

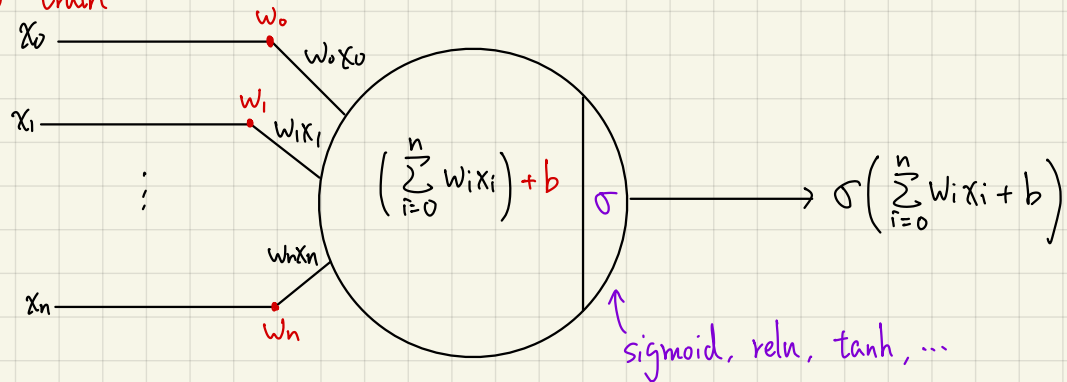
# Training

The most enigmatic procedure in machine learning is training of neural networks, or, in general, parametric families of functions. This is of course an inverse problem, on which we have developed two perspectives above: the optimization perspective (with regularization) and the Bayesian perspective.

Essentially training is described as minimization of loss, i.e. for a given loss function  $L$  on a space of functions  $f$  parametrized by a set of parameters  $\theta \in \Theta$

$$\operatorname{arginf}_{\theta \in \Theta} L(f^\theta)$$

$w, b$  : need to train



$f$ : our neural network

$\Theta$ : weights, biases

$L$ : loss function

Assume that  $\Theta$  is some open subset of points in  $\mathbb{R}^d$  and  $U : \Theta \rightarrow \mathbb{R}$ ,  $\theta \mapsto L(f^\theta)$  a sufficiently regular function with a unique minimum  $\theta^* \in \Theta$ , then one can describe essentially one local and one global method to find the infimum:

1. If  $U$  is strictly convex and  $C^2$  in a neighborhood of the unique minimizer  $\theta^*$ , then

$$d\theta_t = \frac{-\nabla_{\Theta} U(\theta_t) dt}{\text{towards minimum}}$$

converges to  $\theta^*$  as  $t \rightarrow \infty$ . For any  $t \geq 0$  it holds that

$$dU(\theta_t) = -\|\nabla U(\theta_t)\|^2 dt,$$

i.e. the value of  $U$  is strictly increasing along the path  $t \mapsto \theta_t$ . Together with the fact that  $U$  is strictly convex we obtain a convergence of  $\|\theta_t - \theta^*\| \leq C \exp(-\lambda_{\min} t)$  as  $t \rightarrow \infty$ , where  $\lambda_{\min}$  is the minimum of the smallest eigenvalue of the Hessian of  $U$  on  $\Theta$ . This holds remarkably for any starting point  $\theta_0 \in \Theta$  and is the basis of all sorts of gradient descent algorithms.

$$d\theta_t \rightsquigarrow \theta_{t+1} - \theta_t \quad * \theta_t \text{ is a vector, not function}$$

$$\Rightarrow \theta_{t+1} - \theta_t = -\int \nabla U(\theta_t) dt \quad (\text{here, set } \delta=1) = -\delta_t \nabla U(\theta_t) \quad (\text{check the proof of G.D.})$$

Since  $U \in C^2$  and (strictly) convex, the gradient descent algorithm converges to  $\theta^*$ .

Rmk. of  $U$

- The convexity  $\checkmark$  is necessary in the proof the convergence of G-D.
- gradient  $\rightarrow$  (global/local) maximum
- gradient  $\rightarrow$  (global/local) minimum
- $\delta_t$  determines the length  $\nabla U(\theta_t)$ , also called the "learning rate"

2. A far reaching generalization is given by the following consideration: consider  $U$  on  $\Theta$  having a unique minimizer  $\theta^* \in \Theta$ , then the probability measure given by the density with respect to Lebesgue measure on  $\mathbb{R}^d$

$$p_\epsilon := \frac{1}{Z_\epsilon} \exp \left( -\frac{U}{\epsilon} \right)$$

tends in law to  $\delta_{\theta^*}$  as  $\epsilon \rightarrow 0$ . The denominator  $Z_\epsilon$  is just the integral  $\int_\Theta \exp(-U(\theta)/\epsilon) d\lambda(\theta)$  and the above statement nothing else than the fact that the described density function concentrates at  $\theta^*$ . If one manages to sample from the measure  $p_\epsilon d\lambda$ , then one can approximate empirically  $\theta^*$ .

The measure  $p_\epsilon d\lambda$  is the invariant measure of the stochastic differential equation

$$d\theta_t = -\frac{1}{2} \nabla U(\theta_t) dt + \sqrt{\epsilon} dW_t,$$

which is just checked by the following equality

$$\int_\Theta \left( -\frac{1}{2} \nabla U(\theta) \nabla f(\theta) + \frac{\epsilon}{2} \Delta f(\theta) \right) p_\epsilon(\theta) d\lambda(\theta) = 0$$

for all test functions  $f$ .

One method, which generalizes this thought in a time-dependent way, is to sample from a measure concentrating at  $\theta^*$  is to simulate from a stochastic differential equation (for  $\Theta = \mathbb{R}^N$ ) of the type

$$d\theta_t = -\nabla U(\theta_t) dt + \alpha(t) dW_t$$

where  $W$  is an  $N$ -dimensional Brownian motion and the non-negative quantity, called cooling schedule,  $\alpha(t) = O(\frac{1}{\log(t)})$  as  $t \rightarrow \infty$ . For appropriate constants we obtain that  $\theta_t$  converges in law to  $\delta_{\theta^*}$  as  $t \rightarrow \infty$ . This procedure is called simulated annealing and is the fundament for global minimization algorithms of the type 'steepest descent plus noise'.

$U$ : unique minimizer  $\theta^* \in \Theta$

Consider a measure  $p_\varepsilon = \frac{1}{Z_\varepsilon} \exp(-\frac{U}{\varepsilon})$ ,  $Z_\varepsilon = \int_{\Theta} \exp(-\frac{U(\theta)}{\varepsilon}) d\lambda(\theta)$

$\leadsto$  ① concentrate at  $\theta^*$  as  $\varepsilon \rightarrow 0$ .

② if one manages to sample from the measure  $p_\varepsilon d\lambda$ ,  
then one can approximate  $\theta^*$ .

before:  $d\theta_t = -\nabla U(\theta_t) dt$

now:  $d\theta_t = -\nabla U(\theta_t) dt + \alpha(t) dW_t$ ,  $W_t$ : non-negative Brownian motion  
 $\alpha(t)$ : volatility

\* variance of Brownian motion: s-t if  $W_t$  to  $W_s$

\* volatility affects the variance

In [4]:

```
import numpy as np
```

```
# Trajectories of  $d \theta = - \nabla U(\theta) dt + \alpha dW$ 
```

```
N=2000 # time discretization
```

```
theta0=1 # initial value of theta
```

```
T=100 # maturity
```

```
alpha = 0.1 # volatility in Black Scholes
```

```
R=10**5 # number of Trajectories
```

```
def U(theta):
```

```
    return theta**2*(theta-2)**2
```

```
    #return theta**2*(theta-2)**2*np.sin(theta)**2
```

```
def gradU(theta):
```

```
    return 2*theta*(theta-2)**2 + theta**2*2*(theta-2)
```

```
    #return (2*theta*(theta-2)**2 + theta**2*2*(theta-2))*np.sin(theta)**2+2*np.sin(theta)*np.cos(theta)*th
```

```
theta = np.zeros((N,R))
```

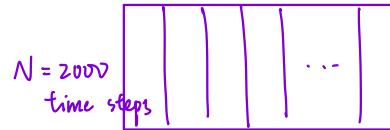
```
theta[0,:]=theta0*np.ones((1,R)) → initial value of each trajectory is 1
```

```
for j in range(N-1):
```

```
    increment = np.random.normal(0, np.sqrt(alpha)*np.sqrt(T)/np.sqrt(N), (1,R))
```

```
    theta[j+1,:] = theta[j,:] + increment - 0.5*gradU(theta[j,:])*(T/N)
```

$R=10000 = \# \text{ trajectories}$

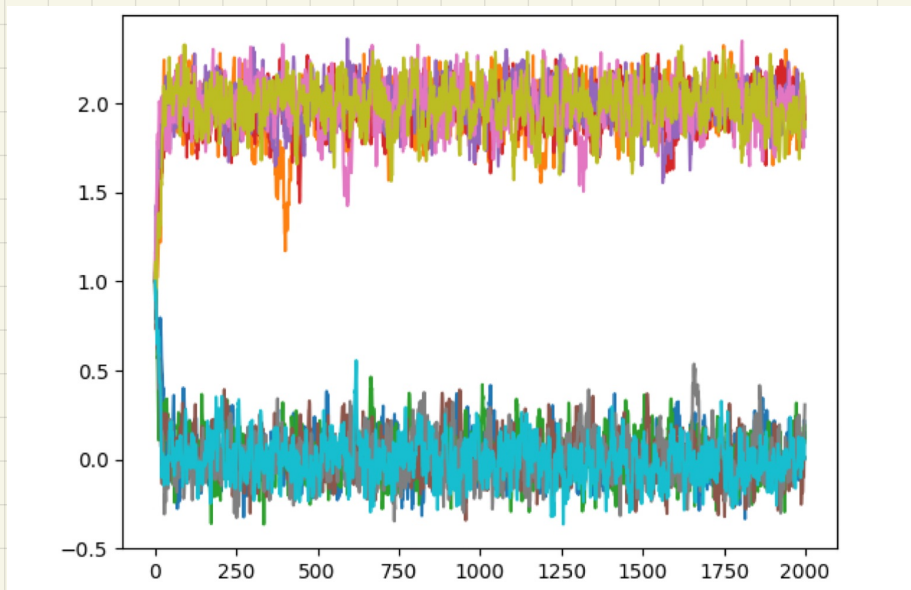


std. deviation

output shape (row vector)

$[N(\mu, \sigma^2), N(\mu, \sigma^2), \dots, N(\mu, \sigma^2)]$

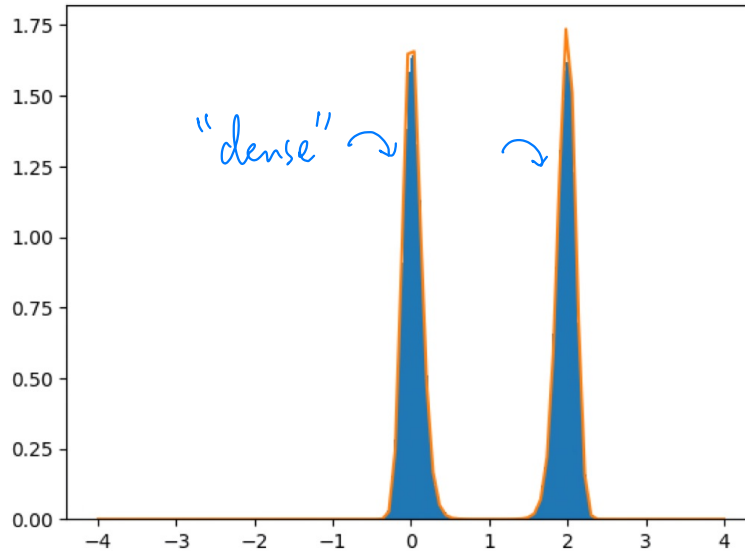
result →  
(前10條 trajectories)



In [6]:

```
a=4
x=np.linspace(-a,a,100)
plt.hist(theta[N-1,:],density=True,bins=100)
plt.plot(x,np.exp(-U(x)/alpha)/(2*a/100*sum(np.exp(-U(x)/alpha))))
plt.show()
```

$$P_i = \frac{\exp(-\frac{U(\theta_i)}{\alpha})}{\sum \exp(-\frac{U(\theta)}{\alpha})}$$



The stochastic gradient descent algorithm essentially says for a function of expectation type  $U(\theta) = E[V(\theta)]$

$$\theta_{n+1} = \theta_n - \gamma_n \nabla V(\theta_n, \omega_n)$$

for independently chosen samples  $\omega_n$  converges in law to  $\theta^*$ . Notice first that all our examples, in particular the ones from mathematical finance, are of expectation type, where the samples  $\omega$  are usually seen as elements from the training data set.

GD

```
for i in range(# epochs):  
     $\theta = \theta - \eta \text{grad}(L, \text{data}, \theta)$ 
```

e.g.  $L = \sum_i \left( \hat{y}^i - \left( b + \sum_j w_j x_j \right) \right)^2$  (MSE)  
          summation over all training examples

$$U(\theta) = E[V(\theta)]$$

rewrite

$$V(\theta) = \mathbb{E}_{n \sim \text{Unif}}[V_n(\theta)] = \mathbb{E}_{n \sim \text{Unif}}[U(\theta, \omega_n)]$$

Let  $N = \text{data size} = \# \text{ training examples.}$

SGD

```
for i in range(# epochs):  
    np.random.shuffle(data)  
    for example in data:
```

or -  $\theta = \theta - \eta \text{grad}(L, \text{example}, \theta)$   
       $\text{example} = \text{np.random.uniform}(1, \text{len}(\text{data}) + 1)$   
       $\theta = \theta - \eta \text{grad}(L, \text{Wexample}, \theta)$

for  $n$  in range  $N$ :

$n \sim \text{Unif}[1, N]$ ,  $\omega_n$ : sample

$$\begin{aligned} \theta &= \theta - \eta \nabla V_n(\theta) \\ &= \theta - \eta \nabla U(\theta, \omega_n) \end{aligned}$$

Lemma  $\mathbb{E}[\nabla U(\omega_n, \theta)] = \nabla U(\theta)$   $\theta \rightarrow \theta^*$



By the central limit theorem and appropriate sub-sampling one can understand

$$\nabla V(\theta, \omega) = \nabla U(\theta) + \text{'Gaussian noise with a certain covariance structure'} \Sigma(\theta)$$

where  $\Sigma(\theta)$  is essentially given by

$$\text{cov}(\nabla V(\theta)) .$$

Please check the paper: **On the Noisy Gradient Descent that Generalizes as SGD**