

# DIFFUSION MODELS

Machine learning  
reading group,  
Wag.

omega sabre

Used in generative AI.

So far we have seen SUPERVISED LEARNING.

- .) Data is of type ( $x = \text{input}$ ,  $y = \text{output}$ ), training set is of type  $\{(x_i; y_i), i=1, \dots, N\}$ .
- .) The result of training is a function  $f$  s.t. for an unseen input  $x$ , the appropriate output is  $y = f(x)$ .
- .) Examples:
  - 1.) Digit recognition: input is an image [5], output is a single digit 5. (Pietro & Marc W7, Poppy W8.)
  - 2.) Property prices: input is a triplet (size, # of rooms, location), output is the value (£).
- .) All algorithms so far have been DETERMINISTIC: Training on two separate occasions using the same training set yields the same function  $f$ . Feeding the same ~~input~~ into  $f$  always gives the same output.

NOTE: Probability measures were used only for averaging, not for sampling; deterministic.

## This week: UNSUPERVISED LEARNING.

- .) Data is not of type  $(x, y)$ .
- .) The result of training is some sort of a generator, that does not necessarily take any input.
- .) Example: generating images of cats.
- .) The algorithm will be STOCHASTIC (non-deterministic)  
 $\rightarrow$  ~~training~~ the generator <sup>will</sup> might generate two distinct images if run twice, even if all other parameters are the same.  
NOTE: the training procedure is still deterministic! training on the same training set always gives the same generator.

Notational differences with the source:

$$\mathbf{x}, \theta \rightsquigarrow \Psi, \mathbf{r}$$

## The general setup

- Assume images of cats are elements of  $\mathbb{R}^d$  for some  $d$ . ~~For example,~~

Example: resolution 800x600, each pixel ~~bit~~ represented by a 32-bit number  $\rightsquigarrow$  simplify to an arbitrary real number.  $\Rightarrow d = 800 \cdot 600$ .

- Assume  $\exists$  a probability density  $\Psi$  on  $\mathbb{R}^d$  s.t. for  $u \in \mathbb{R}^d$ :

$$\Psi(u) \propto P(u \text{ is an image of a cat}).$$

We do not know  $\Psi$ ; we will learn ~~it~~ it from sample data.

- Assuming we know  $\Psi$ , we generate new images as follows (see figure below):

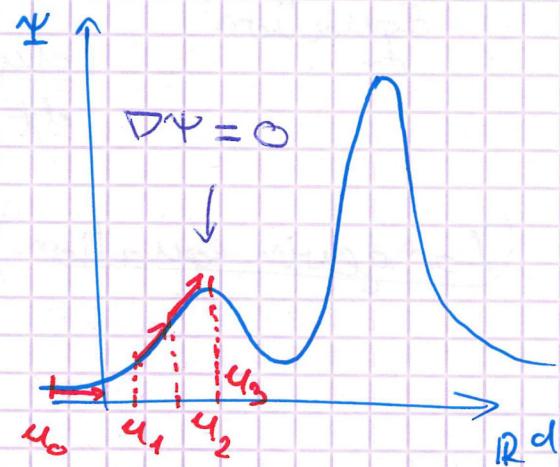
- 1.) Take a random image  $u_0 \in \mathbb{R}^d$ .

This will not be an image of a cat (most likely).

- 2.) Update (transform) the image "in the direction" of  $\nabla \Psi$  to get a new image  $u_1 \in \mathbb{R}^d$ .

This image will ~~have~~ have a higher probability of looking like a cat.

Repeat step 2.



This can be put into an equation:

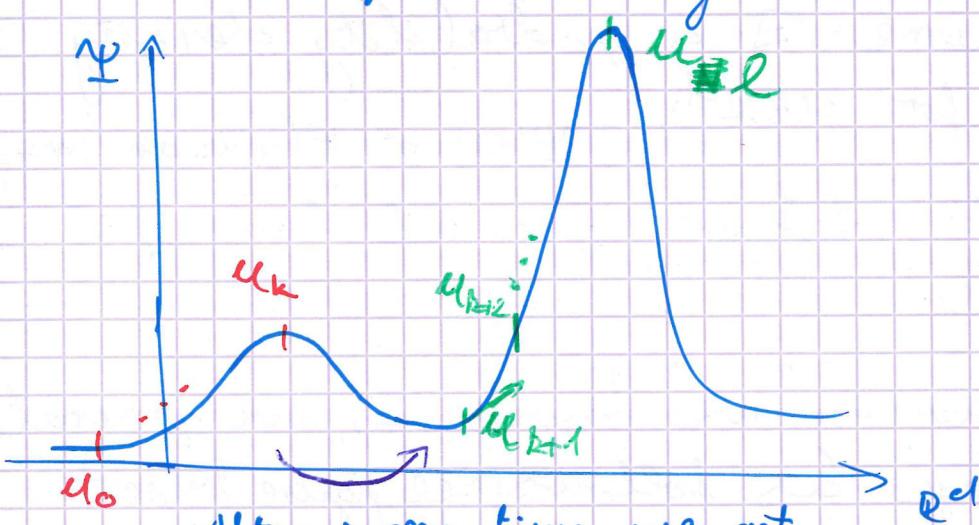
$$\frac{du}{dt} = \nabla \Psi(u) \quad (u(0) = u_0)$$

(Discretising in time gives the steps 1 & 2.)

ISSUE: We get permanently stuck at the first stationary point we encounter.

SOLUTION: Add some noise; this will (eventually) perturb us out of stationary points:

See



After some time we get  
perturbed enough to start moving  
towards a more optimal peak.

Equation:

$$\frac{du}{dt} = \nabla \Psi(u) + \sqrt{2} \frac{dW}{dt}$$

(Langevin equation)

↑  
Brownian/Wiener noise  
(standard).

## ONE MORE ADJUSTMENT: Replace $\nabla \Psi$ with $\nabla \log \Psi$ .

This is called a "score".

- Benefits: 1.) In practice,  $\Psi$  is often a product of factors. Taking log produces sums, which are easier to work with.
- 2.) Computationally more stable in regions where  $\Psi \approx 0$ .

See training set: a sequence of i.i.d. samples.

Get

$$\frac{d\mu}{dt} = \nabla \log \Psi(u)$$

$$\frac{d\mu}{dt} = \nabla \log \Psi(u) + \sqrt{2} \frac{d\Psi}{dt}.$$

To run the generating dynamics, we need to learn  $\nabla \log \Psi$ .

GOAL: Find  $s_\theta: \mathbb{R}^d \rightarrow \mathbb{R}^d$  from a family of functions  $\Omega = \{s_\theta: \mathbb{R}^d \rightarrow \mathbb{R}^d, \theta \in \mathbb{R}^n\}$  st.  $s_\theta \approx \nabla \log \Psi$  in an "appropriate" sense.

One option (not pursued here): Assume  $\exists \Psi_\theta$  prob. density on  $\mathbb{R}^d$  st.  $s_\theta = \nabla \log \Psi_\theta$ .

Find  $\Psi_\theta$  that maximises the likelihood of the samples in the training set (see page 6)  
This is MAXIMUM LIKELIHOOD ESTIMATION!

INSTEAD, we will try to minimise (over  $\theta$ )

$$J(\theta; \Psi) := \int_{\mathbb{R}^d} |s_\theta(u) - \nabla \log \Psi(u)|^2 \Psi(u) du.$$

see note on page 6.

## MINIMISATION

We are given  $N \gg 1$  images of rats:  $\{u_i\}_{i=1}^N$

We treat them as i.i.d. random variables distributed  $\sim \Psi$ .

This cannot be used to estimate  $\nabla \log \Psi$ .

We do not know  $\Psi(u)$ , and we cannot do any serious attempt of ~~to~~ differentiating.

↳ independent identically distributed

NOTE:  $\operatorname{div} \boldsymbol{\sigma}_\Psi = \operatorname{Tr} D \boldsymbol{\sigma}_\Psi$ ,  
 $\nearrow$  D ... Jacobian

SOLUTION: Def-

$$K(\boldsymbol{v}; \Psi) = \int_{\mathbb{R}^d} (|\boldsymbol{\sigma}_{\boldsymbol{v}}(u)|^2 + 2 \operatorname{div} \boldsymbol{\sigma}_{\boldsymbol{v}}(u)) \Psi(u) du.$$

PROP.: Minimising  $J(\boldsymbol{v}; \Psi)$  over  $\boldsymbol{v}$  is equivalent to minimising  $K(\boldsymbol{v}; \Psi)$  over  $\boldsymbol{v}$ .

NOTE (forgot to write on p. 5):

$$J(\boldsymbol{v}; \Psi) = \mathbb{E}_{u \sim \Psi} [|\boldsymbol{\sigma}_{\boldsymbol{v}}(u) - \nabla \log \Psi(u)|^2]$$

This gives probabilistic interpretation to  $J$ .

PROOF of Prop.:

$$\begin{aligned} |\boldsymbol{\sigma}_{\boldsymbol{v}} - \nabla \log \Psi|^2 &= \langle \boldsymbol{\sigma}_{\boldsymbol{v}} - \nabla \log \Psi, \boldsymbol{\sigma}_{\boldsymbol{v}} - \nabla \log \Psi \rangle \\ &= |\boldsymbol{\sigma}_{\boldsymbol{v}}|^2 - 2 \langle \boldsymbol{\sigma}_{\boldsymbol{v}}, \nabla \log \Psi \rangle + |\nabla \log \Psi|^2 \end{aligned}$$

$$\text{Observe: } \langle \boldsymbol{\sigma}_{\boldsymbol{v}}, \nabla \log \Psi \rangle = \langle \boldsymbol{\sigma}_{\boldsymbol{v}}, \frac{\nabla \Psi}{\Psi} \rangle$$

$$\Rightarrow \int_{\mathbb{R}^d} \langle \boldsymbol{\sigma}_{\boldsymbol{v}}, \nabla \log \Psi \rangle \Psi du = \int_{\mathbb{R}^d} \langle \boldsymbol{\sigma}_{\boldsymbol{v}}, \frac{\nabla \Psi}{\Psi} \rangle \Psi du$$

Detour to vector calculus: (Higher-dimensional integration by parts.)

Let  $F: \mathbb{R}^d \rightarrow \mathbb{R}^d$  be a vector field,  
 $g: \mathbb{R}^d \rightarrow \mathbb{R}$  scalar field.

$$\begin{aligned}\operatorname{div}(g \cdot F) &= \sum_{i=1}^d \frac{\partial}{\partial x_i} (g F_i) \\ &= \sum_{i=1}^d \left( \frac{\partial}{\partial x_i} g \right) \cdot F_i + \sum_{i=1}^d g \cdot \frac{\partial}{\partial x_i} F_i \\ &= \langle \nabla g, F \rangle + g \operatorname{div} F\end{aligned}$$

~~$\Rightarrow \int \langle \nabla g, F \rangle dx = \int g \operatorname{div} F dx$~~

Substitute  $g = \Psi$ ,  $F = \nabla \varphi$ .

$$\Rightarrow \int_{\mathbb{R}^d} \langle \nabla \varphi, \nabla \Psi \rangle du = \int_{\mathbb{R}^d} \operatorname{div}(\Psi \cdot \nabla \varphi) du - \int_{\mathbb{R}^d} \Psi \cdot \operatorname{div} \nabla \varphi du$$

By Green's/Divergence thm:

$$\int_{\mathbb{R}^d} \operatorname{div}(\Psi \cdot \nabla \varphi) du = \int_{\partial \mathbb{R}^d} \Psi \langle \nabla \varphi, \hat{n} \rangle dS \stackrel{(*)}{=} 0$$

Improper integral  
surface

$$(*) \quad \lim_{|u| \rightarrow \infty} \Psi(u) = 0 \quad (\text{assume})$$

Since  $\nabla \varphi$  should be  $\approx \nabla \log \Psi$ ,  
we can conclude  $\Psi \nabla \varphi(u) \rightarrow 0$  as  $|u| \rightarrow \infty$ .

$$\Rightarrow \int_{\mathbb{R}^d} \langle \nabla \varphi, \nabla \Psi \rangle du = - \int_{\mathbb{R}^d} \Psi \cdot \operatorname{div} \nabla \varphi du$$

Independent of  $\varphi$ , so can be ignored when minimizing. Assuming

$$J(u; \Psi) = \int_{\mathbb{R}^d} (|\nabla \varphi(u)|^2 + 2 \operatorname{div}(\nabla \varphi)) \Psi(u) du + \int_{\mathbb{R}^d} |\nabla \log \Psi(u)|^2 du$$

□

17

This is useful since  $K$  can be estimated from  $\{u_i\}_{i=1}^N$ :

$$K(r; \Psi) \approx \frac{1}{N} \sum_{i=1}^N |\nabla_\Psi(u_i)|^2 + 2 \operatorname{div}(\nabla_\Psi u_i) \xrightarrow[N \rightarrow \infty]{a.s.} K(r; \Psi)$$

a.s. means "almost surely",  
i.e. with probability 1.

Follows from the strong law of large numbers under some assumptions.

We now consider another approach.

Def. (for  $\sigma > 0$ ):

$$p_\sigma(u) := \frac{1}{(2\pi\sigma)^{d/2}} e^{-\frac{1}{2\sigma^2} \|u\|^2}$$

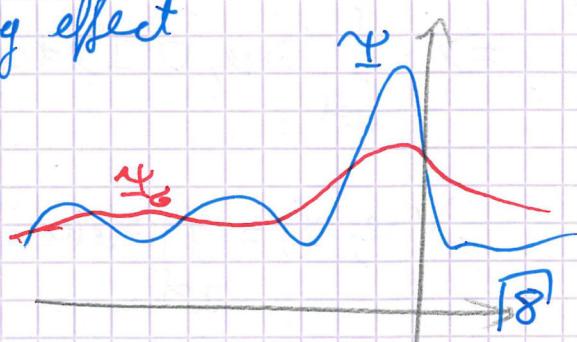
Multivariate  
(d-dim) Gaussian  
prob. density.

$$\begin{aligned} \Psi_\sigma(u) &:= (p_\sigma * \Psi)(u) = \int_{\mathbb{R}^d} p_\sigma(u-w) \Psi(w) dw \\ &= \int_{\mathbb{R}^d} p_\sigma(u-w) \Psi(w) dw. \end{aligned}$$

Interpretation of  $\Psi_\sigma$ : the value of  $\Psi_\sigma(u)$  is a "weighted average" of  $\Psi$ , where the weight is given by  $p_\sigma(w)$ . The averaging effect increases with  $\sigma$ .

$\Psi_\sigma \xrightarrow{\sigma \rightarrow \infty}$  "mostly constant function"  
(TAKE WITH 3 PINCHES OF SALT)

$\Psi_\sigma \xrightarrow{\sigma \rightarrow 0} \Psi$



We will minimize  $J(\mathbf{w}; \Psi_0)$  instead of  $J(\mathbf{v}; \Psi)$ . In practice, this is done for a range of  $\delta$ 's.

### MOTIVATION:

1.) As mentioned earlier it is challenging to properly learn  $\log \Psi$  in the regions where  $\Psi \approx 0$ . This is because all our samples (training set) are from areas where  $\Psi \gg 1$ .

This poses a problem: when starting the generative algorithm from a random point (where  $\Psi \approx 0$ ), we do not know well how to update the image. Our knowledge of  $\nabla \log \Psi$  is ~~bad~~ poor in this region of  $\mathbb{R}^d$ . We thus learn  $\nabla \log \Psi_\delta$  for  $\delta$  large first, and progressively decrease  $\delta$ .

2.) Minimising  $J(\mathbf{v}; \Psi_\delta)$  is easier than minimising  $L(\mathbf{v}; \Psi)$ . Problem with  $L(\mathbf{v}; \Psi)$ :  $\text{div}(\nabla \mathbf{v})$  is costly to compute.

We use the following proposition.

PROP.: Minimising  $J(\mathbf{v}; \Psi_\delta)$  over  $\mathbf{v}$  is  $\Leftrightarrow$  to minimising  $L(\mathbf{v}; \Psi_\delta)$  over  $\mathbf{v}$ , where

$$L(\mathbf{v}; \Psi_\delta) := \iint_{\mathbb{R}^d \times \mathbb{R}^d} |S_{\Psi_\delta}(v) - \nabla v \log p_\delta(w|v)|^2 p_\delta(w|v) \Psi(w) dw dv.$$

PROOF: See the book (Prop. 13.13, p. 181 in the my version). GibHub/drive

Again,  $L(\mathbf{v}; \Psi_\delta)$  is easily estimated from the training set:

$$L(\mathbf{v}; \Psi_\delta) \approx \frac{1}{N} \sum_{i=1}^N \iint_{\mathbb{R}^d} |S_{\Psi_\delta}(v_i) - \nabla v \log p_\delta(w|v_i)|^2 p_\delta(w|v_i) \Psi(w) dw.$$