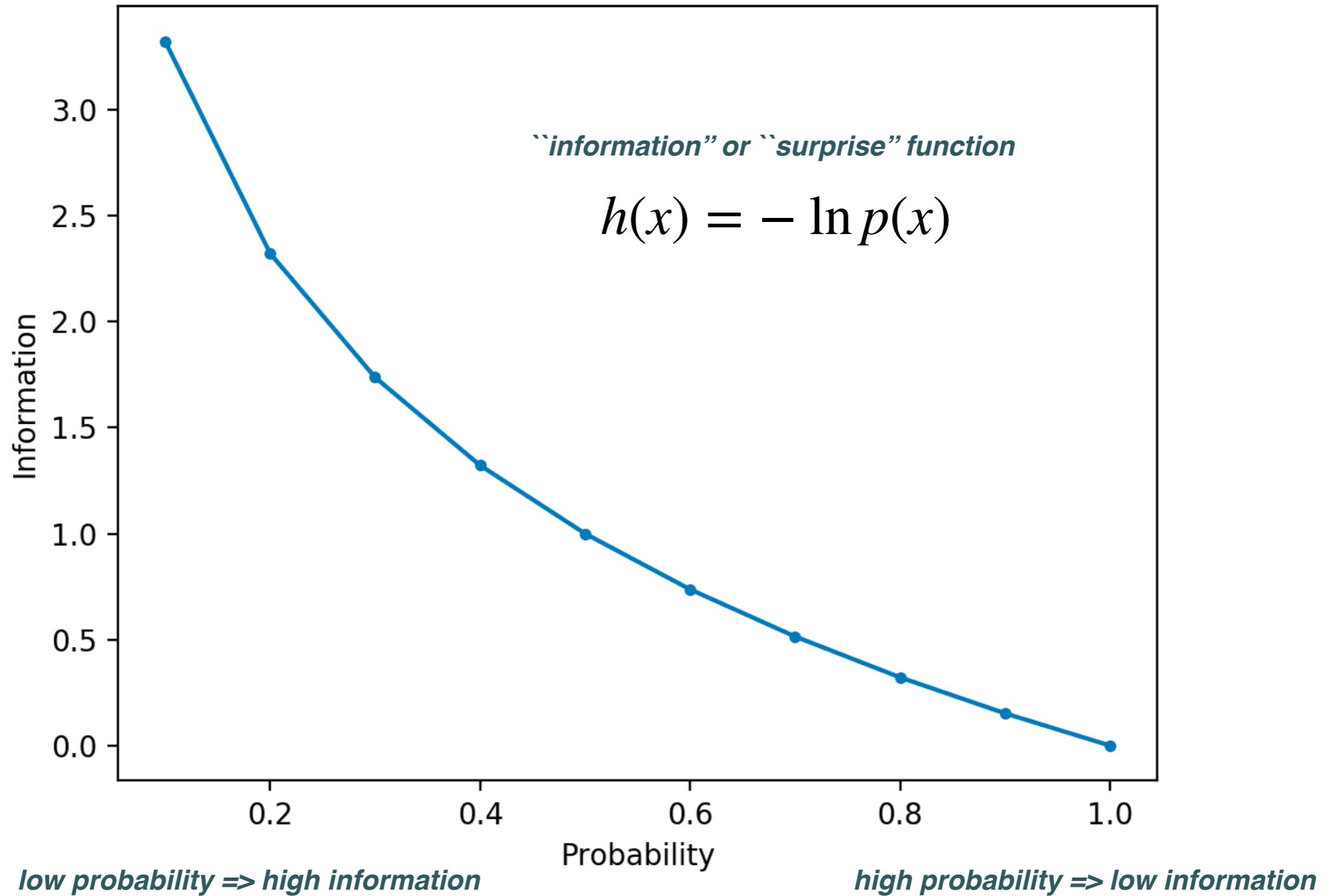
*the broader the distribution, the higher the disorder and thus the higher the entropy*



*the highest possible value of the entropy corresponds to a uniform distribution*

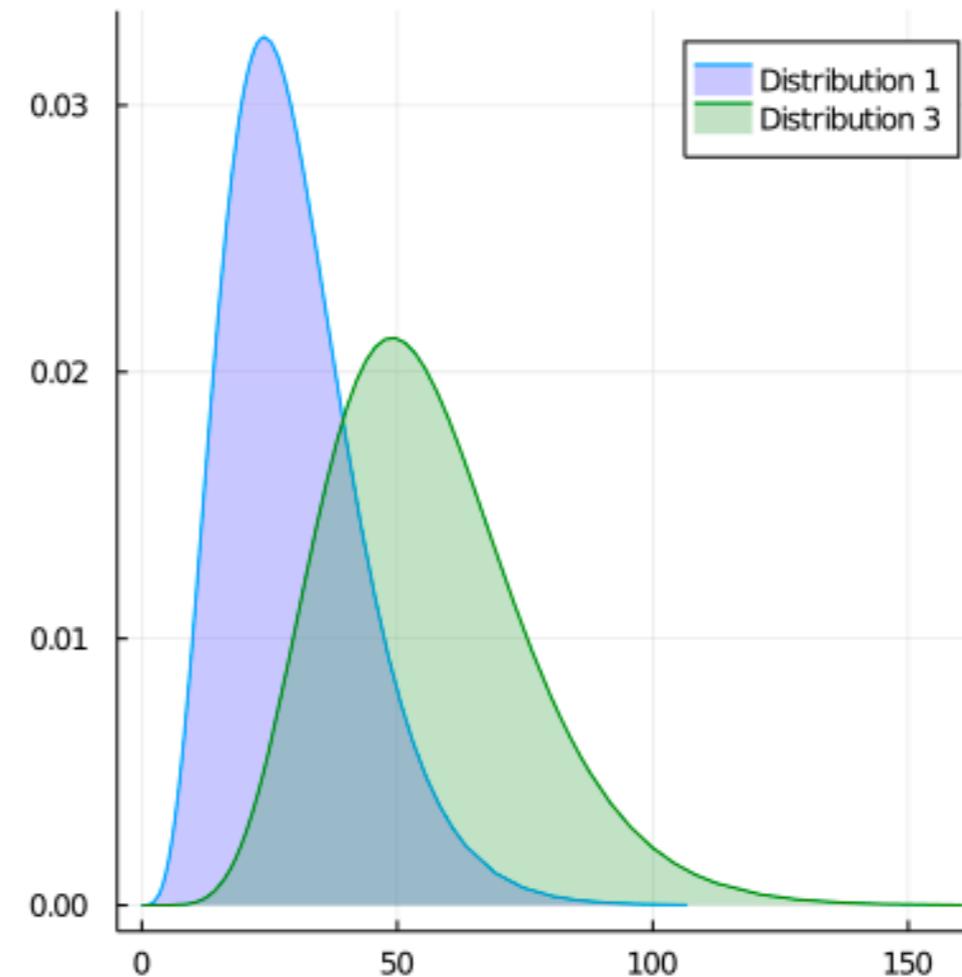
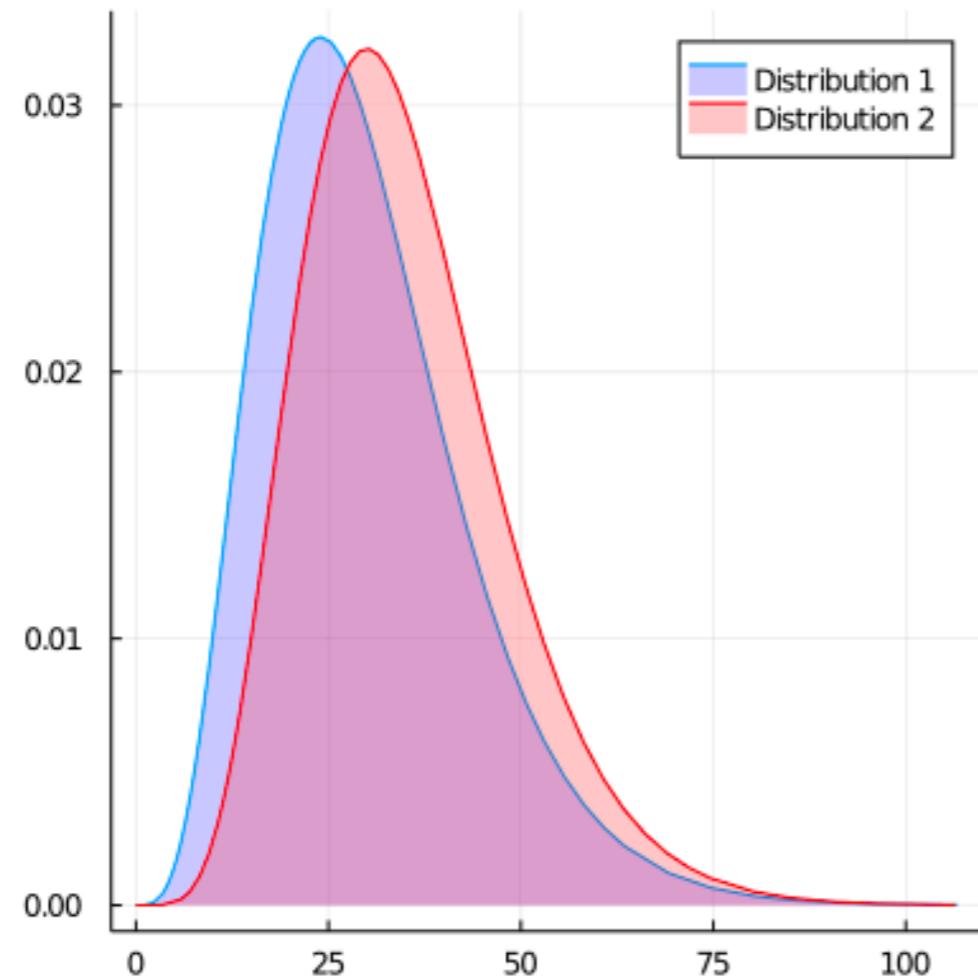
# Information Theory

Probability vs Information



# Similarity of probability distributions

In many ML applications we need to assess the **degree of similarity** between two probability distributions



*Q1: which quantities can we use to determine if Distribution #2 is closer to Distribution #1 or to Distribution #3?*

*Q2: is similarity an absolute or a relative concept?*

# Similarity of probability distributions

In many ML applications we need to assess the **degree of similarity** between two probability distributions

We can use the concepts of Information Theory and compute the relative entropy between two distributions

$$H[x] = - \sum_x p(x) \ln p(x) \simeq \int dx p(x) \ln p(x)$$

Entropy of stochastic variable  $x$ , distributed following  $p(x)$

$$\int dx p(x) \ln q(x)$$

Relative entropy of  $q(x)$  with respect to  $p(x)$

now compute difference between actual entropy and relative entropy

$$\int dx p(x) \ln p(x) - \int dx p(x) \ln q(x) = \int dx p(x) \ln \frac{p(x)}{q(x)}$$

the smaller this quantity, the more similar the two probability distributions  $p(x)$  and  $q(x)$  are

# Similarity of probability distributions

This quantity is the **Kullback-Leibler divergence**

$$D_{\text{KL}}(p \parallel q) = \int dx p(x) \ln \frac{p(x)}{q(x)}$$

*Let's develop some intuition: what do you think we should find when comparing two Gaussian distributions?*

# Similarity of probability distributions

This quantity is the **Kullback-Leibler divergence**

$$D_{\text{KL}}(p \parallel q) = \int dx p(x) \ln \frac{p(x)}{q(x)}$$

$$p(x) = \frac{1}{\sqrt{2\pi\sigma_p^2}} e^{-(x-x_p)^2/2\sigma_p^2} \quad q(x) = \frac{1}{\sqrt{2\pi\sigma_q^2}} e^{-(x-x_q)^2/2\sigma_q^2}$$

$$D_{\text{KL}}(p \parallel q) = \ln \frac{\sigma_q}{\sigma_p} + \frac{1}{(2\pi\sigma_p^2)^{1/2}} \int_{-\infty}^{\infty} dx e^{-(x-x_p)^2/2\sigma_p^2} \times \left( -\frac{(x-x_p)^2}{2\sigma_p^2} + \frac{(x-x_q)^2}{2\sigma_q^2} \right)$$

$$D_{\text{KL}}(p \parallel q) = \ln \frac{\sigma_q}{\sigma_p} + \frac{1}{2\sigma_q^2} \left( (x_q - x_p)^2 + (\sigma_p^2 - \sigma_q^2) \right)$$

*The closer the means and the standard deviations, the smaller the KL divergence for two gaussian distributions*

# Similarity of probability distributions

In many ML applications we need to assess the **degree of similarity** between two probability distributions

The **Kullback-Leibler divergence** is a measure of the **similarity** between two probability distributions  $p(\mathbf{x})$  and  $q(\mathbf{x})$  and plays an important role in machine learning models

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

*note that KL divergence is not symmetric (as usual with statistical distances)*

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

*Jensen-Shannon divergence (symmetric version of KL divergence)*

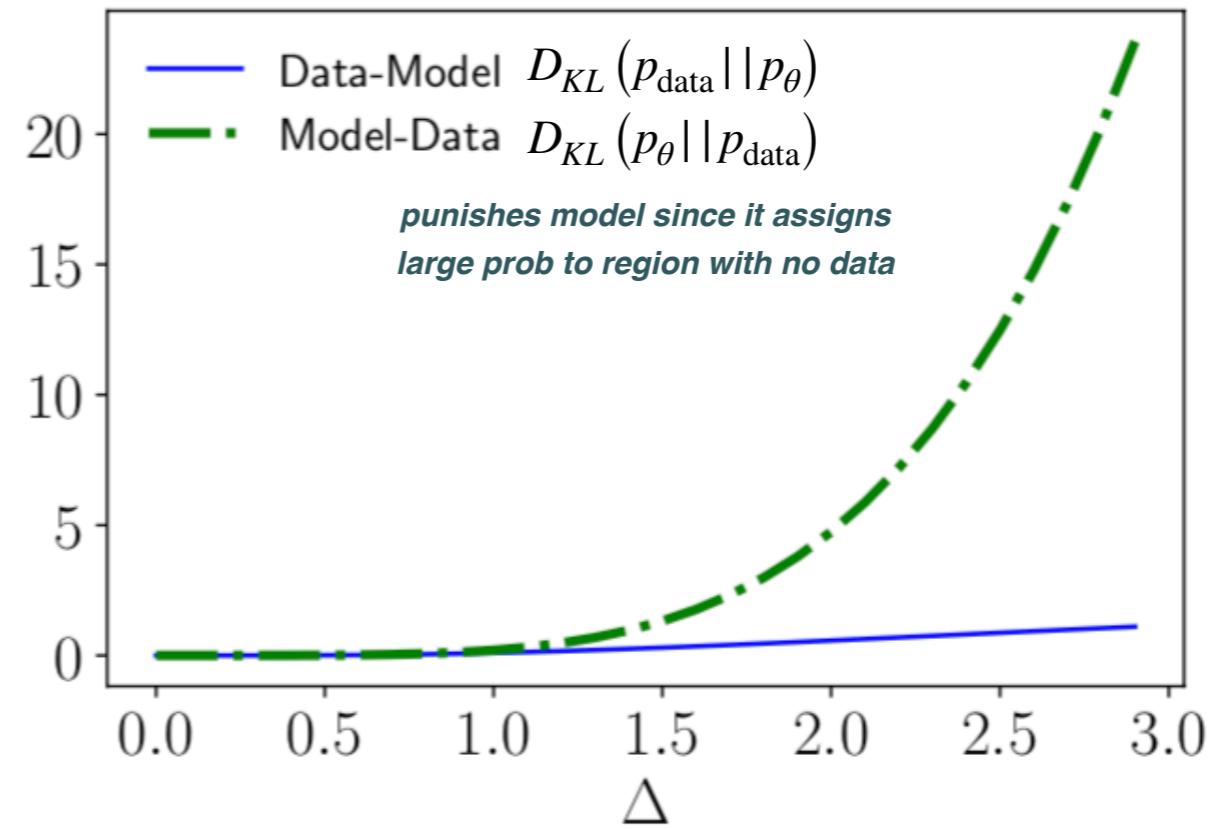
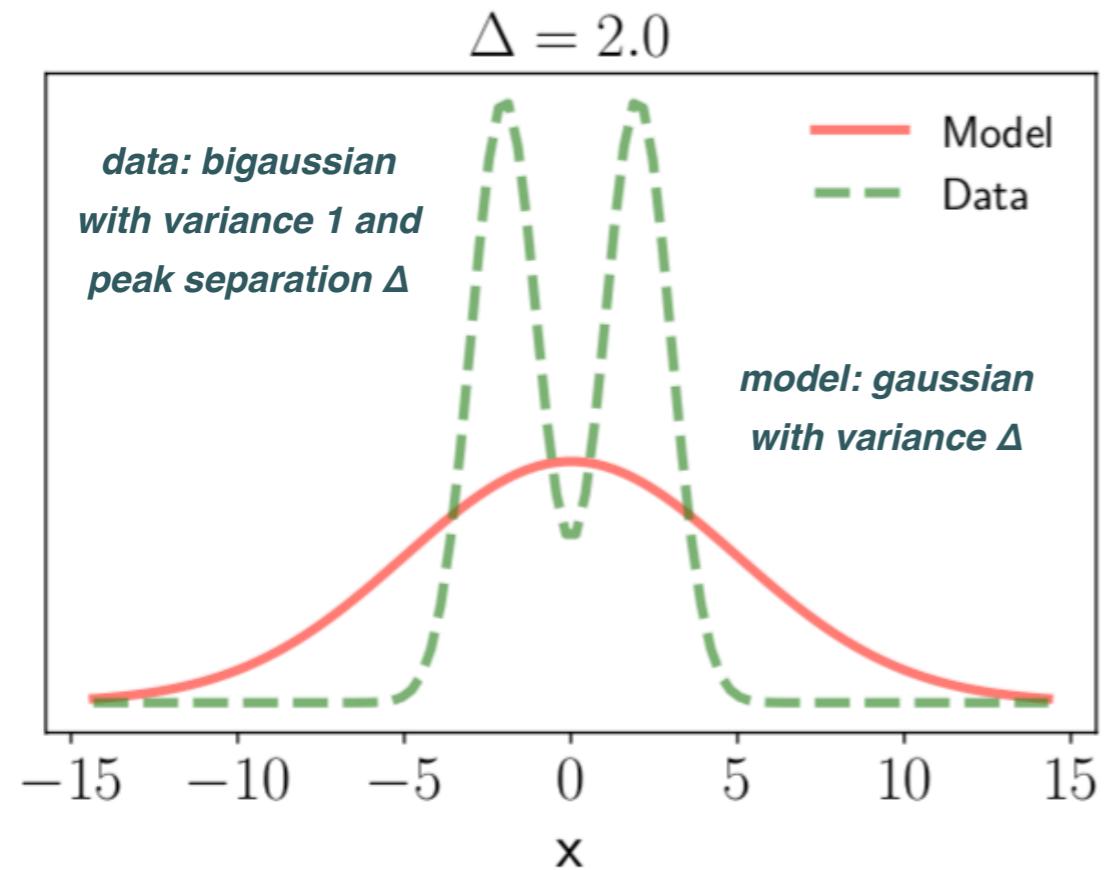
The KL-divergence is **positive-definite**, and only vanishes when  $p(\mathbf{x})=q(\mathbf{x})$

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

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

# Similarity of probability distributions

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

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

accounts for non-trivial features beyond the lowest moments of the distribution

*identifies when model gives too much weight to regions without data!*

# Statistical distances and similarity

another popular statistical distance is the **Kolmogorov-Smirnov statistic**

Assume we have  $n$  observations of the random variable  $x$ . The **empirical distribution function** is

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n I_{[-\infty, x]}(x_i)$$

*indicator function*

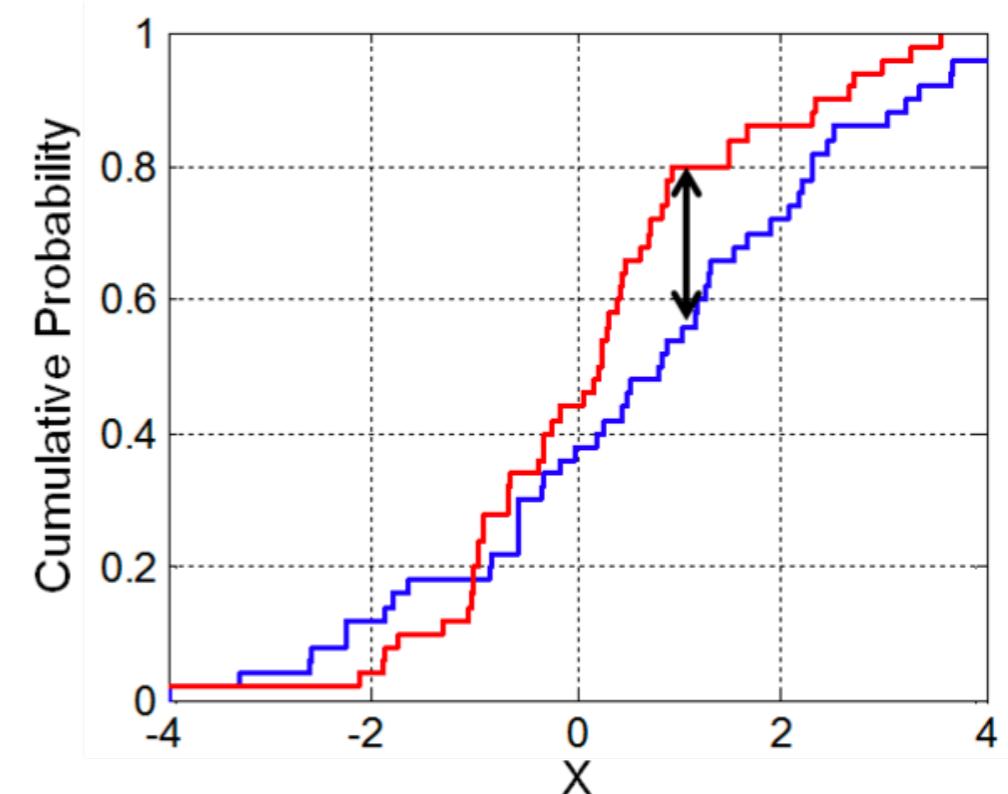
$$I_{[-\infty, x]}(x_i) = 1, \text{ for } x_i \leq x, \quad I_{[-\infty, x]}(x_i) = 0, \text{ for } x_i > 0$$

we can then define the **KS-statistic** between two probability distributions as

$$\text{KS} = \sup_x [F_{1,n}(x) - F_{2,m}(x)]$$

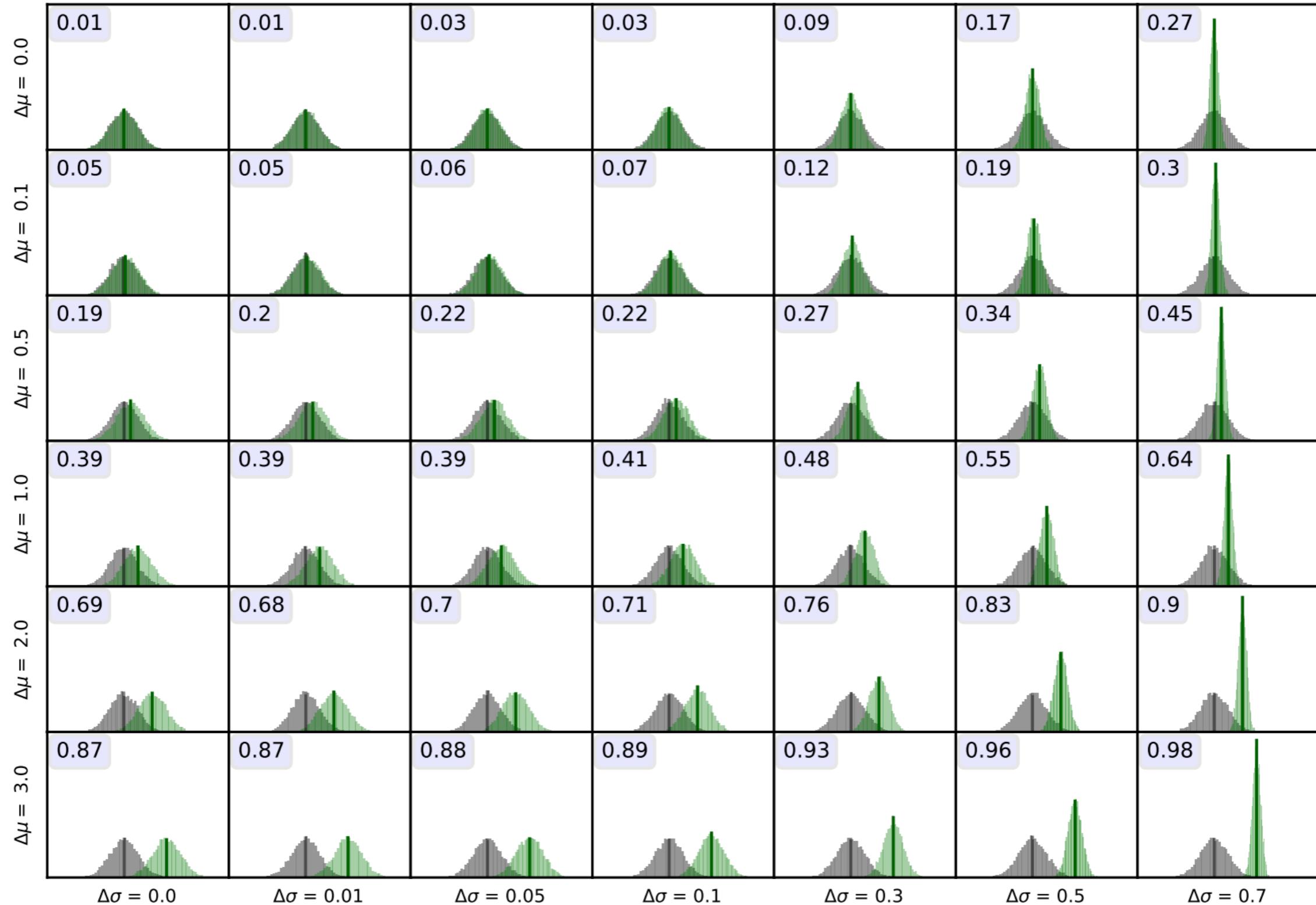
*EDF sample 1  
( $n$  elements)*      *EDF sample 2  
( $m$  elements)*

*the KS statistics is sensitive to all possible differences between the two distributions (which might be excessive in some cases)*



# Statistical distances and similarity

KS statistic for different means and std devs ( $N=10000, \mu_0=0, \sigma_0=1$ )



# **Generative Models**

# Generative Models

**Q:** what is the key property of **generative models** which makes them distinct from **discriminative models?**

# Generative Models

**Q:** what is the key property of **generative models** which makes them distinct from **discriminative models**?

*discriminative models: tell part  
muffins from chihuahuas*



*generative models: determine  
the probability distribution  
associated to “muffin” ...*



# Generative Models

**Q:** what is the key property of **generative models** which makes them distinct from **discriminative models**?

*discriminative models*

$$p(\mathcal{C}_k | x)$$

*probability that instance  $x$  belongs to class  $C_k$*

*generative models*

$$p(x | \mathcal{C}_k)$$

*probability distribution of  $x$  associated to the class  $C_k$*

What we can do with generative models that is not possible with discriminative models?

# Generative Models

let us go back to our discussion of **Decision Theory**. There we saw that a general classification problem can be separated into two distinct steps:

*posterior class probabilities*

- The **inference stage**, where a set of input examples is used to train a model for  $p(\mathcal{C}_k | \mathbf{x})$
- The **decision stage**, where the information on these posterior probabilities is used to make **optimal class assignments**

So far we focused on **discriminative models**, where some criterion (e.g. minimise misclassification) is used together with the posterior probabilities to assign each new instance to a class

Here discuss **generative models** which aim to model the **distribution in the space of inputs  $\mathbf{x}$** . The name stems because using this distribution one can **generate synthetic data points** in the input space

one benefit of this approach is that we access the marginal density in the space of input data,  $p(\mathbf{x})$ , which is specially useful to detect new data points that have low probability in the model: **outlier, anomaly, or novelty detection**

# Generative Models

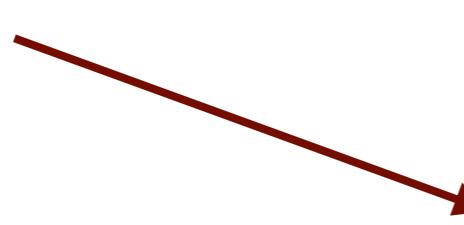
**generative models** can be better understood in the context of **classification**: assume we have a set of  $n$  instances  $\{\mathbf{x}\}$  with  $p$  features each that we want to classify into categories  $C_k$

- First of all, one must **solve the inference problem** and determine the **class-conditional probabilities** for each class individually

$$p(x | \mathcal{C}_k) \longrightarrow \text{probability that given a class } C_k \text{ the instance } x \text{ is found in it}$$

- Then separately infer the **prior class probabilities**  $p(\mathcal{C}_k)$

- Determine the **posterior class probabilities** by means of Bayes' Theorem

$$p(\mathcal{C}_k | x) = \frac{p(x | \mathcal{C}_k) p(\mathcal{C}_k)}{p(x)}$$


*probability that  $x$  belongs to class  $C_k$*

$$p(x) = \sum_k p(x | \mathcal{C}_k) p(\mathcal{C}_k) \longrightarrow \text{probability to generate the instance } x$$

# Generative Models

$$p(\mathcal{C}_k | \mathbf{x}) = \frac{p(\mathbf{x} | \mathcal{C}_k) p(\mathcal{C}_k)}{p(\mathbf{x})}$$

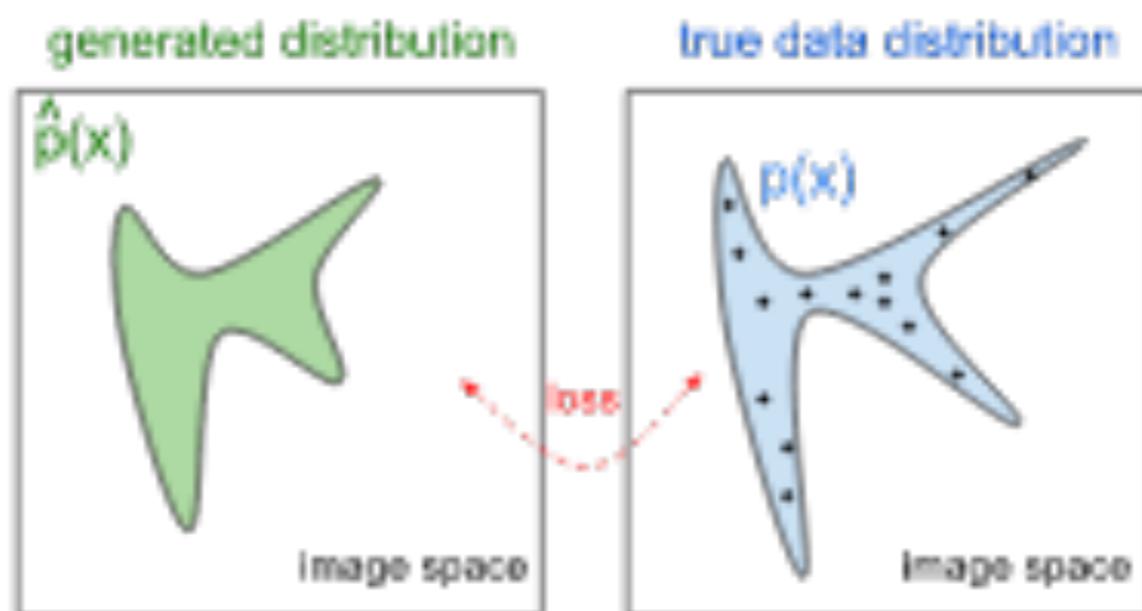
*probability that  $\mathbf{x}$  belongs to class  $C_k$*

$$p(\mathbf{x}) = \sum_k p(\mathbf{x} | \mathcal{C}_k) p(\mathcal{C}_k)$$

*probability to generate the instance  $\mathbf{x}$*

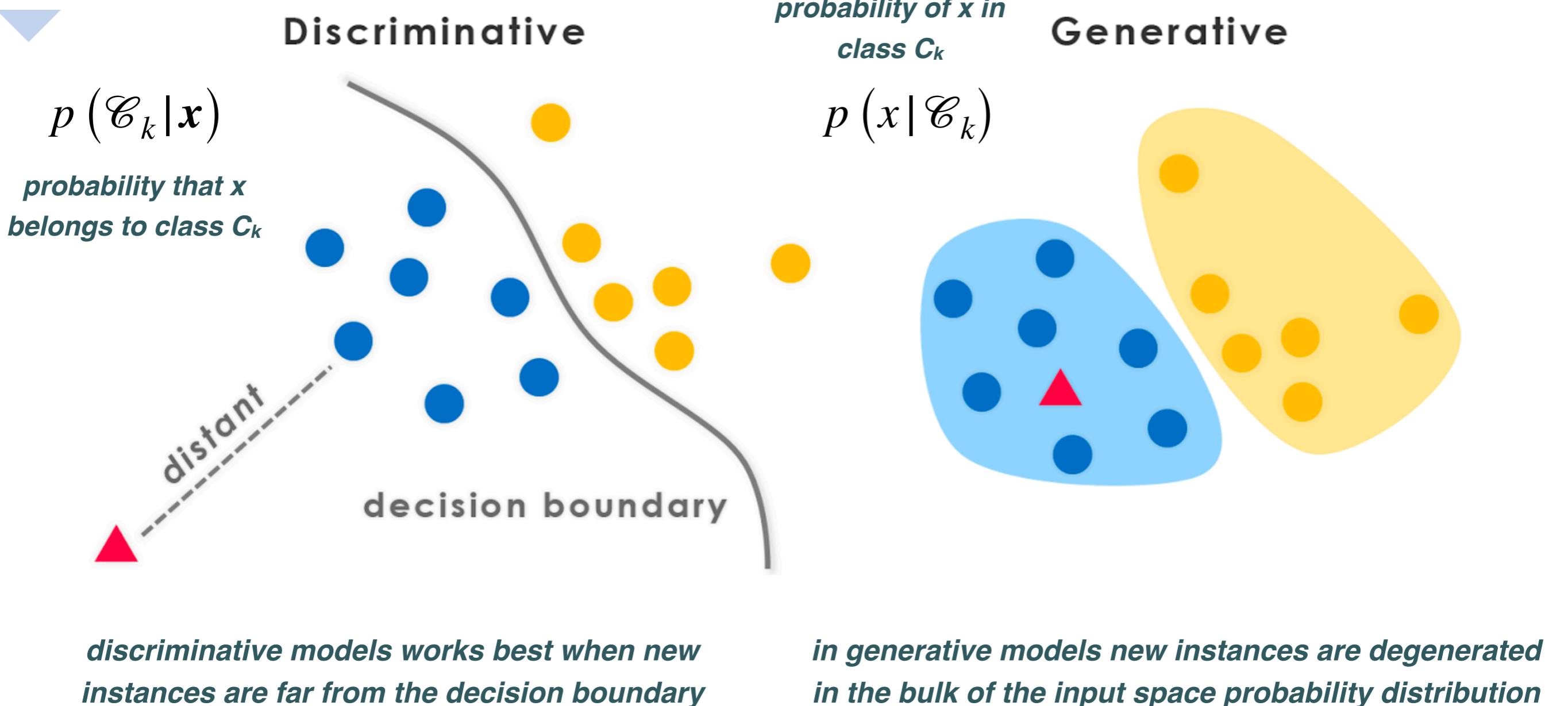
given that we now have the posterior probabilities, we can use **Decision Theory** to determine the class membership for each new instance  $\mathbf{x}$

the key aspect of generative models is that we have access to  $p(\mathbf{x})$ , the probability distribution in the input data space: sampling from  $p(\mathbf{x})$  one can then **generate new synthetic data points** in input space



# Generative Models

*Compare graphically Discriminative and Generative Models for classification*



# Probabilistic generative models

based on these ideas let us construct an explicit **probabilistic generative model** in the context of classification problems with **two categories**. In this case Bayes' Theorem reads

*posterior class  
probabilities*

$$p(\mathcal{C}_k | \mathbf{x}) = \frac{p(\mathbf{x} | \mathcal{C}_k) p(\mathcal{C}_k)}{p(\mathbf{x})}$$

$$p(\mathcal{C}_1 | \mathbf{x}) = \frac{p(\mathbf{x} | \mathcal{C}_1) p(\mathcal{C}_1)}{p(\mathbf{x} | \mathcal{C}_1) p(\mathcal{C}_1) + p(\mathbf{x} | \mathcal{C}_2) p(\mathcal{C}_2)}$$

which can be simplified as

$$p(\mathcal{C}_1 | \mathbf{x}) = \frac{1}{1 + \exp(-a)} = \sigma(a) \quad a = \ln \frac{p(\mathbf{x} | \mathcal{C}_1) p(\mathcal{C}_1)}{p(\mathbf{x} | \mathcal{C}_2) p(\mathcal{C}_2)}$$

in terms of the **logistic sigmoid function** which we discussed in the context of neural nets

# Probabilistic generative models

Note that once we have trained such a model we can do two things:

- Use **posterior probabilities** with **decision theory** to determine **decision boundaries**

$$p(\mathcal{C}_k | \mathbf{x}) + \text{decision theory} \rightarrow \text{class boundaries}$$

- Use the **probability distribution in the space of input data** to generate synthetic samples

$$p(\mathbf{x}) \rightarrow \text{generate new instances of } \mathbf{x}$$

For  $K > 2$  classes the posterior probabilities are given by

$$p(\mathcal{C}_k | \mathbf{x}) = \frac{\exp(a_k)}{\sum_j \exp(a_j)} \quad a_k = \ln(p(\mathbf{x} | \mathcal{C}_k)p(\mathcal{C}_k))$$



*normalised exponential of softmax function*

*(smoothed version of the maximum function)*

# Generative Models

we can summarise the main ideas **underlying generative models** as follows

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*

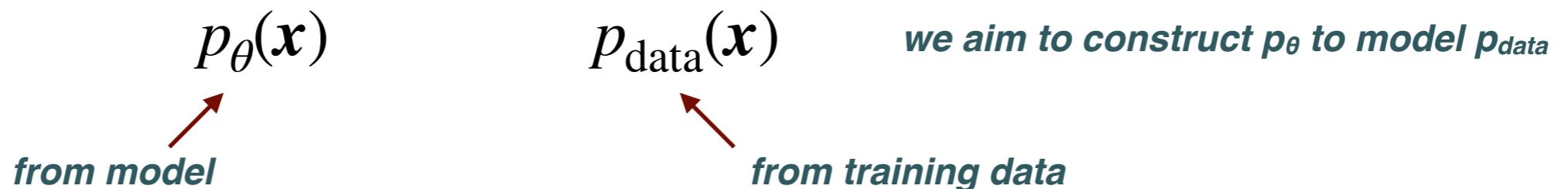
*e.g. generate new samples of a given phase  
of the Ising model*

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

# **Adversarial Learning and Generative Adversarial Networks**

# Maximising similarity

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

Kullback-Leibler divergence

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

Entropy

*Expected value of model probability given the data*

# Maximising similarity

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

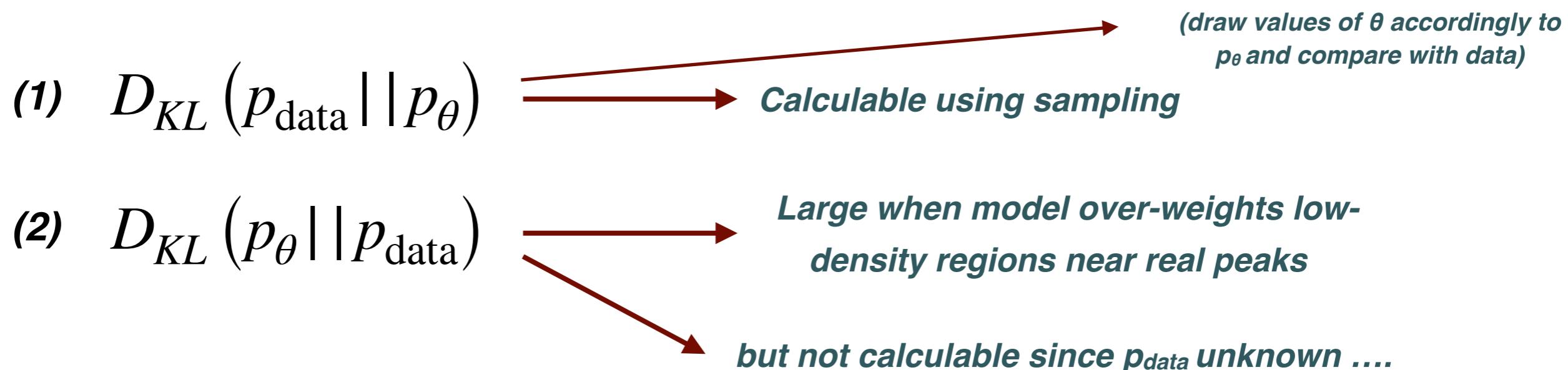
$$\langle \log p_\theta(x) \rangle_{\text{data}} = S_p[p_{\text{data}}] - D_{KL}(p_{\text{data}} \parallel p_\theta)$$

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

by minimising the KL divergence between two distributions, we can build a model that reproduces the **input data distribution**

# Adversarial Learning

A subtle point concerns which of the two versions of the KL-divergence to minimise:



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

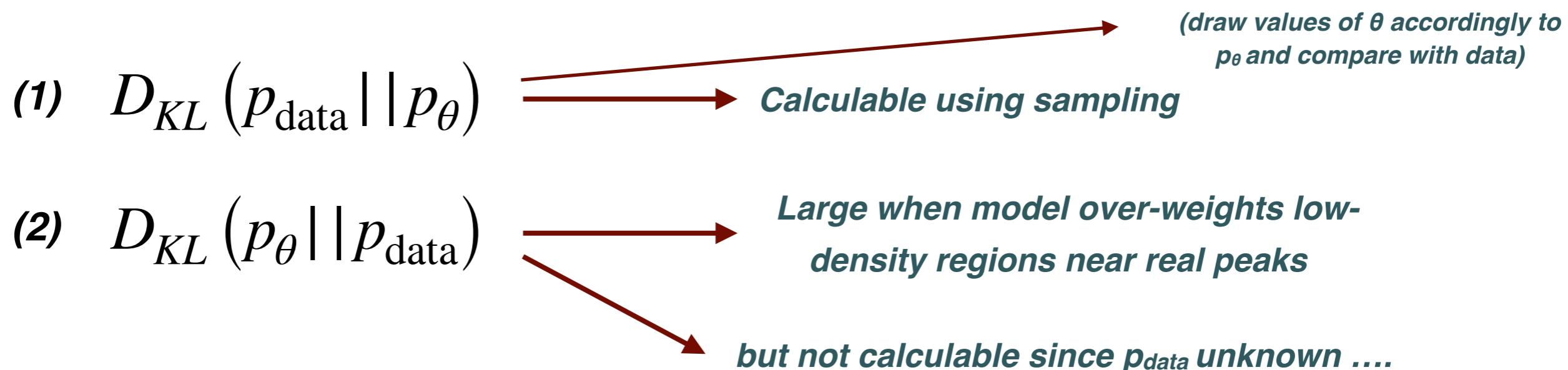
By punishing the model for generating points that can be easily discriminated from the data, Adversarial Learning decreases the **weight of regions in the model space that are far away from data points**, regions that inevitably arise when maximising the likelihood

# Adversarial Learning

**Q:** do you think that Adversarial Learning belongs to the family of **Supervised** or of **Unsupervised** machine learning algorithms?

# Adversarial Learning

A subtle point concerns which of the two versions of the KL-divergence to minimise:



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

**Adversarial Learning** is a type of **Unsupervised Machine Learning**: there is **no cost function** and the examples are **unlabelled**: the task here is to learn how to **generate new examples** using the existing ones as input!

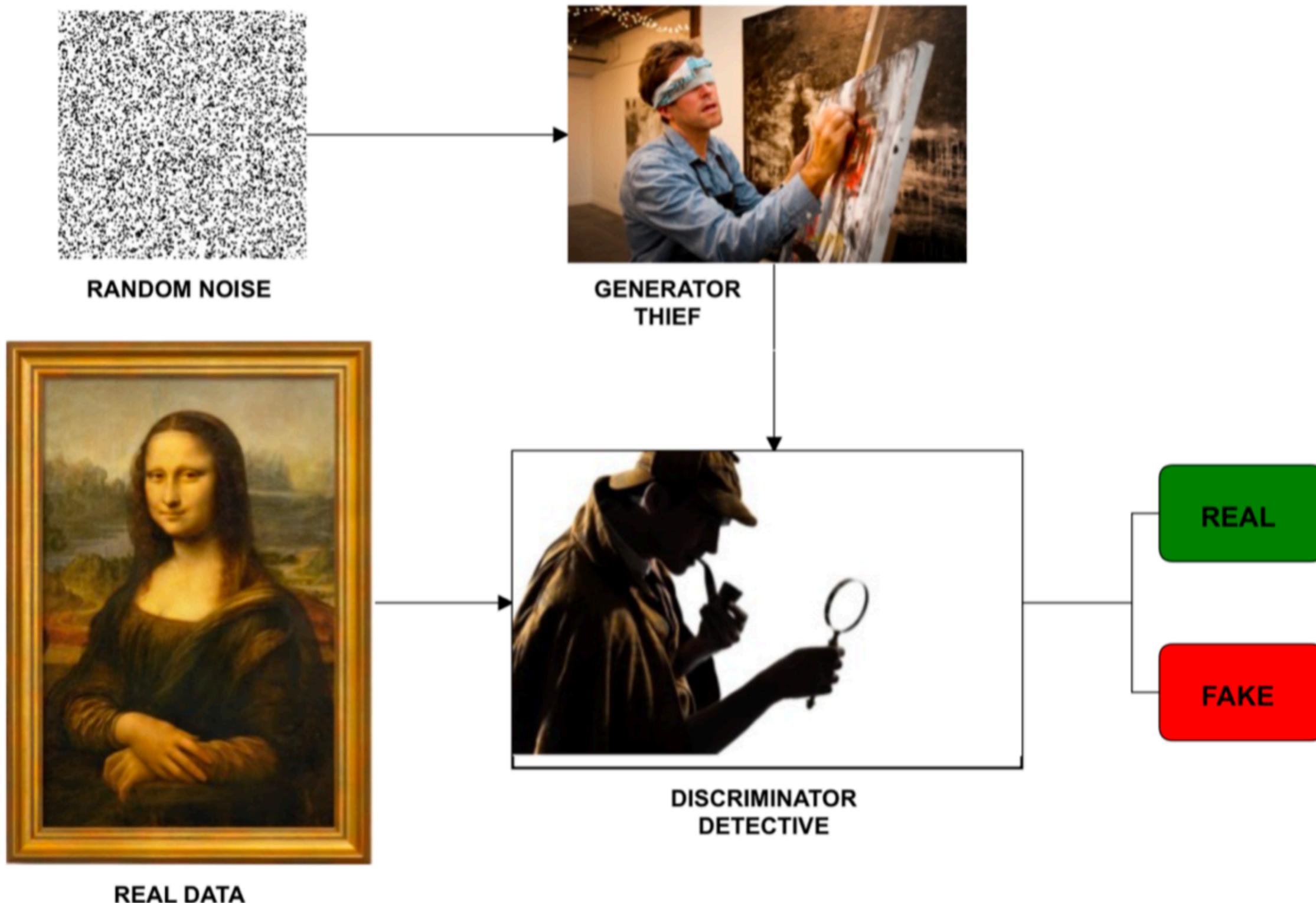
# Generative adversarial networks

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**

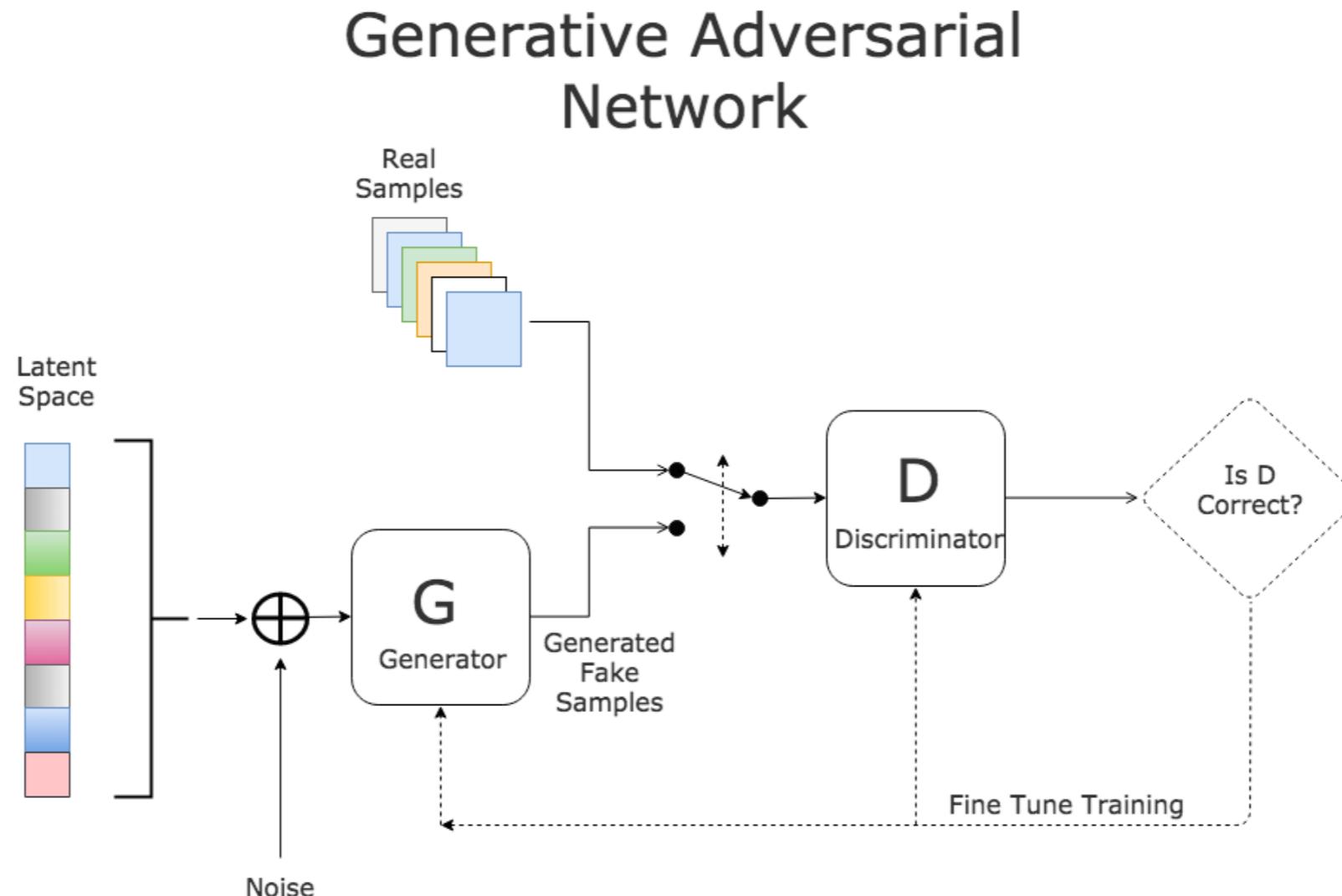
the generator network **G** should be trained to maximise the probability that the discriminator network **D** makes a mistake: that is, **G** should generate pseudo-data samples that are virtually **indistinguishable** from the actual data

# Generative adversarial networks



# Generative Adversarial Networks

- 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.



# 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 D      output of D when input a real data sample      output of D when input a "fake" data sample produced by G  
NN params of 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**  $\{z_n\}_{n=1}^N$

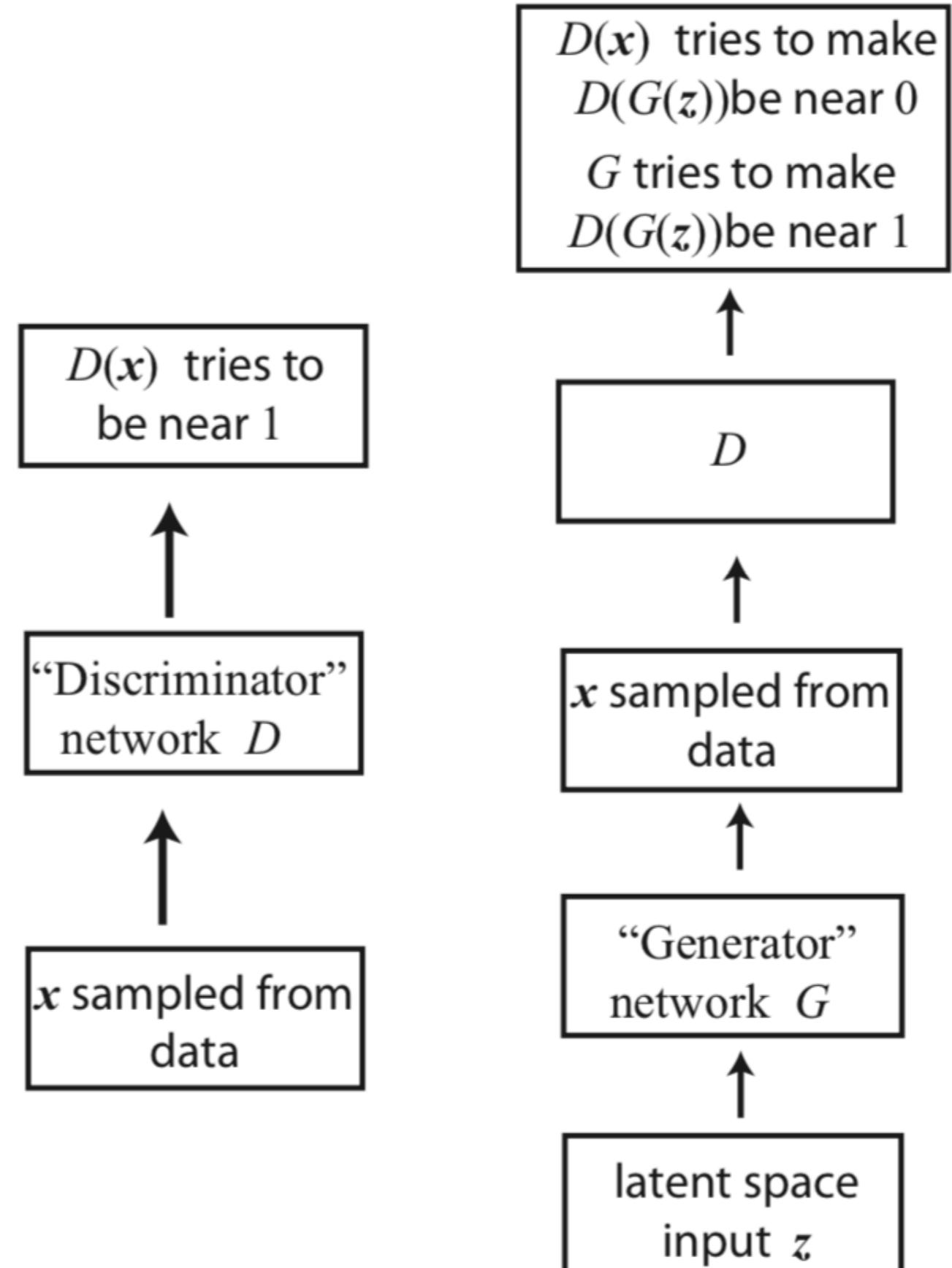
$$C(\theta_D, \theta_G) = \frac{1}{N} \sum_{n=1}^N \log(1 - D(G(z_i)))$$

*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

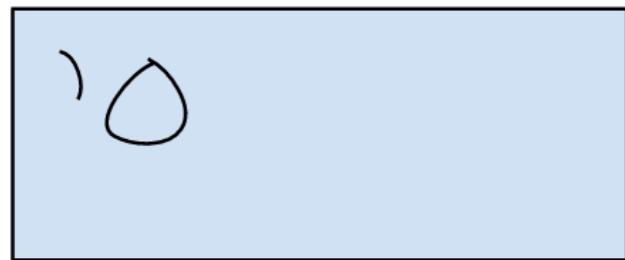
the generator and discriminator are sequentially trained and iterated until convergence is achieved, at this point **D** cannot tell apart the pseudo-data from **G** from the real data



# GAN training

*initial training*

Generated Data



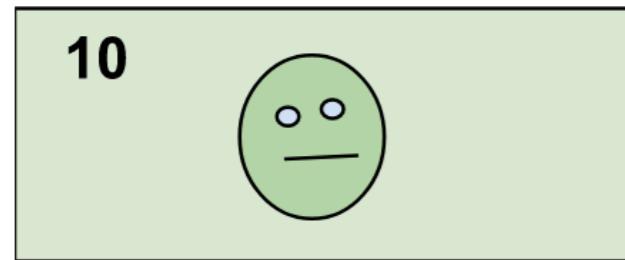
Discriminator

FAKE      REAL



Real Data

*intermediate training*



FAKE      REAL



*long training*



REAL      REAL



*convergence!!*

# Image generation with GANs



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

# Image generation with GANs

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

# Image generation with GANs



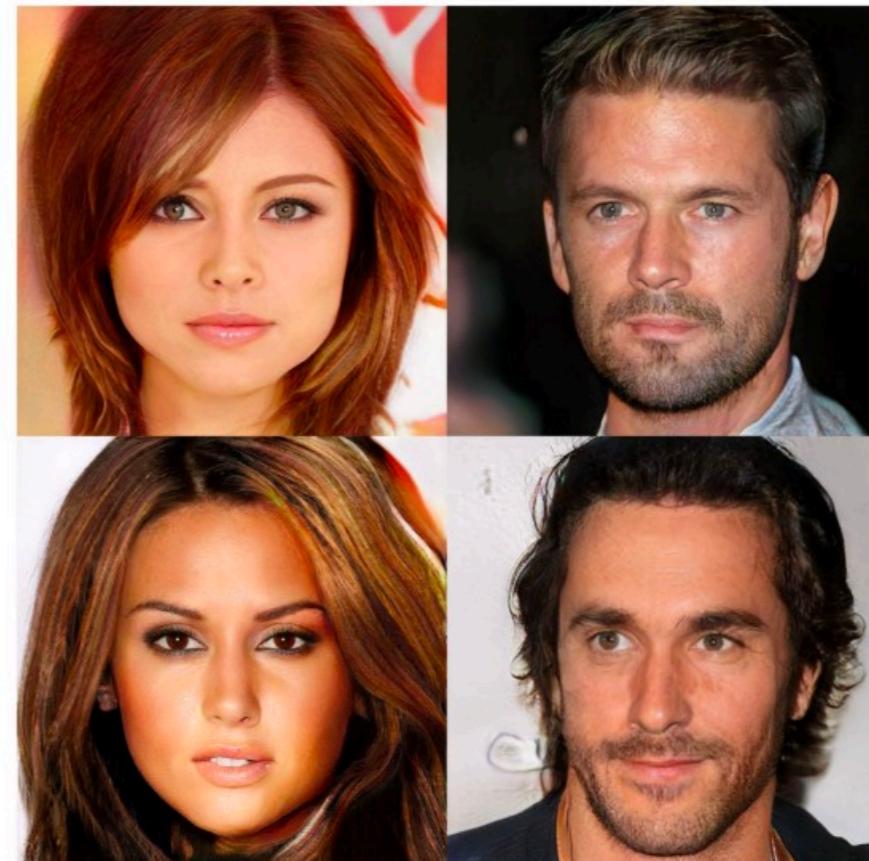
DCGAN  
11/2015



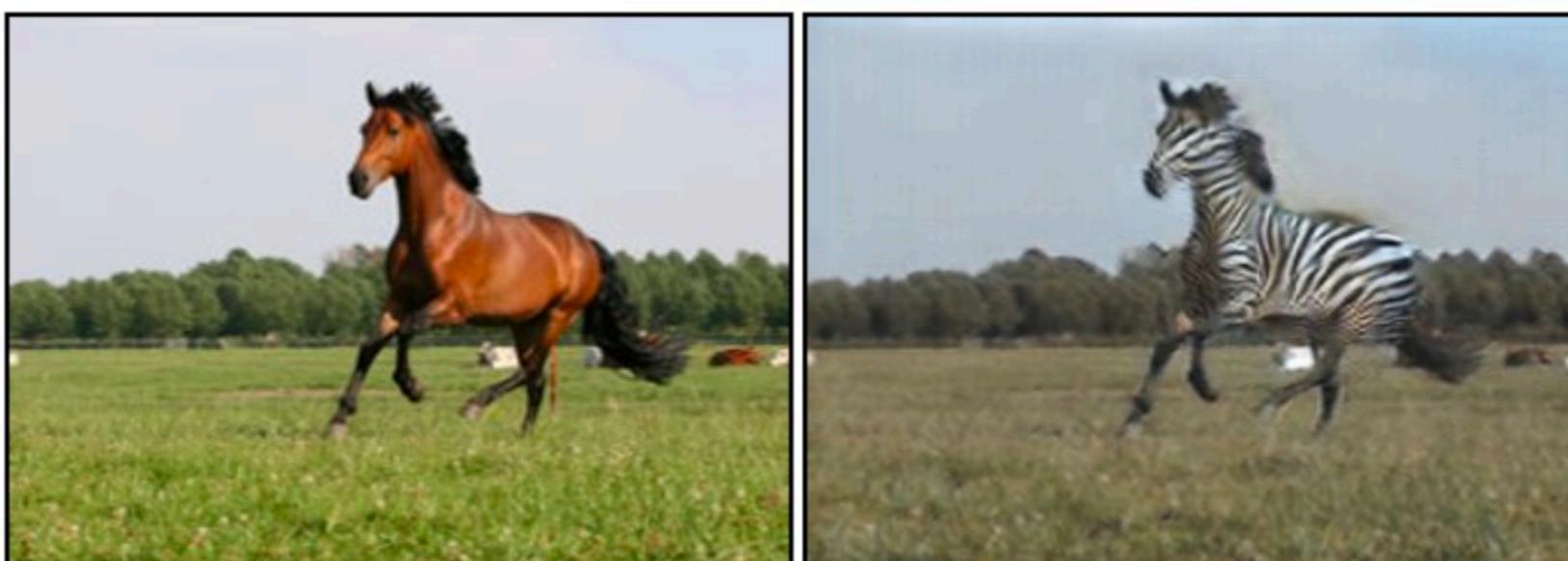
EBGAN-PT  
9/2016



BEGAN  
3/2017  
128 × 128



Progressive GAN  
10/2017  
1024 × 1024



horse → zebra

# Image generation with GANs

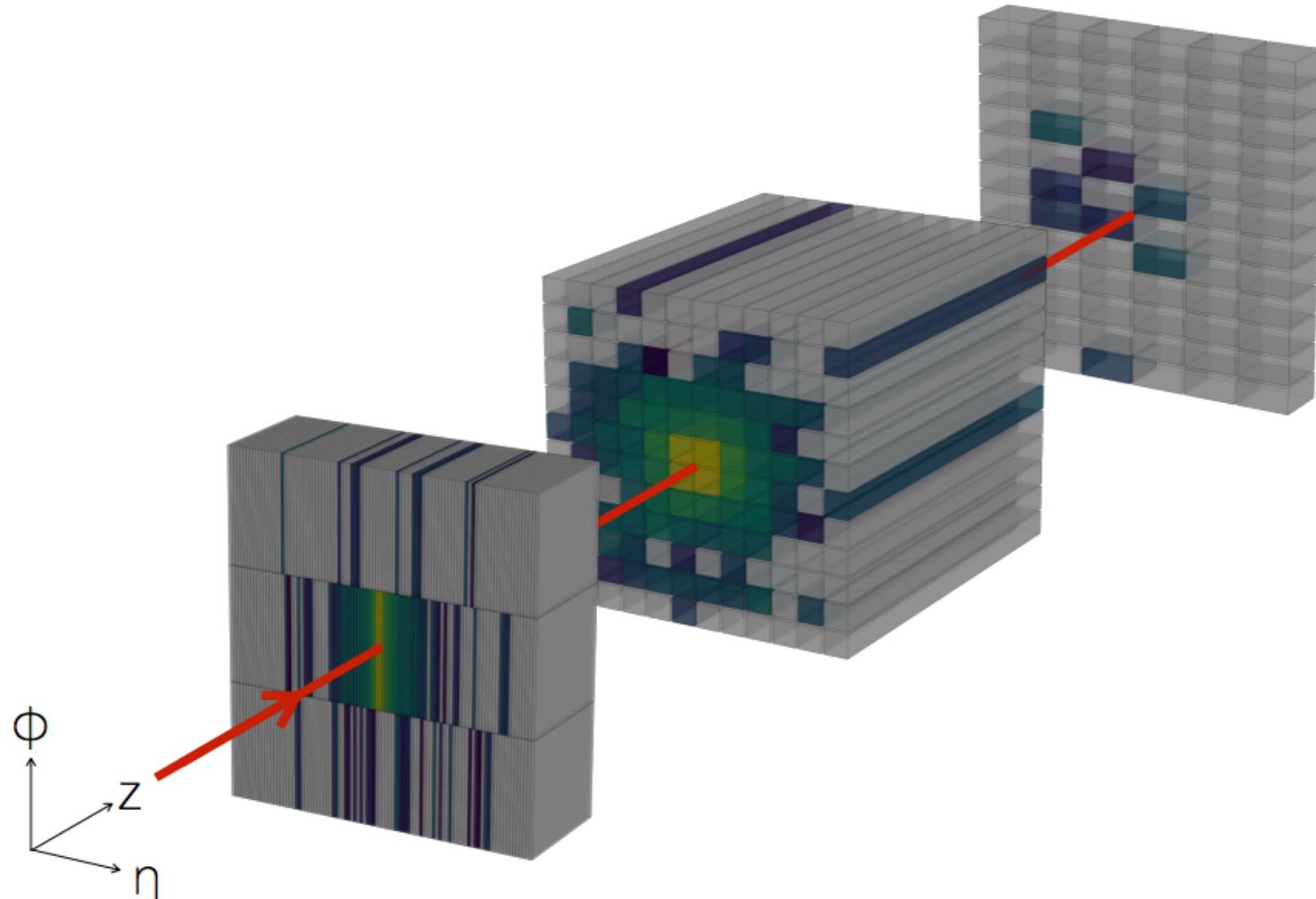


<https://deeppdreamgenerator.com>

# Generative Adversarial Networks

# GANs for detector simulation in HEP

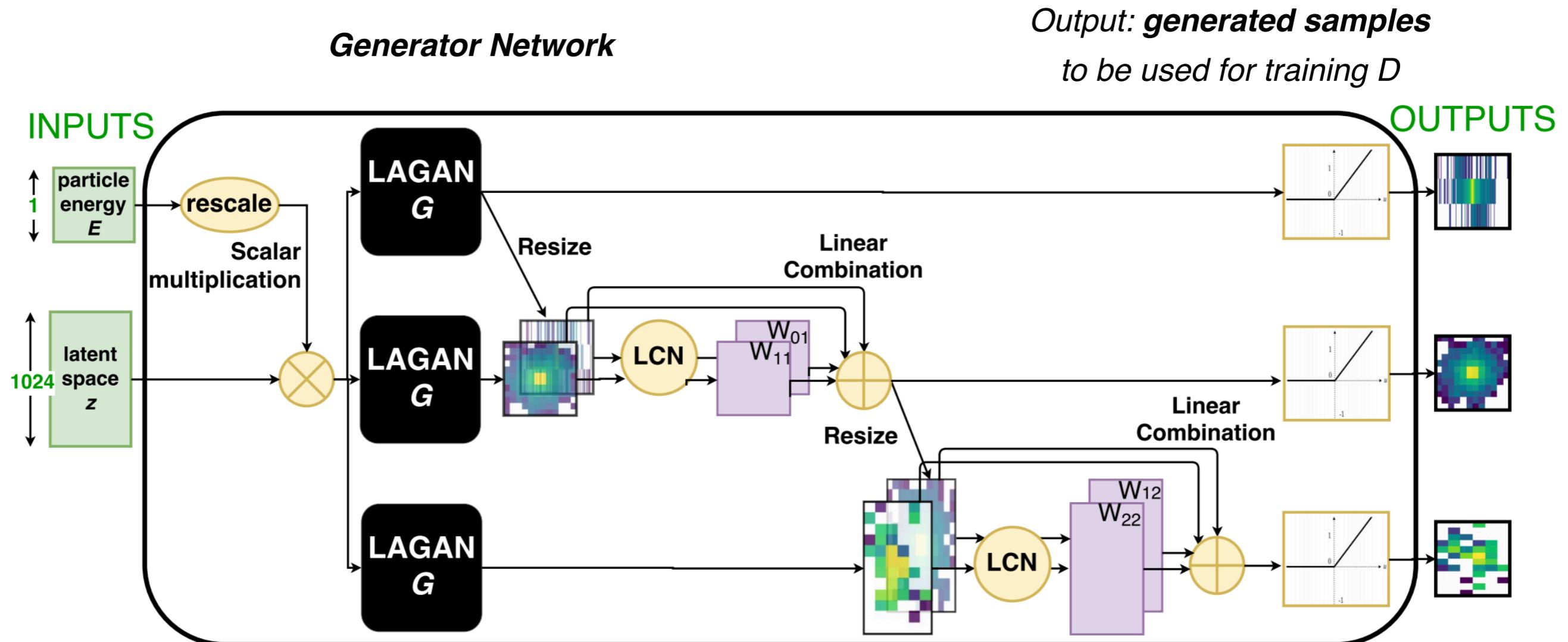
- Modelling accurately the response of detectors with the **propagation of high energy particles** is an essential task for present and future HEP experiment
- Detector simulation at the LHC** is a very CPU-intensive task, dominated by modelling of particle showers inside calorimeters
- Generative Adversarial Networks** can speed up detector simulation by orders of magnitude



*Task: to efficiently model the **propagation of high energy particles** (and their interaction) within the layers of electromagnetic and hadronic calorimeters*

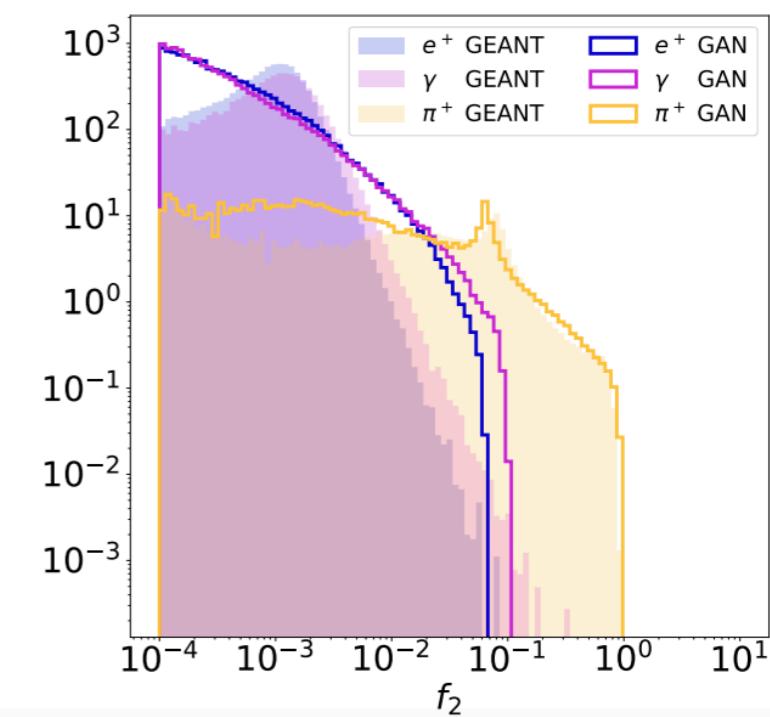
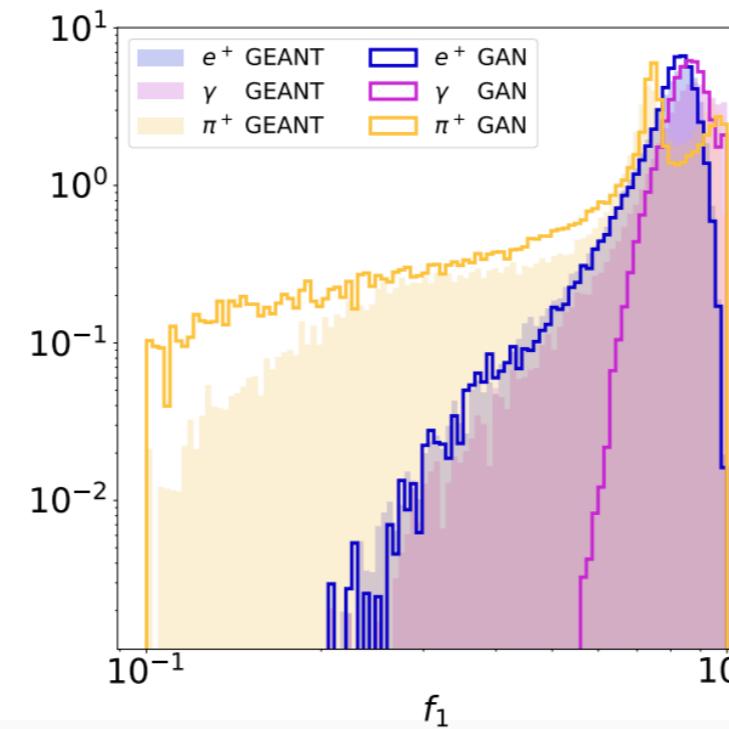
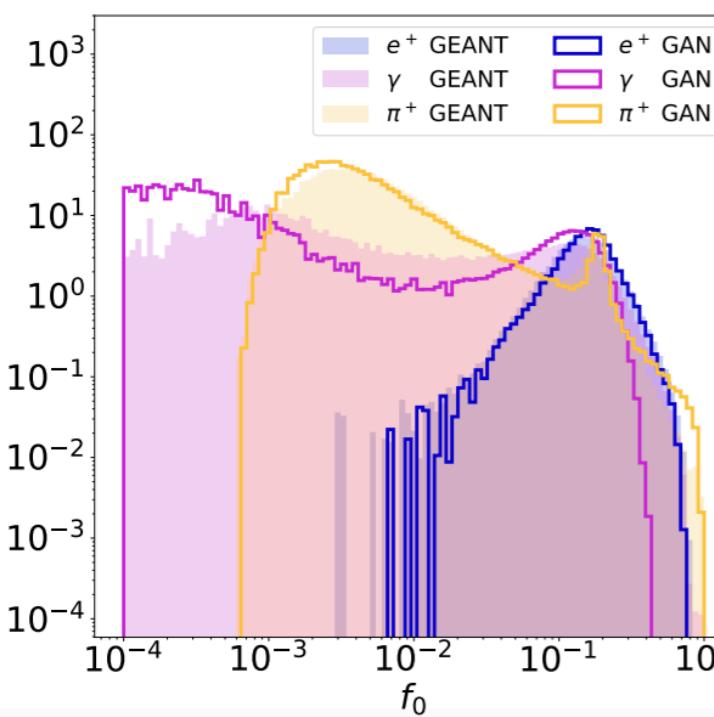
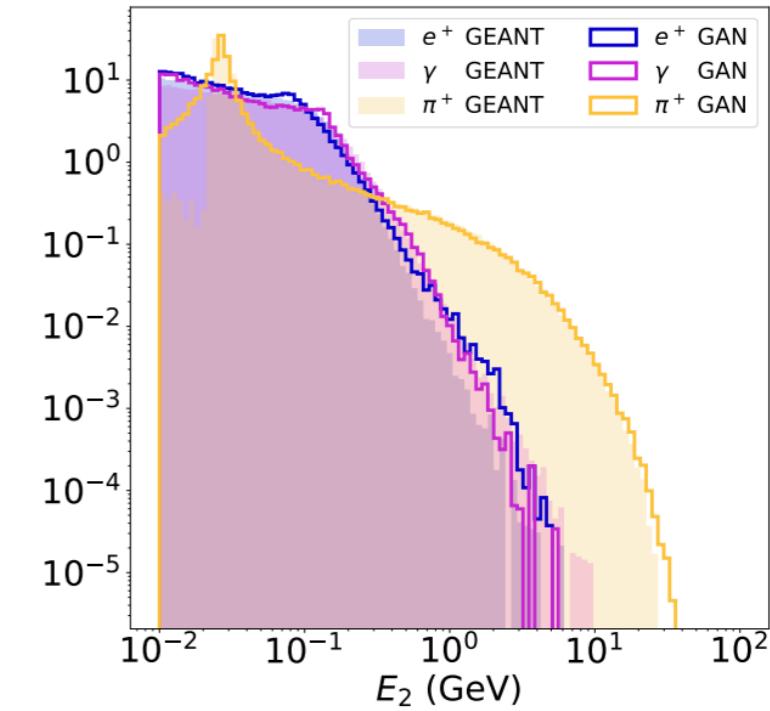
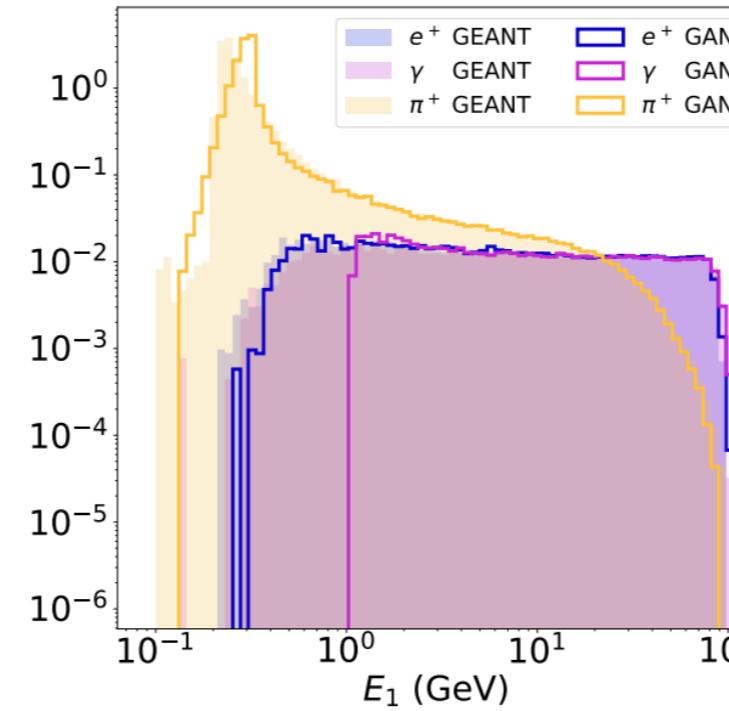
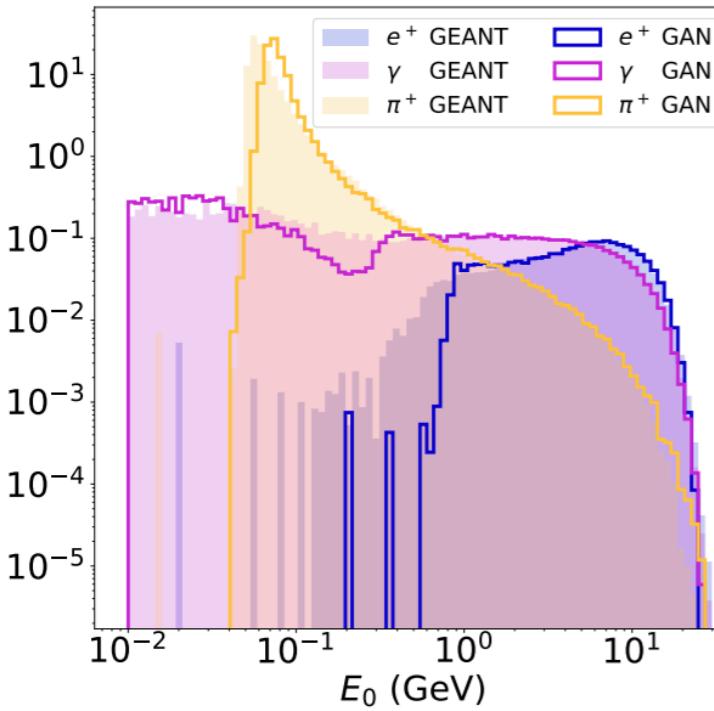
# GANs for detector simulation in HEP

- Use GANs as a tool to **speed up full simulation of particle showers** in a HEP calorimeter
- The generator G learns a map from **a latent space** to space of **generated samples** for training
- Carefully understanding the **underlying physics of particle propagation** in a detector is crucial to optimise the training strategy, e.g. relationships between neighbouring detector layers



Paganini et al. 17

# GANs for detector simulation in HEP

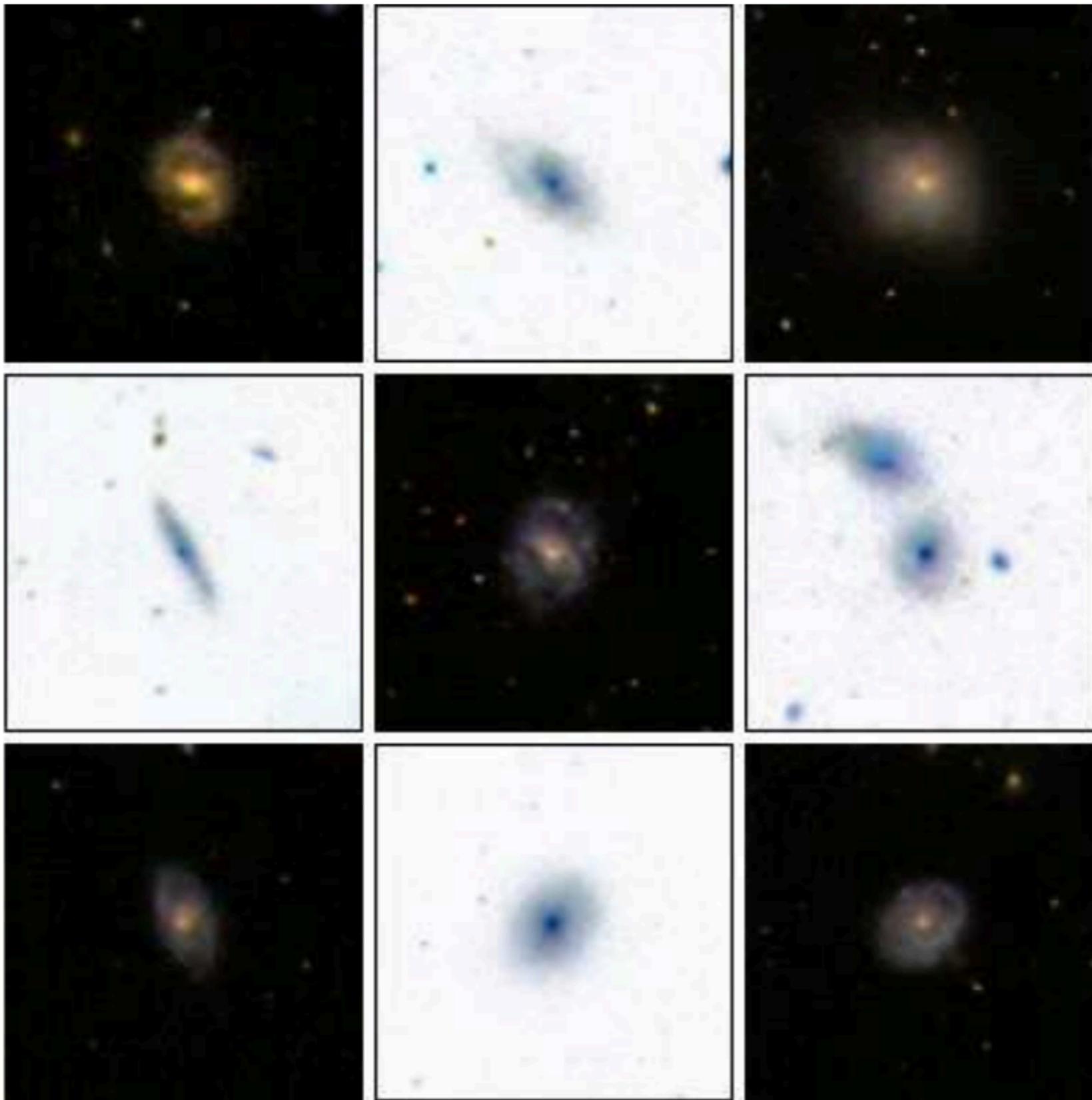


# GANs for detector simulation in HEP

Simulator	Hardware	Batch Size	ms/shower
GEANT4	CPU	N/A	1772
CALOGAN	CPU	1	13.1
		10	5.11
		128	2.19
		1024	2.03
CALOGAN	GPU	1	14.5
		4	3.68
		128	0.021
		512	0.014
		1024	0.012

*Speed-up by several orders of magnitude, specially when running in GPUs*

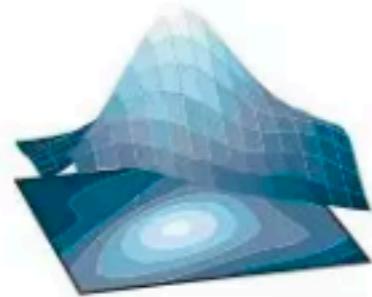
# Synthetic Galaxies from GANs



*Useful for applications  
where a large number  
of synthetic images of  
galaxies speed up  
analysis*

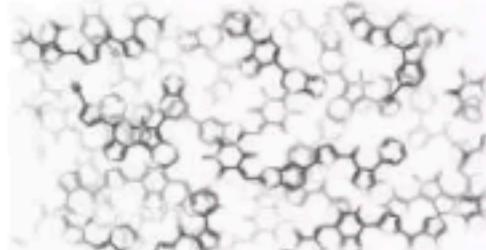
# GANs in chemistry

## Functional space



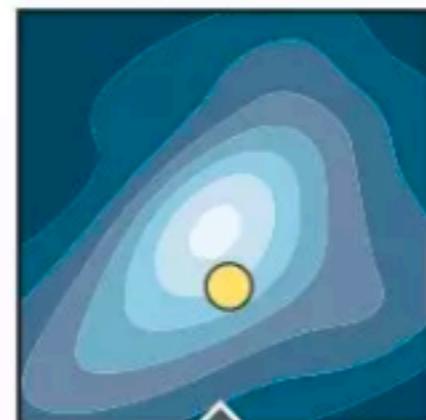
Desired properties (redox potential, solubility, toxicity)

## Chemical space

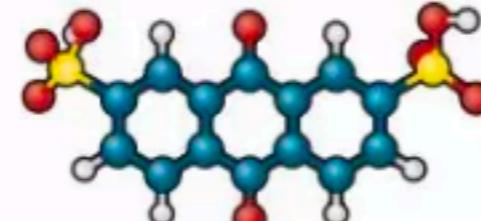


(Drug-like, photovoltaics, polymers, dyes)

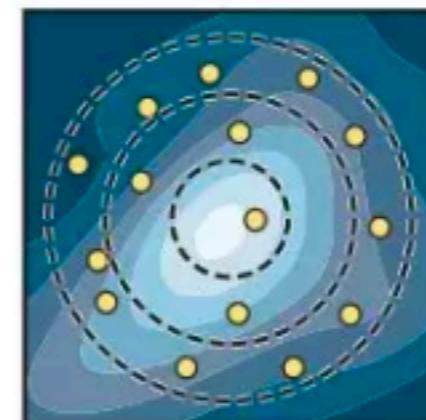
## Direct



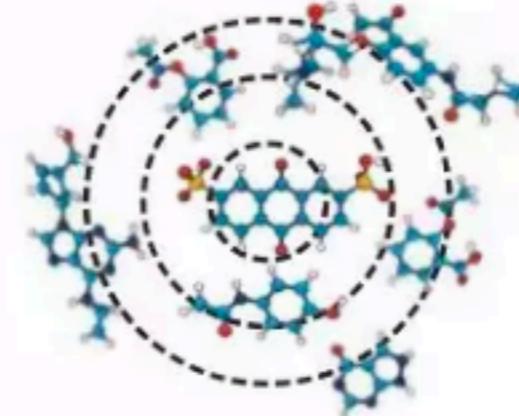
Experiment or simulation (Schrödinger equation)



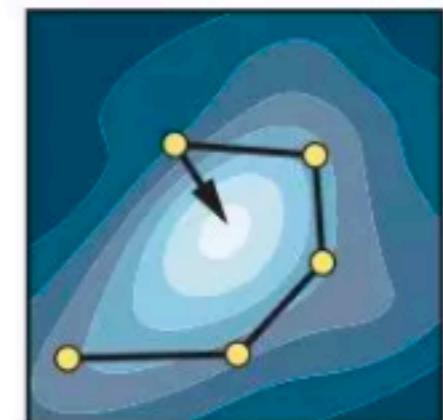
## Inverse



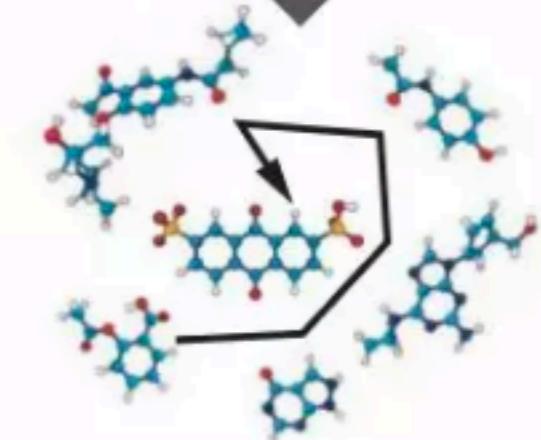
High-throughput virtual screening (e.g., with 3 filtering stages)



## Inverse



Optimization, evolutionary strategies, generative models (VAE, GAN, RL)



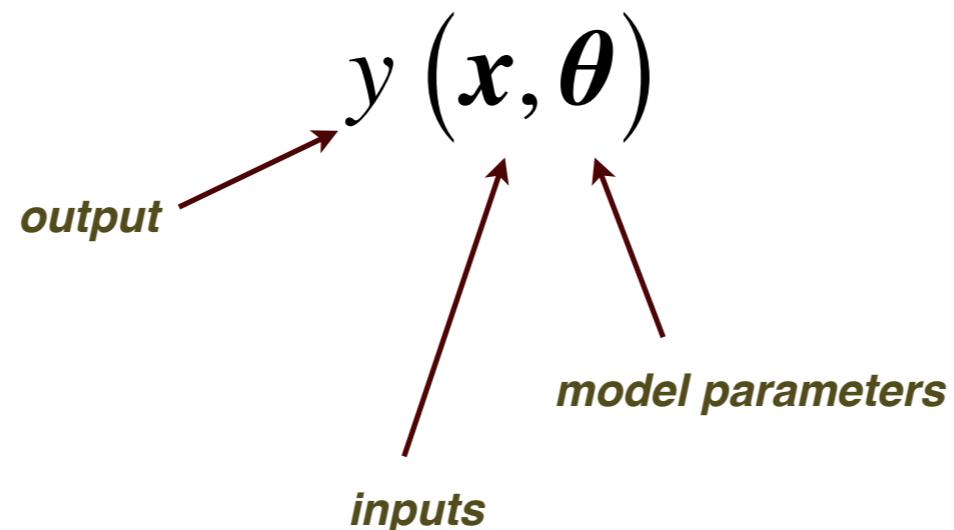
*Start from molecule, evaluate space of properties*

*Start from space of properties, design new molecules*

# **Kernel Methods**

# Kernel Methods

most of the ML models we have considered for regression and classification are based on the **parametric mapping** (linear or non-linear) between inputs and outputs



after the training, the input examples are **discarded** and not relevant for subsequent **predictions**: the information contained in the original training examples is now **entirely contained** on the model parameters

# Kernel Methods

most of the ML models we have considered for regression and classification are based on the **parametric mapping** (linear or non-linear) between inputs and outputs

$$y(x, \theta) \leftarrow \text{info on training samples entirely contained in model params}$$

were the **training examples are discarded** after the model parameters (or their posterior distribution) have been determined

here we introduce models where the training dataset is also used in the **prediction phase**

# Feature space mapping

A key concept of **kernel methods** is that of a **feature space mapping**

This is a transformation, in general non-linear, between the **input data space  $X$**  and a **feature space  $F$**  of the same or different dimensions:

$$F = \{\phi(x) : x \in X\}$$

*Example of linear feature-space transformation*

$x$	$\phi(x)$
<i>Height</i>	<i>Height+Weight</i>
<i>Weight</i>	<i>Weight-Height</i>
<i>Age</i>	<i>Age + 2*Weight</i>

*Example of non-linear feature-space transformation*

$x$	$\phi(x)$
<i>Height</i>	<i>Height/Weight</i>
<i>Weight</i>	<i>Weight+Height</i>
<i>Age</i>	<i>Age * Weight</i>

such mapping can be chosen to **increase the efficiency** of the ML model training

*note that the dimensions of feature space can be much larger than of data space*

# Feature space mapping

most of the ML models we have considered for regression and classification are based on the **parametric mapping** (linear or non-linear) between inputs and outputs

$$y(x, \theta) \leftarrow \text{info on training samples entirely contained in model params}$$

were the **training examples are discarded** after the model parameters (or their posterior distribution) have been determined

here we introduce models where the training dataset is also used in the **prediction phase**

the dimensions of the **feature space mapping** are in general different from those of input space

$$x \longrightarrow \phi(x)$$

*Height*

*Height/Weight*

*Weight*

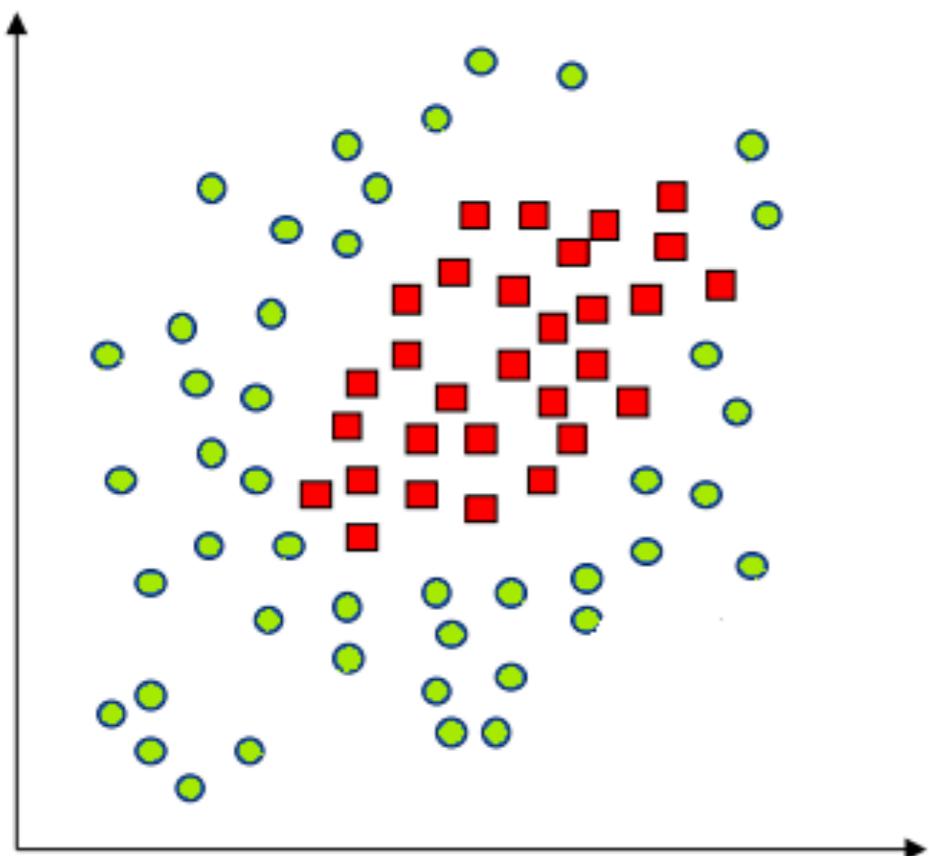
*Weight + Height*

*Age*

*Age \* Weight*

*Height \* Weight+Age*

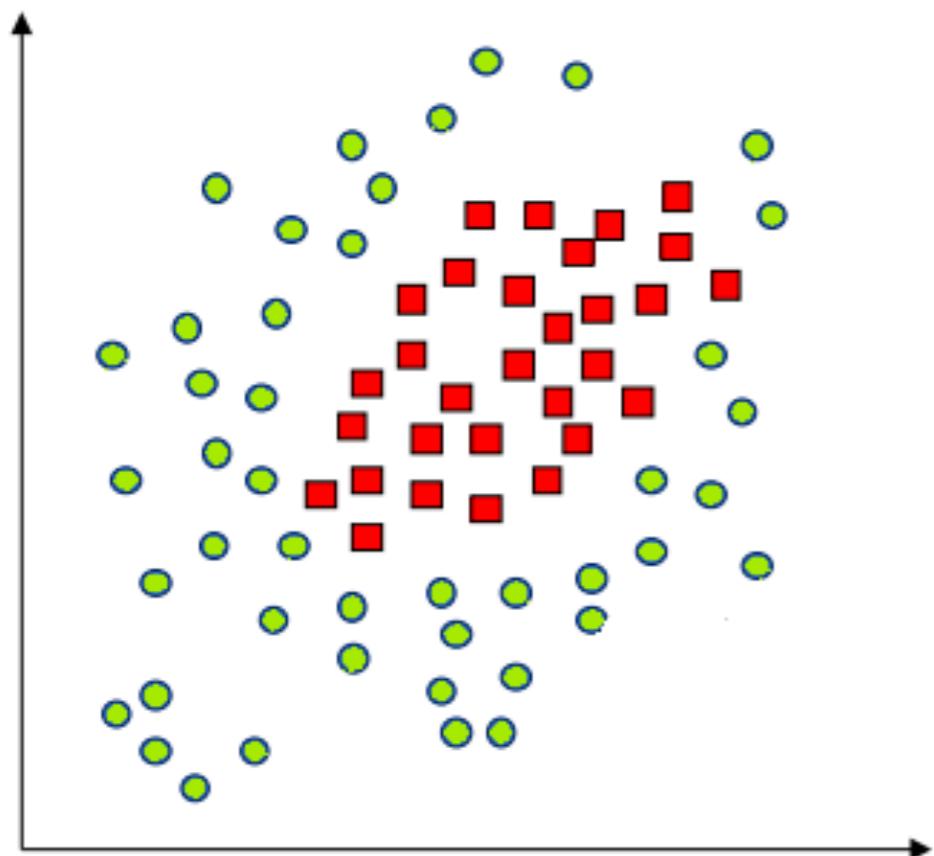
# Feature space mapping



*Data Space ( $p=2$ ):  
linear classification impossible*

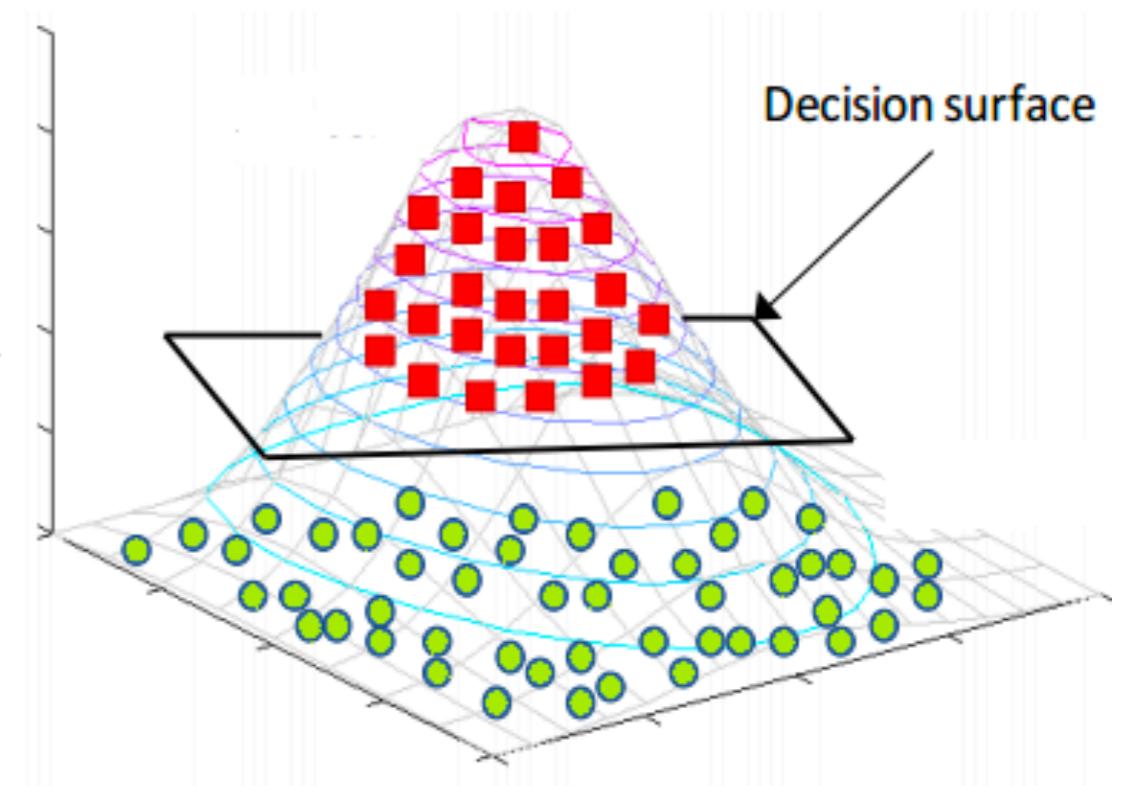
*Q: can we go to a higher-dimensional  
“feature space” where linear classification  
is possible?*

# Feature space mapping



*Data Space ( $p=2$ ):  
linear classification impossible*

**kernel**



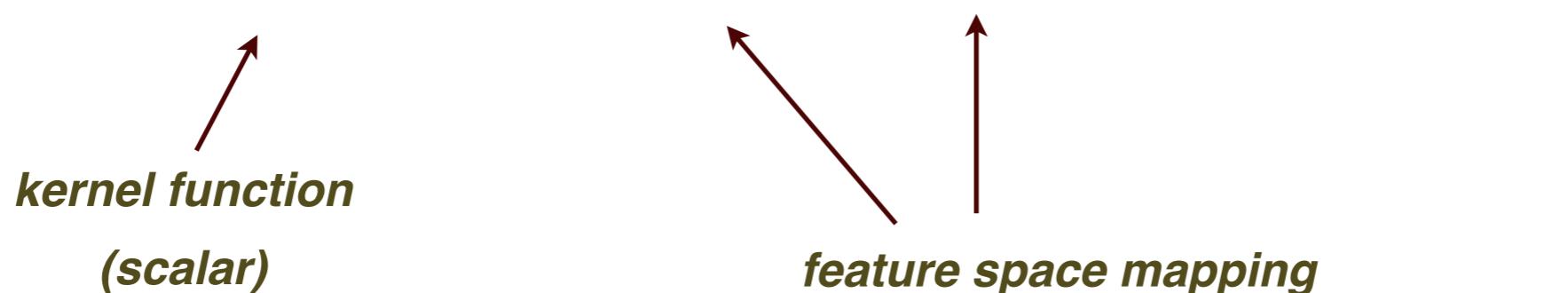
*Feature space ( $p=3$ ):  
linear classification possible!*

# Kernel Methods

the key ingredient of **kernel methods** is the **kernel function** evaluated at the training points

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$$

*inner product  
in feature space*



*kernel function  
(scalar)*

*feature space mapping*

the simplest kernel is the **linear kernel** (identity mapping in feature space)

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}' = \sum_{k=1}^p x_i x'_i$$

*sum over p,  
dimension of input data space*

# Kernel Methods

in kernel methods, the training dataset is also used in the **prediction phase**

The advantages of kernel methods arise when choice special, non-linear kernels

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$$

for example, assume a two-dimensional problems (**p=2**) with the following feature mapping

$$\phi(\mathbf{x}) = (x_1 x_2, x_1 - x_2)$$

then the corresponding kernel will be given by

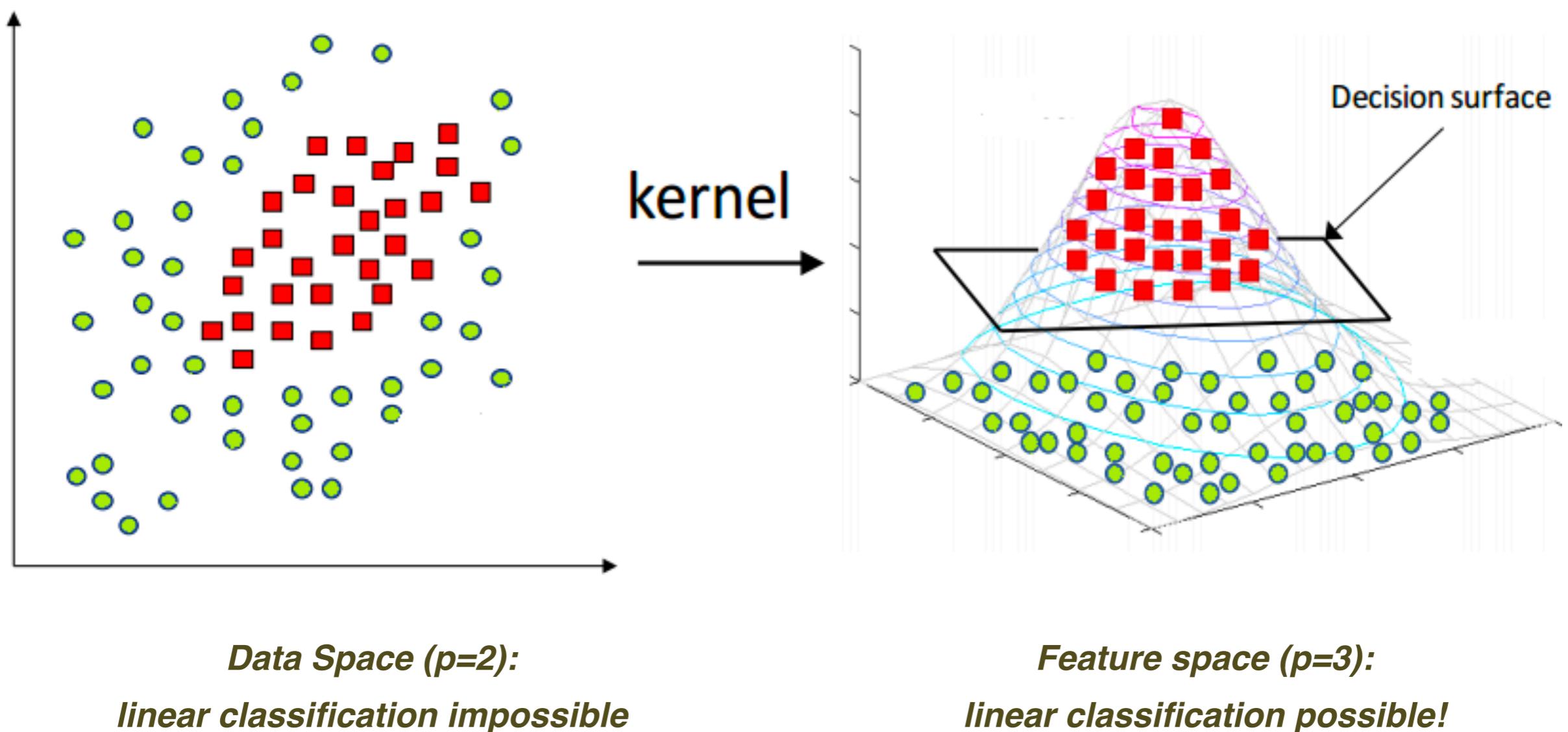
$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}') = x_1 x_2 x'_1 x'_2 + (x_1 - x_2)(x'_1 - x'_2)$$

as we will see, kernel methods exploit the **kernel trick**, whereby they can work directly at the level of kernels **without the need to specify the feature space mapping**

# Kernel Methods

Why working with kernels is advantageous?

Assume we work in a data space of **dimension p**, but classification requires a feature space of **dimension  $p' > p$**



# Kernel Methods

Why working with kernels is advantageous?

Assume we work in a data space of **dimension p=3**, but classification requires a feature space of **dimension p=9**

$$\phi(x) = (x_1^2, x_1x_2, x_1x_3, x_2x_1, x_2^2, x_2x_3, x_3x_1, x_3x_2, x_3^2)^T$$

$$\phi(x)^T \phi(y) = \sum_{i,j=1}^3 x_i x_j y_i y_j \quad \text{scales as } O(n^2)$$

# Kernel Methods

Why working with kernels is advantageous?

Assume we work in a data space of **dimension p=3**, but classification requires a feature space of **dimension p=9**

$$\phi(x) = (x_1^2, x_1x_2, x_1x_3, x_2x_1, x_2^2, x_2x_3, x_3x_1, x_3x_2, x_3^2)^T$$

$$\phi(x)^T \phi(y) = \sum_{i,j=1}^3 x_i x_j y_i y_j \quad \text{scales as } O(n^2)$$

The same result can be obtained working with an appropriate **kernel**

$$k(x, y) = (x^T y)^2 = (x_1 y_1 + x_2 y_2 + x_3 y_3)^2 = \phi(x)^T \phi(y)$$

*scales as O(n)*

**efficient and computationally advantageous** method to transform data to higher dimensions!

# Kernel Methods

Kernel methods are used in a **wide variety of Machine Learning** applications:

- ✿ Pattern analysis
- ✿ Clustering
- ✿ Dimensional reduction
- ✿ Classification
- ✿ Anomaly detection

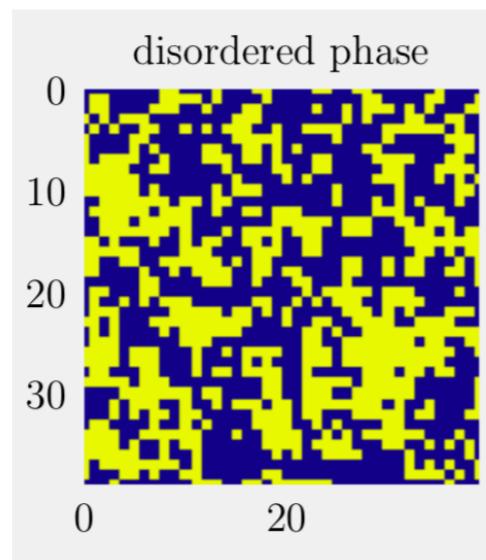
A kernel is a **similarity function**: it measures how similar two inputs are

working with **kernels** offers many advantages over working with **feature vectors**,  
which in general have infinite dimensionality

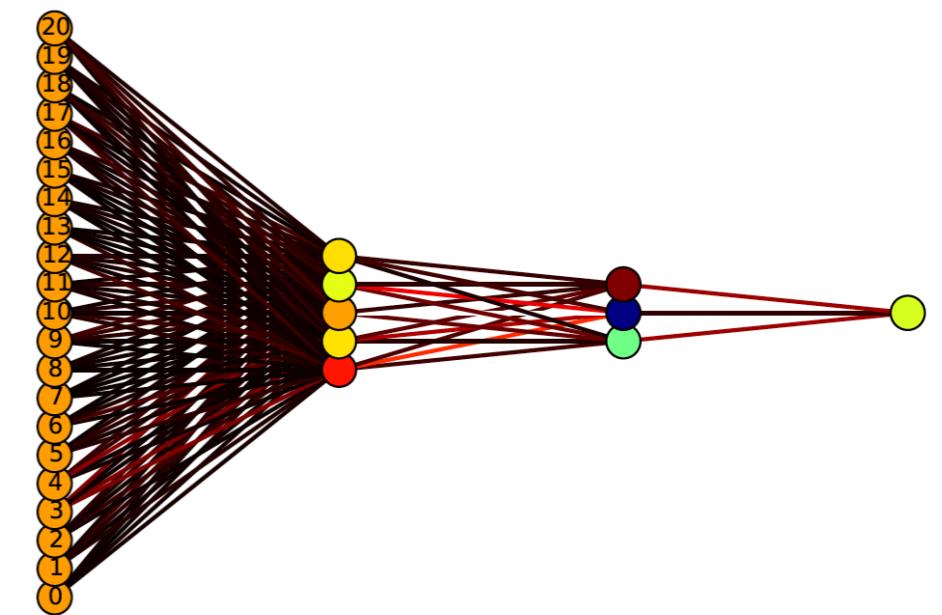
# Kernel Methods

Classifier working on **feature space**

labelled training examples => Feature map => input to classifier

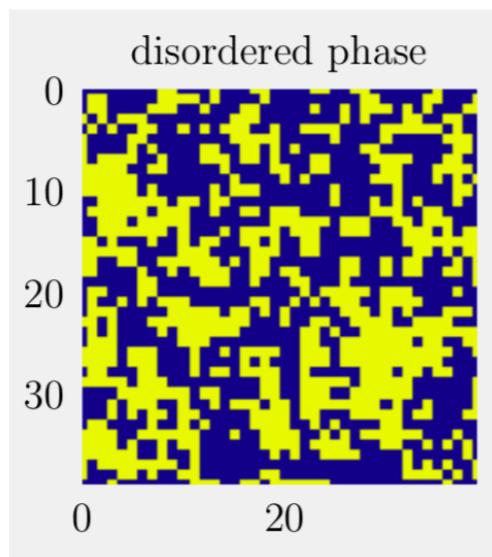


*array with spin values  
+ label indicating phase*



Classifier working with **kernel method**

labelled training examples + kernel (similarity measure) => Classifier



$$k(\mathbf{x}, \mathbf{x}', \alpha, \beta, \dots)$$

*optimise hyper-params of the kernel  
to achieve good classification*

*how to select a suitable kernel?*

# Support Vector Machines

# Support Vector Machines

SVM are **kernel methods** popular for classification, regression, and novelty detection

A peculiar feature of SVM is that the model parameters are found by a solution of a convex optimisation problem: any local solution is also a **global optimum**

The advantages of support vector machines include:

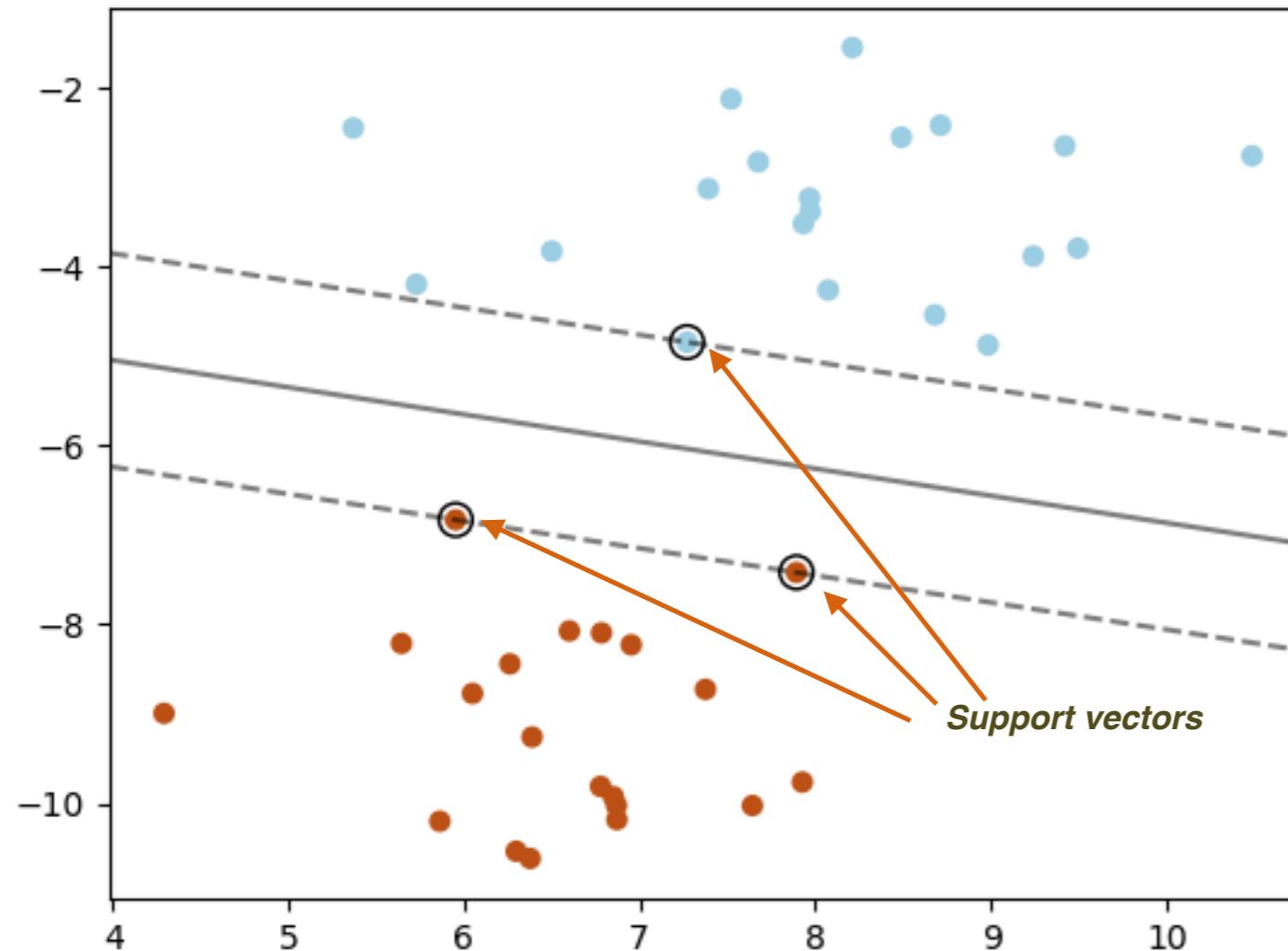
- Effective in **high dimensions**.
- Effective when number of dimensions is greater than number of samples.
- Includes on a **subset of training points** in the decision function (**support vectors**)
- Different **kernels functions** can be used for the figure of merit

*kernel methods operate in a high-dimensional,  
implicit feature space without computing the coordinates in data space*

# Support Vector Machines

SVM are ML algorithms popular for classification, regression, and novelty detection

A peculiar feature of SVM is that the model parameters are found by a solution of a convex optimisation problem: any local solution is also a **global optimum**



*In SVM classification, the support vectors are samples that lie in the margin boundaries (at least for linearly separable problems)*

# Support Vector Machines

As other **classification algorithms**, the problem that SVM try to solve is the following. Starting from a training dataset

$$\{\mathbf{x}_i^T, t_i\}, \quad i = 1, \dots, N, \quad t_i \in (-1, 1) \quad \text{binary classification}$$

our goal is to find the **model parameters** such that the prediction from

$$\text{sign}(\boldsymbol{\theta}^T \boldsymbol{\phi}(\mathbf{x}) + b) \quad \begin{array}{l} \text{linear in model parameters,} \\ \text{but non-linear in feature space} \end{array}$$

is correct for most of the samples. Here  $\boldsymbol{\phi}$  is the **feature-space transformation**, where the *feature space* is where our data lives

*Example of a feature-space transformation*

$\mathbf{x}$	$\boldsymbol{\phi}(\mathbf{x})$
<i>Height</i>	<i>Height/Weight</i>
<i>Weight</i>	<i>Weight+Height</i>
<i>Age</i>	<i>Age * Weight</i>

# Support Vector Machines

the problem then that SVM are aiming to solve is

$$\min_{\theta, b, \xi} \left( \frac{1}{2} \theta^T \theta + C \sum_{i=1}^n \xi_i \right)$$

*inverse of margin between the two categories*

subject to the **additional constraint** that

$$Class\ label\ of\ sample\ i \longrightarrow t_i \times (\theta^T \phi(x_i) + b) \geq 1 - \xi_i \longleftarrow tolerance\ (regulator)$$

which essentially aims to **maximise the margin between the two categories** while including a penalty when a sample is misclassified or within the margin boundary, where the term proportional to  $C$  plays the role of a **regulator**

SVM are an example of a **constrained optimisation problem**, where one tries to minimise a quadratic function subject to **linear constraints**. They are most efficiently solved by means of the Lagrange multiplier method

# Support Vector Machines

In which respect SVD is a **kernel method**? So far it seems to work only in **feature space** ...

we can reformulate the problem such that the feature space disappears completely!

many ML models for regression and classification can be reformulated in a  
**dual representation** where the kernel functions arise naturally

in doing so, we will also show in which respect kernel methods are special in that they  
**explicitly include the training examples** in the final decision model

# Dual representations

many ML models for regression and classification can be reformulated in a **dual representation** where the kernel functions arise naturally

consider a linear regression model with a regularised least-squares error function

$$E_{\text{tr}}(\boldsymbol{\theta}) = \frac{1}{2} \sum_{n=1}^N (\boldsymbol{\theta}^T \phi(\mathbf{x}_n) - t_n)^2 + \frac{\lambda}{2} \boldsymbol{\theta}^T \boldsymbol{\theta}$$

*inner product  
in feature space*

*sum over  
training samples*

*regulator*

*output from  
training examples*

*same structure  
as SVD!*

the model parameters are determined **analytically** by requiring the vanishing of the gradient

$$\boldsymbol{\theta} = \sum_{n=1}^N a_n \phi(\mathbf{x}_n), \quad a_n = -\frac{1}{\lambda} (\boldsymbol{\theta}^T \phi(\mathbf{x}_n) - t_n)$$

*recall that inner products take place in feature space*

one can formulate a dual representation of the problem in terms of the **parameter vector**

$$\mathbf{a} = (a_1, a_2, \dots, a_N)^T \quad \mathbf{t} = (t_1, t_2, \dots, t_N)^T$$

# Dual representations

with some algebra one can show that the original figure of merit of the model

$$E_{\text{tr}}(\boldsymbol{\theta}) = \frac{1}{2} \sum_{n=1}^N (\boldsymbol{\theta}^T \boldsymbol{\phi}(x_n) - t_n)^2 + \frac{\lambda}{2} \boldsymbol{\theta}^T \boldsymbol{\theta}$$

admits a **dual representation** in terms of **kernel function**

$$E_{\text{tr}}(\boldsymbol{a}) = \frac{1}{2} \boldsymbol{a}^T \mathbf{K} \boldsymbol{K} \boldsymbol{a} - \boldsymbol{a}^T \mathbf{K} \boldsymbol{t} + \frac{1}{2} \boldsymbol{t}^T \boldsymbol{t} + \frac{\lambda}{2} \boldsymbol{a}^T \mathbf{K} \boldsymbol{a}$$

where  $\mathbf{K}$  is known as the **Gram matrix**, an  $N \times N$  symmetric matrix with entries

$$K_{nm} = k(x_n, x_m) = \boldsymbol{\phi}(x_n)^T \boldsymbol{\phi}(x_m)$$

*matrix in space of  
input data*

given by the **kernel function** evaluated over two of the **input training examples**

the crucial property of this dual representation is that now the predictions for new inputs will explicitly **depend on the training examples**

*model params nor  
feature space maps  
do not appear  
in this formulation*

*SVD can be expressed entirely in terms of kernels*

# Dual representations

the crucial property of this dual representation is that now the predictions for new inputs will explicitly **depend on the training examples**

$$y(\mathbf{x}) = \mathbf{k}(\mathbf{x})^T (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t}$$

*output of trained model*      *new input*      *depends on the N original inputs*

$$\mathbf{k}(\mathbf{x}) = (k(\mathbf{x}_1, \mathbf{x}), k(\mathbf{x}_2, \mathbf{x}), \dots)$$

in the dual formulation we invert a  $N \times N$  matrix (data space) rather than a  $M \times M$  one (feature space), so this does not appear to be advantageous ...

the main benefit is being able to **work directly with kernel functions** and bypass a choice of feature map

this allows one to work in feature spaces of **very high (even infinite) dimensionality**

# Support Vector Machines



<https://www.youtube.com/watch?v=5zRmhOUjjGY>

# The kernel trick

recall that the kernel function is constructed from a **feature space mapping**

$$k(\mathbf{x}_n, \mathbf{x}_m) = \boldsymbol{\phi}(\mathbf{x}_n)^T \boldsymbol{\phi}(\mathbf{x}_m) = \sum_{i=1}^M \phi_i(\mathbf{x}_n)\phi_i(\mathbf{x}_m)$$

one can also construct kernel functions directly, provided they can be expressed as above

e.g., for a 2D input space  $k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z})^2 = (x_1 z_1 + x_2 z_2)^2 = \boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\phi}(\mathbf{z})$

$$\boldsymbol{\phi}(\mathbf{x}) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$$

note that feature space dimension ( $M=3$ ) is bigger than input space dimensionality ( $d=2$ )

there exist more general methods to construct acceptable kernels e.g. using **basis functions**

a popular choice is Gaussian kernel, built from an **infinite-dimensional feature map**

$$k(\mathbf{x}, \mathbf{z}) = \exp\left(-||\mathbf{x} - \mathbf{z}||^2/2\sigma^2\right)$$

# The kernel trick

a popular choice is Gaussian kernel, built from an **infinite-dimensional feature map**

$$k(x, z) = \exp\left(-||x - z||^2/2\sigma^2\right)$$

we can demonstrate that this is a **valid kernel** as follows: given two valid kernels

$$k_1(x, z), \quad k_2(x, z)$$

we can construct **new valid kernels** from them using the following rules

$$k(x, z) = f(x)k_1(x, z)f(z) \quad k(x, z) = \exp(k(x, z))$$

so now we can check that the Gaussian kernel is indeed a *bona-fide* kernel by expanding

$$k(x, z) = \exp\left(-x^T x/2\sigma^2\right) \exp\left(-x^T z/2\sigma^2\right) \exp\left(-z^T z/2\sigma^2\right)$$

*linear kernel*

*the feature map associated to this kernel has infinite dimensionality!*

# Summary

- ➊ Many problems in ML require assessing the **similarly between probability distributions**
- ➋ **Information theory** provides several quantities (e.g. KL divergence) to quantify this
- ➌ Generative models are design to evaluate **new instances of the training data**, with features as close as possible to the real ones
- ➍ Generative Adversarial Networks are a particular implementation of a generative model where **two networks compete againts each other** to evaluate the generative probability