

Computational intelligence: set of neuro-inspired methodologies to address complex real-world problems to which traditional methodologies are unclear

- combines 3 frameworks:
  - ANN
  - EA
  - Fuzzy Logic → understand natural language → expensive



Deep learning: part of ML methods, based on ANN with representation learning  
- both supervised or non-supervised

Differences with brain:
 

- 1) dynamic - plastic
- 2) analog

 → Not just different # of neurons (not so important)

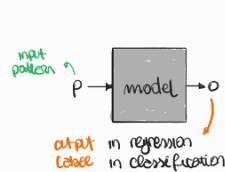
Turing Test → used to understand if a machine shows human-level intelligence

- judge interacts with two boxes and needs to decide which one is a machine
- doesn't mean that it has consciousness but that it can have a human conversation by mimicking
- Chat GPT-4o has this abilities

AGI (Artificial Capable Intelligence) → whether the AI is able to accomplish something concrete (like doubling an initial sum of money)

AFI is the next step but it doesn't have an associated step

LEAST SQUARES



Why not using 1-NN instead of a ANN? <sup>with a look-up table</sup>

- it cannot interpolate and generalize
  - it cannot perform denoising
- ↳ k-NN would mitigate the problem interpolating

But
 

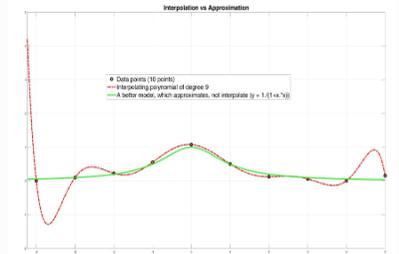
- 1) affords all complexity of inference phase
- 2) since it is model free we cannot compress or exploit regularities in D

Training set  $\{P_i, y_i\}_{i=1}^{N_{obs}}$  where  $P_i \in \mathbb{R}^D$   
↳ exact in classification or output in regression

TRS ∪ TSS = learning set while TRS ∩ TSS =  $\emptyset$

$$P = \begin{bmatrix} P_1^T & \dots & P_N^T \\ \vdots & & \vdots \\ P_1^T & \dots & P_N^T \end{bmatrix} \in \mathbb{R}^{N \times D} \quad \vec{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} \in \mathbb{R}^N$$

typically  $N \gg D$

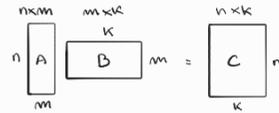


Interpolation: estimate  $f$  via  $D = \langle P, \vec{y} \rangle$  exactly

Approximation: estimate  $f$  via  $\langle P, \vec{y} \rangle$  approximately, while sacrificing some level of accuracy

LLS (Linear Least Squares)

$f(p) = b + vp$       $b \in \mathbb{R}$       $v \in \mathbb{R}^{1 \times D} \rightarrow$  row vector



We want to minimize  $E = \frac{1}{2} \|f(p) - y\|^2$  → minimizing this over the dataset is called LLS

Define  $X = \begin{bmatrix} 1 & \dots & 1 \\ P_1^T & \dots & P_N^T \\ \vdots & & \vdots \\ P_1^T & \dots & P_N^T \end{bmatrix} \in \mathbb{R}^{(N+1) \times N}$       $y = \vec{y}^T \in \mathbb{R}^{1 \times (N+1)}$

$\Rightarrow w = [b, v] \in \mathbb{R}^{1 \times (D+1)}$  can be obtained by  $\vec{w} = y X^T (X X^T)^{-1} \in \mathbb{R}^{1 \times (D+1)}$      time complexity  $O(N(D+1)^2)$

↳ **Penrose Pseudo-Inverse**

- $X X^T : [D+1] \times [N] \cdot [N] \times [D+1] \rightarrow \mathcal{O}(N(D+1)^2)$
- inverting  $\rightarrow \mathcal{O}((D+1)^3)$
- $X^T A : [N] \times [D+1] \cdot [D+1] \times [D+1] \rightarrow \mathcal{O}(N(D+1)^2)$

$L_w = \frac{1}{N} \sum (x_i w - y_i)^2 = \|wX - y\|_2^2 = (wX - y)(wX - y)^T = (wX - y) \begin{pmatrix} x_i^T w - y_i \\ \vdots \\ x_N^T w - y_N \end{pmatrix} = wX X^T w^T - \underbrace{wX y^T - y X^T w^T}_{\text{SCALARS}} + y y^T = wX X^T w^T - 2wX y^T + y y^T$

$\nabla_w L_w = \frac{\partial}{\partial w} [wX X^T w^T] - \frac{\partial}{\partial w} [2wX y^T] = 2wX X^T - 2(y X^T)^T = 0$       $\Rightarrow w = y X^T (X X^T)^{-1}$

↳ continuous + unconstrained

$$\frac{\partial}{\partial x} x^T A x = \frac{\partial}{\partial x} \sum_i \sum_j m_{ij} x_i x_j = \sum_j m_{1j} x_j + \sum_i m_{i1} x_i = \sum_k (m_{1k} + m_{k1}) x_k$$

$\rightarrow \frac{\partial}{\partial x} x^T A x = (A + A^T) x$       $\frac{\partial}{\partial x} x^T A x = 2Ax$       $\downarrow$       $k$ -th component of  $(A + A^T)x$

WLLS If we want to weight each  $(x_i, y_i)$  pair using  $\delta_i$ :  $w = y \Delta X^T (X \Delta X^T)^{-1}$       $\Delta = \text{diag}(\delta)$   
• weight less outliers

GRADIENT DESCENT LLS

→ The formulas must use the whole dataset at once, making it inefficient to compute     → OFFLINE COMPUTATION

$E_n = \frac{1}{2} (y_n - y_n)^2 = \frac{1}{2} (b + v p_n - y_n)^2$       $\frac{\partial E_n}{\partial b} = b + v p_n - y_n \in \mathbb{R}$       $\frac{\partial E_n}{\partial v} = [b + v p_n - y_n] p_n^T = (y_n - y_n) p_n^T$       $\rightarrow$  complexity of one iteration of one sample:  $O(D+1)$

↳ heterogeneous cost      $\nabla_w E_n = \frac{\partial E_n}{\partial w} = [b + v p_n - y_n] p_n^T$       $\rightarrow$   $\frac{\partial E_n}{\partial v} = (y_n - y_n) p_n^T$

• Jacobian Used when loss function is a vector

Step 1) Initialize randomly  $b$  and  $\bar{v}$

Step 2) Set a learning rate  $\eta$  → slow convergence or overshooting

for  $i$  to #iter:

for  $i$  to  $N$ :

$$\text{Compute } \hat{y}_i = b + \bar{v} \cdot x_i$$

$$b' = b - \eta \frac{\partial E_n}{\partial b}$$

$$\bar{v}' = \bar{v} - \eta \frac{\partial E_n}{\partial \bar{v}}$$

• Single batch or more mini-batches  $\bar{g}_i = \frac{1}{|B|} \sum_{x \in B} \frac{\partial E_n}{\partial \bar{v}}$

↓ peak performance when the batch size saturate GPU memory

For  $K$  iterations and whole dataset  $O(KN(D+1))$   
- typically smaller than closed formula

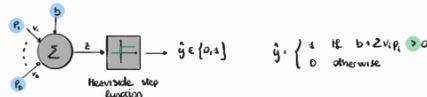
**BENEFITS:**

- efficient for large system of equations → no pseudo inverse computation
- lower memory requirements → not whole  $X$  in memory
- parallelizable

**CLASSIFICATION**

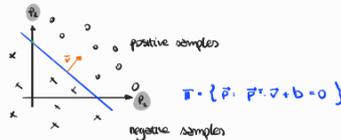
- 2 class classifier (2C-CP) is a function that maps  $x \in \mathbb{R}^D$  to  $y \in \{w_1, w_2\}$
- in multiclass exclusive multi class (MEMC-CP)  $y \in \{w_1, \dots, w_k\} \subset \mathbb{R}$
- in multiclass classification problem (ML-CP) → used in computer vision, like in YOLO  $y \in \mathcal{P}^2$  (any subset, even empty)

**MLP** → should the artificial neuron be a perfect replica of the biological neuron?  
- No, since we are not trying it to understand the brain but to learn new solutions to independent problems → airplanes do not flap their wings



The perceptron is the first ANN ever discovered (1958)

- $|w_i|$  signifies the importance of feature  $i$  in the decision
- $\bar{v}$  is  $\perp$  to the decision boundary, that is shifted by  $b$



The rule for learning  $b$  and  $\bar{v}$  is called **Delta Rule**

- 1) Initialize  $v, b$  randomly
- 2) Compute  $\hat{y}_i = \text{Heaviside}(b + \bar{v} \cdot x_i)$

$$b \leftarrow b - \eta (y_i - \hat{y}_i)$$

$$\bar{v} \leftarrow \bar{v} - \eta (y_i - \hat{y}_i) \bar{x}_i$$

•  $0 < \eta < 1$

• problem  $\frac{d \text{Heaviside}(x)}{dx} = \delta(x)$  → cannot use backpropagation if not been found by trial and error → we need a C' antiderivative

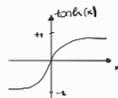
• Problem: It cannot shatter a two classes problem in 2D (when the points are not linearly separable) → eg XOR problem

**ACTIVATION FUNCTION** → their role is to introduce non-linearities

- parameter-less
- parametered → learnable (slows down the convergence)

**Sigmoid (logistic)**  $\sigma(x) = \frac{1}{1+e^{-x}}$

**Hyperbolic tangent (tanh)**  $\tanh(x) = \frac{\sinh(x)}{\cosh(x)} = \frac{e^x - e^{-x}}{e^x + e^{-x}}$



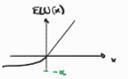
**rectified linear unit (ReLU)**  $\text{ReLU}(x) = \max\{0, x\}$

**leaky ReLU**  $\text{LReLU}(x) = \begin{cases} x & \text{if } x \geq 0 \\ \alpha x & \text{otherwise} \end{cases}$

$\alpha$  = small constant or learnable param (robust ReLU)

**Exponential Linear Unit (ELU)**

$\text{ELU}(x) = \begin{cases} x & \text{if } x \geq 0 \\ \alpha(e^x - 1) & \text{otherwise} \end{cases}$



**Linear Unit**  $\text{LU}(x) = x$

**Swish**  $\text{swish}(x) = x \cdot \sigma(x) = \frac{x}{1+e^{-x}}$

**SiGN**  $\text{sign}(x) = \begin{cases} +1 & \text{if } x > 0 \\ -1 & \text{otherwise} \end{cases}$

• used in spiking NN → almost everywhere differentiable

**MAX OUT**  $\text{maxout}(x_1, x_2) = \max\{x_1, x_2\}$   
• multi-input  $\mathbb{R}^n \rightarrow \mathbb{R}$

**Softmax**  $\text{softmax}(x_i) = \frac{e^{x_i}}{\sum_j e^{x_j}}$

↓ per-layer estimation instead of per-node activation → each node a single output

**Sigmoid Derivative**  $\sigma'(z) = \frac{1}{1+e^{-z}} \cdot \frac{e^{-z}}{(1+e^{-z})^2} = \sigma(z) \cdot (1 - \sigma(z))$

**Tanh Derivative**  $\tau'(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} \cdot \frac{e^z + e^{-z}}{(e^z + e^{-z})^2} = \tau(z) \cdot (1 - \tau(z)^2)$

**Logistic Regression** → 2 mutually exclusive class. same as perceptron but with a sigmoid. Accept as class 1 if  $\sigma(z) \geq \frac{1}{2}$

before a loss called **BINARY CROSS ENTROPY (BCE)**  $\text{BCE} = -y \log(\hat{y}) - (1-y) \log(1-\hat{y})$  → when the two classes are mutually exclusive

$\hat{y}(p) = \sigma(p) = \frac{1}{e^{-(b+\bar{v} \cdot x)} + 1}$

before  $\bar{w} = [b, \bar{v}] \in \mathbb{R}^{1 \times (D+1)}$   
 $\bar{x} = [1, x_1, \dots, x_D]^T$

→  $\hat{y}(x) = \sigma(\bar{w} \cdot \bar{x}) = \frac{1}{1+e^{-\bar{w} \cdot \bar{x}}}$

$\mathcal{L} = -y \log \hat{y} - (1-y) \log(1-\hat{y})$

$\frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w_1}$

$\frac{\partial \mathcal{L}}{\partial \hat{y}} = -\frac{y}{\hat{y}} + \frac{1-y}{1-\hat{y}}$

$\frac{\partial \hat{y}}{\partial z} = \sigma(z) \cdot (1 - \sigma(z)) \cdot \hat{y} \cdot (1 - \hat{y})$

$\frac{\partial z}{\partial w_j} = \frac{\partial \bar{w} \cdot \bar{x}}{\partial w_j} = x_j$

$$\rightarrow \frac{\partial \mathcal{L}}{\partial w_j} = \left[ -\frac{y_j}{\hat{y}_j} + \frac{1-y_j}{1-\hat{y}_j} \right] \hat{y}_j(1-\hat{y}_j) x_j = -(1-\hat{y}_j)y_j x_j + \hat{y}_j(1-y_j) x_j = -y_j x_j + \hat{y}_j x_j = (\hat{y}_j - y_j) x_j$$

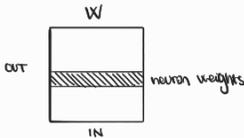
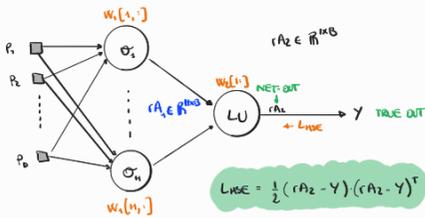
$$\rightarrow \frac{\partial \mathcal{L}}{\partial w} = (\hat{y} - y) \vec{x}$$

In case of a minibatch  $X \in \mathbb{R}^{b \times (b+1)}$ ,  $X = \begin{bmatrix} \text{first column of ones} \\ \vdots \\ \text{pattern matrix } P \in \mathbb{R}^{b \times b} \end{bmatrix}$   $\rightarrow \frac{\partial \mathcal{L}(w)}{\partial w} = (\hat{y} - y) X$  (row vectors)

The problem is convex in  $\vec{w}$ , since BCE is convex in  $z$ :  $-y \log(\sigma(z)) - (1-y) \log(1-\sigma(z))$  and  $z$  is linear in  $\vec{w}$   
 $\rightarrow$  converges always to the global solution independently from initialization of  $b$  and  $w$

No closed solution exists  $\nabla_w \mathcal{L} = (\hat{y} - y) \vec{x} = [O(WX) - y] \vec{x} \rightarrow$  this can be solved, but not in the case of mini-batch

**HLP - MSE LIN**  $\rightarrow$  use a layer of neurons with non-linear activations and a final neuron with LU. Implement the loss as MSE  
 - if I need multi-output regression I could add more final neurons



**FWD**

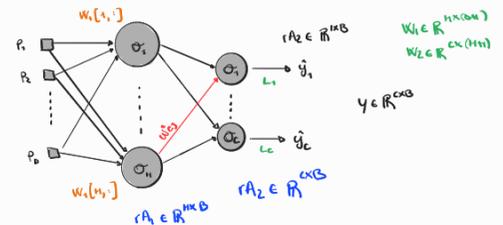
$$\begin{aligned} rA_0 &= P^T \in \mathbb{R}^{b \times b} \\ A_0 &= [\vec{1}_b, rA_0] \in \mathbb{R}^{(b+1) \times b} \\ rZ_1 &= W_1 A_0 \in \mathbb{R}^{n \times b} \\ rA_1 &= \sigma(rZ_1) \\ A_1 &= [\vec{1}_n, rA_1] \in \mathbb{R}^{(n+1) \times b} \\ rZ_2 &= W_2 A_1 \in \mathbb{R}^{c \times b} \\ rA_2 &= \text{Tr}(rZ_2) \in \mathbb{R}^{c \times b} \end{aligned}$$

**BWD**

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial W_1} &= \frac{\partial \mathcal{L}}{\partial rZ_1} \frac{\partial rZ_1}{\partial W_1} = \frac{\partial \mathcal{L}}{\partial rZ_1} A_0^T \in \mathbb{R}^{n \times (b+1)} \\ \frac{\partial \mathcal{L}}{\partial rZ_1} &= \frac{\partial \mathcal{L}}{\partial rA_1} \frac{\partial rA_1}{\partial rZ_1} = \frac{\partial \mathcal{L}}{\partial rA_1} \circ \sigma'(rZ_1) \in \mathbb{R}^{n \times b} \\ \frac{\partial \mathcal{L}}{\partial rA_1} &= \frac{\partial \mathcal{L}}{\partial A_1} \frac{\partial A_1}{\partial rA_1} = \begin{pmatrix} 0, I_n \end{pmatrix} \frac{\partial \mathcal{L}}{\partial A_1} \in \mathbb{R}^{(n+1) \times b} = \text{red} \left( \frac{\partial \mathcal{L}}{\partial A_1} \right) \\ \frac{\partial \mathcal{L}}{\partial W_2} &= \frac{\partial \mathcal{L}}{\partial rZ_2} \frac{\partial rZ_2}{\partial W_2} = (rA_2 - Y) A_1^T \in \mathbb{R}^{(c+1) \times b} \\ \frac{\partial \mathcal{L}}{\partial rZ_2} &= \frac{\partial \mathcal{L}}{\partial rA_2} \frac{\partial rA_2}{\partial rZ_2} = W_2^T (rA_2 - Y) \in \mathbb{R}^{n \times c} \end{aligned}$$

considering multiple expts without non-linearity  $rA_0 = P^T$ ,  $rZ_1 = W_1 A_0$ ,  $rA_1 = \sigma(rZ_1)$ ,  $rZ_2 = W_2 A_1$ ,  $rA_2 = \text{Tr}(rZ_2) \Leftrightarrow rA_2 = W_2 \cdot W_1 \cdot P^T \Leftrightarrow rA_2 = W \cdot P^T$   
 $\rightarrow$  not completely free when using SGB and Dropout  
 - generalizable with reduction

shallow NN  $\rightarrow$  one hidden or no one (+ initial + final)  $\begin{bmatrix} 1 & 0 \\ \vdots & \vdots \\ 1 & 0 \end{bmatrix}$   
 deeper networks  $\rightarrow$  hierarchical representations of the input



**MLP - BCE SIG**

With  $C$  classes (not exclusive) we need  $C$  output neurons with a sigmoid activation

$$L = \sum_{c=1}^C L_c \quad \text{with} \quad L_c = - \sum_{b=1}^b \left[ y_c^b \log(rA_{2c}^b) + (1-y_c^b) \log(1-rA_{2c}^b) \right]$$

$\hookrightarrow$  BCE for neuron  $c$

b-th pattern:  $\frac{\partial L}{\partial rA_{2c}^b} = \frac{\partial L_c}{\partial rA_{2c}^b} = \frac{\partial L_c}{\partial rZ_{2c}^b} = -\frac{y_c^b}{rA_{2c}^b} + \frac{1-y_c^b}{1-rA_{2c}^b}$

$$rA_{2c} = \sigma(rZ_{2c}) \rightarrow \frac{\partial rA_{2c}^b}{\partial rZ_{2c}^b} = \begin{cases} rA_{2c}^b (1-rA_{2c}^b) & \text{if } k=c \\ 0 & \text{otherwise} \end{cases} \Rightarrow \frac{\partial L}{\partial rZ_{2c}^b} = \sum_{k=1}^c \frac{\partial L}{\partial rA_{2c}^k} \frac{\partial rA_{2c}^k}{\partial rZ_{2c}^b} = \frac{\partial L}{\partial rA_{2c}^b} \frac{\partial rA_{2c}^b}{\partial rZ_{2c}^b} = -\frac{y_c^b}{rA_{2c}^b} \cdot rA_{2c}^b (1-rA_{2c}^b) + \frac{1-y_c^b}{1-rA_{2c}^b} \cdot rA_{2c}^b (1-rA_{2c}^b)$$

$$\frac{\partial L}{\partial W_{1j}} = \sum_{k=1}^c \frac{\partial L}{\partial rZ_{2c}^k} \frac{\partial rZ_{2c}^k}{\partial W_{1j}} = \sum_{k=1}^c (rA_{2c}^k - y_c^k) \frac{\partial rZ_{2c}^k}{\partial W_{1j}} = (rA_{2c}^b - y_c^b) rA_{1j}^b$$

$\hookrightarrow W_2 \in \mathbb{R}^{c \times (n+1)}$

$$\frac{\partial L}{\partial A_{1j}} = \sum_{k=1}^c \frac{\partial L}{\partial rZ_{2c}^k} \frac{\partial rZ_{2c}^k}{\partial A_{1j}} = \sum_{k=1}^c \frac{\partial L}{\partial rZ_{2c}^k} W_{2kj} \rightarrow \frac{\partial L}{\partial A_{1j}} = W_2^T \frac{\partial L}{\partial rZ_{2c}}$$

$$\frac{\partial L}{\partial rZ_{1j}} = \sum_{k=1}^c \frac{\partial L}{\partial rA_{1k}} \frac{\partial rA_{1k}}{\partial rZ_{1j}} = \frac{\partial L}{\partial rA_{1j}} \sigma'(rZ_{1j}) \quad \frac{\partial L}{\partial rA_{1j}} = \text{red} \left( \frac{\partial L}{\partial A_{1j}} \right)$$

$$\frac{\partial L}{\partial W_1} = \frac{\partial L}{\partial rZ_1} \frac{\partial rZ_1}{\partial W_1} = \frac{\partial L}{\partial rZ_1} A_0^T \quad \mathbb{R}^{(n+1) \times b} \cdot \mathbb{R}^{b \times (b+1)} = \mathbb{R}^{(n+1) \times (b+1)}$$

resemble the MSE-LIN  $\leftarrow$

$$\begin{aligned} rA_0 &= P^T \in \mathbb{R}^{b \times b} \\ A_0 &= \text{ext}(rA_0) \in \mathbb{R}^{(b+1) \times b} \\ rZ_1 &= W_1 A_0 \in \mathbb{R}^{n \times b} \\ rA_1 &= \sigma(rZ_1) \in \mathbb{R}^{(n+1) \times b} \\ A_1 &= \text{ext}(rA_1) \in \mathbb{R}^{(n+1) \times b} \\ rZ_2 &= W_2 A_1 \in \mathbb{R}^{c \times b} \\ rA_2 &= \sigma(rZ_2) \in \mathbb{R}^{c \times b} \end{aligned} \quad \begin{aligned} \frac{\partial \mathcal{L}}{\partial W_1} &= \frac{\partial \mathcal{L}}{\partial rZ_1} \frac{\partial rZ_1}{\partial W_1} = \frac{\partial \mathcal{L}}{\partial rZ_1} A_0^T \in \mathbb{R}^{(n+1) \times (b+1)} \\ \frac{\partial \mathcal{L}}{\partial rZ_1} &= \frac{\partial \mathcal{L}}{\partial rA_1} \frac{\partial rA_1}{\partial rZ_1} = \frac{\partial \mathcal{L}}{\partial rA_1} \circ \text{red} \left( \frac{\partial \mathcal{L}}{\partial A_1} \right) \circ \sigma'(rZ_1) \in \mathbb{R}^{(n+1) \times b} \\ \frac{\partial \mathcal{L}}{\partial rA_1} &= \frac{\partial \mathcal{L}}{\partial A_1} \frac{\partial A_1}{\partial rA_1} = W_1^T \frac{\partial \mathcal{L}}{\partial rZ_2} \in \mathbb{R}^{(n+1) \times b} \\ \frac{\partial \mathcal{L}}{\partial rZ_2} &= \frac{\partial \mathcal{L}}{\partial rA_2} \frac{\partial rA_2}{\partial rZ_2} = \frac{\partial \mathcal{L}}{\partial rA_2} A_1^T \in \mathbb{R}^{c \times (n+1)} \\ \frac{\partial \mathcal{L}}{\partial rA_2} &= rA_2 - Y \in \mathbb{R}^{c \times b} \end{aligned}$$

MLP-CESM → classes assumed to be mutually exclusive ⇒ we categorical CE and softmax activation

softmax is a vectorial function  $\mathbb{R}^c \rightarrow \mathbb{R}^c$  with  $a_k = \frac{e^{z_k}}{\sum_j e^{z_j}}$

$r_{z_c}^b = \text{softmax}(r_{z_c}^b) \rightarrow \frac{\partial r_{z_c}^b}{\partial z_c^b} = \begin{cases} r_{z_c}^b (1 - r_{z_c}^b) & \text{if } c=k \\ -r_{z_c}^b \cdot r_{z_k}^b & \text{if } c \neq k \end{cases}$

if  $c=k \rightarrow \frac{e^{z_c^b}}{\sum_j e^{z_j^b}} \cdot \frac{1}{\sum_j e^{z_j^b}} = \frac{e^{z_c^b}}{(\sum_j e^{z_j^b})^2} = r_{z_c}^b (1 - r_{z_c}^b)$   $[\frac{\partial}{\partial z} (\frac{f}{g}) = \frac{\partial f}{\partial z} \frac{1}{g} + f \cdot \frac{-1}{g^2} \frac{\partial g}{\partial z}]$

if  $c \neq k \rightarrow -\frac{e^{z_c^b}}{(\sum_j e^{z_j^b})^2} \cdot e^{z_k^b} = -r_{z_c}^b \cdot r_{z_k}^b$

$L = \sum_c \sum_b -y_c^b \log(r_{z_c}^b)$  one-hot encoding

$\frac{\partial L}{\partial r_{z_c}^b} = -\frac{y_c^b}{r_{z_c}^b}$

$\frac{\partial r_{z_c}^b}{\partial z_c^b} = \begin{cases} r_{z_c}^b (1 - r_{z_c}^b) & \text{if } c=k \\ -r_{z_c}^b \cdot r_{z_k}^b & \text{if } c \neq k \end{cases} \Rightarrow \frac{\partial L}{\partial z_c^b} = \sum_{k=1}^c \frac{\partial L}{\partial r_{z_k}^b} \frac{\partial r_{z_k}^b}{\partial z_c^b} = \frac{\partial L}{\partial r_{z_c}^b} \frac{\partial r_{z_c}^b}{\partial z_c^b} + \sum_{k \neq c} \frac{\partial L}{\partial r_{z_k}^b} \frac{\partial r_{z_k}^b}{\partial z_c^b}$

$= -\frac{y_c^b}{r_{z_c}^b} r_{z_c}^b (1 - r_{z_c}^b) + \sum_{k \neq c} -\frac{y_k^b}{r_{z_k}^b} \cdot -r_{z_k}^b \cdot r_{z_c}^b$

$= -y_c^b + r_{z_c}^b \sum_{k=1}^c y_k^b = r_{z_c}^b - y_c^b$

$\frac{\partial L}{\partial w_{c_j}^b} = \sum_k \frac{\partial L}{\partial z_k^b} \frac{\partial z_k^b}{\partial w_{c_j}^b} = \sum_k (r_{z_k}^b - y_k^b) \frac{\partial z_k^b}{\partial w_{c_j}^b} = (r_{z_c}^b - y_c^b) r_{A_1}^b$

$\frac{\partial L}{\partial A_1^b} = \sum_{k=1}^c \frac{\partial L}{\partial z_k^b} \frac{\partial z_k^b}{\partial A_1^b} = \sum_{k=1}^c \frac{\partial L}{\partial z_k^b} w_{k1}^b$

$\frac{\partial L}{\partial z_j^b} = \sum_{k=1}^H \frac{\partial L}{\partial A_k^b} \frac{\partial A_k^b}{\partial z_j^b} = \frac{\partial L}{\partial A_j^b} \sigma'(r_{z_j}^b)$

$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial z_c^b} \frac{\partial z_c^b}{\partial w_i} = \frac{\partial L}{\partial z_c^b} A_0^T$

$r_{A_0} = p^T \in \mathbb{R}^{D \times D}$

$A_0 = \text{act}(r_{A_0}) \in \mathbb{R}^{(1+b) \times b}$

$r_{z_1} = W_1 A_0 \in \mathbb{R}^{H \times b}$

$r_{A_1} = \sigma(r_{z_1}) \in \mathbb{R}^{H \times b}$

$A_1 = \text{act}(r_{A_1}) \in \mathbb{R}^{(1+H) \times b}$

$r_{z_2} = W_2 A_1 \in \mathbb{R}^{c \times b}$

$r_{A_2} = \text{softmax}(r_{z_2}) \in \mathbb{R}^{c \times b}$

$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial z_1} \frac{\partial z_1}{\partial w_1} = \frac{\partial L}{\partial z_1} A_0^T \in \mathbb{R}^{H \times (1+b)}$

$\frac{\partial L}{\partial z_1} = \frac{\partial L}{\partial A_1} \frac{\partial A_1}{\partial z_1} = \text{vec}(\frac{\partial L}{\partial A_1}) \odot \sigma'(r_{z_1}) \in \mathbb{R}^{H \times b}$

$\frac{\partial L}{\partial A_1} = \frac{\partial L}{\partial z_2} \frac{\partial z_2}{\partial A_1} = W_2^T \frac{\partial L}{\partial z_2} \in \mathbb{R}^{(1+H) \times b}$

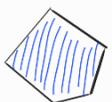
$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial z_2} A_1^T \in \mathbb{R}^{c \times (1+H)}$

$\frac{\partial L}{\partial z_2} = r_{A_2} - y \in \mathbb{R}^{c \times b}$

CONSIDERATION WITH C=2

1 hidden layer: n hyperplanes defined by  $w_i \cdot x + b_i = 0$  (since ReLU activation)

the output is  $\sigma = \sum_{i=1}^n w_i \text{ReLU}(z_i)$   $z_i = w_i \cdot x + b_i$  → Intersection of active hyperplanes that is a convex region (POLYTOPE)



2 hidden layers: I can represent union of convex polytopes also disconnected → no more convex  
- each of them can be mapped differently

3 hidden layers: can represent structure with holes



Insufficient training data can lead to overfitting, while excessive increase in training time

- 1) Held out position (train-Val-Test) → performances depends only on the randomness of a single position  
- only a fraction of data is used for training
- 2) K-Fold cross validation → Divide Train+Val in K partition, then K-1 for training and the last for validation  
- reduce significantly the variance of Val performance estimate
- 3) Stratified K-Fold → important for imbalanced dataset
- 4) Leave one out cross validation → K=# of samples  
- too expensive and not used in DL
- 5) Time series splitting → here random sampling is invalid  
- cannot use future data to train and then test on the past since it violates causality
  - non-random split (sequential)
    - $[t_0, t_i]$  training
    - $[t_i, t_j]$  val
    - $[t_j, t_n]$  test

• rolling cross validation where val/test window moves forward → preserves causality

Weight Initialization → improper init will lead to vanishing or exploding gradients

• XAVIER initialization: given a layer  $n_{in}$  = # of inputs,  $n_{out}$  = # of outputs → initialize  $W$  such that it has a constant variance across layers  
 keeps the variance of the input and of backprop consistent

FWD: If  $M_x = 0$  and  $\sigma_x = \text{Var}(x)$  →  $\text{Var}(y_j) = n_{in} \text{Var}(W) \text{Var}(x)$  since  $y_j = W^T x + b_j$  →  $\text{Var}(y_j) = \sum_{i=1}^{n_{in}} \text{Var}(W_{ij} x_i) = \sum_{i=1}^{n_{in}} \text{Var}(W_{ij}) \text{Var}(x_i)$   
 • to maintain  $\text{Var}(y_j) = \text{Var}(x)$  we need  $n_{in} \text{Var}(W) = 1$

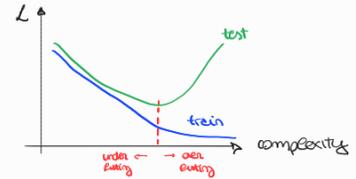
BWD:  $n_{out} \text{Var}(W) = 1$

• To satisfy both conditions we set  $\text{Var}(W) = \frac{2}{n_{in} + n_{out}}$  → Applicable if the function is symmetric around zero (tanh or linear)

Batch Norm  $\hat{x} = \gamma \frac{x - M_B}{\sqrt{\sigma_B^2 + \epsilon}} + \beta$  → mitigates covariance shifts during training → stabilize training + speed up  
 for a single neuron  $\gamma$  and  $\beta$  allows the network to learn an optimal range for each feature during inference use  $M_B$  and  $\sigma_B$  learned during training

DropOut → regularization technique by randomly dropping out a subset of neurons during training

- avoids being over-reliant on a neuron output
- in inference scale the output by  $\frac{1}{1-p}$  to maintain consistency

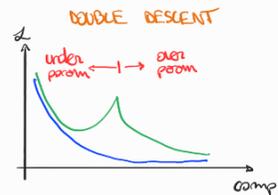


Bias-Variance Trade-Off  
 - overfitting when test error  $\gg$  train error  
 - underfitting when training error stops decreasing

rule 2012 → rule of thumb: keep the # of params  $\leq \frac{|D|}{5}$   
 • 5 observation per free-parameter

Why large models don't overfit (hypothesis):

- Lottery Hypothesis: because they contain subnetworks that are selected during training and can generalize well
- Implicit regularization of SGD that promotes simpler models (min  $\sum w_i^2$  eq)
- The effective model capacity is limited by SGD + regularization + augmentation + early stopping
- Batch Norm has some regularization effects by preventing to rely too heavily on a single feature



sample from  $p(x, y)$   $K$  datasets, then  $A(D_i) \rightarrow \hat{f}_i(x_0)$   
 [  $y = f(x) + \epsilon$  ]  
 learning function  $A(D_i) \rightarrow \hat{f}_i(x_0)$   
 evaluation point  $x_0$

→ expected prediction error  $EPE_A(x_0) = \mathbb{E}_{D_1, p(y|x_0)} [(y - \hat{f}_n(x_0))^2] \approx \frac{1}{K} \sum_{i=1}^K \mathbb{E}_{y \sim p(y|x_0)} (y_i - \hat{f}_n(x_0))^2$   
 $= \sigma_\epsilon^2 + (\mathbb{E}_n[\hat{f}_n(x_0)] - f(x_0))^2 + \mathbb{E}_n[(\hat{f}_n(x_0) - \mathbb{E}_n[\hat{f}_n(x_0)])^2]$   
 irreducible error  $\sigma_\epsilon^2$   
 randomness on  $y$  that doesn't depend on  $x$

→ typically decreasing one will increase the other (but regularization will decrease var, but increase bias)

• more params allows to approximate  $f(x_0)$  correctly, but increases variance and is due to noise ( $\epsilon$ )

only linear relations

Pearson statistics measures the linear dependency of two vectors  $R \in [-1, +1]$

$$R = \frac{\text{cov}(y, \hat{y})}{\sigma_y \cdot \sigma_{\hat{y}} \rightarrow \sqrt{\text{Var}(\hat{y})}}$$

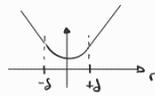
Error function: is a more general term that can be used to describe different discrepancy metrics (MAE, MSE, ...)  
 - used as a measure to assess the model performance

Loss function: what it is optimized

- MAE =  $\frac{1}{N} \sum |y_i - \hat{y}_i|$  Differentiable almost everywhere ( $\hat{y}_i \neq y_i$ )  
 → used in regression since it is more robust to outliers

$$L_S(r) = \begin{cases} \frac{1}{2}r^2 & |r| \leq \delta \\ \delta(|r| - \frac{1}{2}\delta) & |r| > \delta \end{cases}$$

HUBER LOSS



$L_S(r) \in C^1$  → robust to outliers  
 - larger  $\delta$  more  $L_2$ , smaller more  $L_1$

$$HL = \frac{1}{N} \sum_{i=1}^N \max(0, 1 - y_i \hat{y}_i)$$

HINGE LOSS

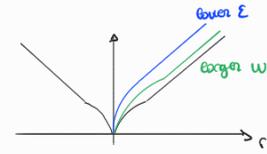
$\hat{y}_i \in \{-1, +1\}$

→ used in SVM for binary classification → penalize  $\hat{y}_i$  closer to the decision boundary

$$KLD = \frac{1}{N} \sum y_i \log \frac{y_i}{\hat{y}_i} \quad D_{KL}(P \parallel Q) = \sum_{x \in X} p(x) \log \frac{p(x)}{q(x)} = \mathbb{E}_{p(x)} \left[ \log \frac{p(x)}{q(x)} \right]$$

$$L_{\text{wing}}(r) = \begin{cases} w \ln(1 + \frac{|r|}{c}) & |r| < w \\ |r| - c & |r| \geq w \end{cases}$$

with  $c = w - w \ln(1 + \frac{w}{\epsilon})$   
 - params  $w, \epsilon > 0$   
 - continuous in  $|r| = w$   
 - reduce outliers since grows linearly  
 - steeper gradient around zero



REGULARIZATION → parameter surface will have fewer local minima and a shifted global minima

**L1 (Lasso)**  $+\lambda \sum |w_i|$  → encourages sparsity some  $w_i = 0$  (eg. when some input features are redundant) → feature selection ^ interpretable model

$$+\log P(w|X, y) = \log P(y|X, w) + \log P(w) = \underbrace{-\frac{1}{2\sigma^2} \|y - Xw\|^2}_{\text{regression}} - \underbrace{\sum \frac{|w_i|}{b}}_{\text{prior}}$$

MAP  $\Leftrightarrow P(y|X, w) = \mathcal{N}(Xw, \sigma^2) \wedge P(w_i) = \frac{1}{2b} e^{-\frac{|w_i|}{b}}$

**L2 (ridge)**  $+\lambda \sum w_i^2$  → encourages smaller weights  
 - penalization on all terms rather than  $L_1 \Rightarrow$  smoother model  
 - better generalization

In regression if define  $E = \sum_{i=1}^m (y_i - p(x_i))^2 + \lambda w^T w \Rightarrow p = (A^T A + \lambda I_m)^{-1} A^T y$

**Offline Learning**

- entire training processed at once (batches)
- repeat multiple epochs - weights update when all points are seen
- requires going through all data for each epoch

**Online Learning**

- individual points or small batches are processed sequentially
- weights are updated after each point
- learns continuously when new data is available
- handles streaming of data, adapt to changing environments
- efficient for large dataset, but can stable than batch learning

### GRADIENT DESCENT

- learning rate can be manually tuned (starting small and in case increase until divergence) or by a grid search or we could use a scheduler in order to not have it fixed during training:

- 1) step decay: reduce by a factor  $\times$  every  $y$  epoch
  - 2) exponential decay over epochs
  - 3) Adaptive methods (like Adam, RMSprop, ...) → automatically adjust lr based on gradient history
- validation set performance decreasing lr if doesn't improve

- if the training loss oscillates too much  $\Rightarrow$  decreases lr, while if too slow to decrease increase it

SGD  $\Theta_{t+1} = \Theta_t - \eta \nabla_{\Theta} \mathcal{L}(\Theta_t, x_i, y_i)$

mini-batch

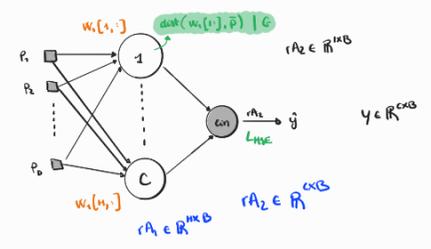
**AdaGrad**  $\theta_0^i = 0 \Rightarrow g_t = \nabla_{\theta} \mathcal{L}(\theta_{t-1})$   $\Rightarrow \theta_t^i = \theta_{t-1}^i - \frac{\eta}{\sqrt{G_t^i + \epsilon}} g_{t,i}$   
 - typically in batches  
 - each parameter has its own learning rate based on its past updates  
 - doesn't work for long strings  
 - never decreases  
 - avoid zero division

**RMSProp** root mean square propagation  
 $v_{t+1}^i = \rho v_t^i + (1-\rho) g_t^2(\theta_i)^2 \rightarrow$  gradient moving average, solves AdaGrad problem  $\rho \rightarrow$  decay rate of exponential decay  
 $\theta_{t+1}^i = \theta_t^i - \eta \frac{\nabla_{\theta} \mathcal{L}(\theta_t)_i}{\sqrt{v_{t+1}^i + \epsilon}}$   
 - variance of the gradients

**Adam**  $m_{t+1}^i = \beta_1 m_t^i + (1-\beta_1) g_t(\theta_i)$   $\hat{m}_{t+1}^i = \frac{m_{t+1}^i}{1-\beta_1^{t+1}}$   
 $v_{t+1}^i = \beta_2 v_t^i + (1-\beta_2) g_t^2(\theta_i)$   $\hat{v}_{t+1}^i = \frac{v_{t+1}^i}{1-\beta_2^{t+1}}$   
 $\Rightarrow \theta_{t+1}^i = \theta_t^i - \eta \frac{\hat{m}_{t+1}^i}{\sqrt{\hat{v}_{t+1}^i + \epsilon}}$   
 - bias correction otherwise forward pass for initialization  
 - moving avg of squared grad (second moment)

**RADIAL BASIS FUNCTION NETWORKS**  $\rightarrow$  3 layers network - used for regression problems  
 $\varphi(\vec{x}) = \varphi(\|\vec{x} - \mathbf{c}\|)$   $\varphi: [0, \infty) \rightarrow \mathbb{R}$   
 - depends only on the distance to a center

Typically the activation function of the hidden layer is a **gaussian** with higher activation if the feature is near the center  
 $G(x|\mu, \sigma) = e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$   $\rightarrow$  the normalization is not interesting  
 $G_2(x|\mu, \sigma) = e^{-\frac{1}{2\sigma^2}(\vec{x}-\vec{\mu})^T(\vec{x}-\vec{\mu})}$   $\Rightarrow z(w, x) = (w-x)^2$  and  $G_2(z|\epsilon) = e^{-z^2/\epsilon^2}$   
 - perceptron function  
 - gaussian, not optimized



For multivariate case:  $z(\vec{w}, \vec{x}) = \|\vec{w} - \vec{x}\|_2 = \sqrt{\sum_{d=1}^D (x_d - w_d)^2}$   
 - centroid

I could use mini-batches  $W_i \in \mathbb{R}^{C \times D}$ ,  $r_{A0} \in \mathbb{R}^{D \times B} \Rightarrow r_{2i} = \text{dist}(W_i, r_{A0}) \in \mathbb{R}^{C \times B}$  with  $r_{2i}(i, j) = \text{dist}(W_i[i, :], r_{A0}[i, :])$   
 -  $i$ -th row  
 -  $j$ -th sample in mini-batch

$\rightarrow$  Weights of the hidden layer axes in the same dimension of our input  
 $\rightarrow$  we know how to initialize the centroids based on a clustering algorithm

$p(x) = \frac{1}{(2\pi)^{D/2} |Z|^{D/2}} \exp(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu))$  we contain the matrix to be diagonal and with constant variance  
 $G(\vec{x}|\vec{\mu}, \sigma) = \exp(-\frac{1}{2\sigma^2} (x-\mu)^T (x-\mu))$   
 $G_2(z) = \exp(-z^2/\epsilon^2)$   
 -  $\epsilon = \frac{1}{\sqrt{2\sigma}}$

Typically trained offline and with a single batch  $\rightarrow$  If we fixed  $W_i$  using clustering algorithm (suboptimal) the MSE w.r.t  $W_2$  is linear and can be computed with the pseudo-inverse or F could we online learning with backprop  $W_1$  and  $W_2$

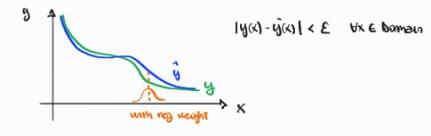
$$r_{2i} = \text{dist}(W_i, \rho) = \begin{bmatrix} \|W_i^1 - \rho\| & \dots & \|W_i^C - \rho\| \\ \vdots & & \vdots \\ \|W_i^1 - \rho\| & \dots & \|W_i^C - \rho\| \end{bmatrix} \quad r_{A0} = G_2(\text{dist}(W_i, \rho)) \quad \Rightarrow \hat{y}(\hat{\rho}) = W_2^T + \sum_{i=1}^C W_2^i G_2(\|W_i^i - \rho\|)$$

-  $i$ -th centroid  
 -  $i$ -th row as  $W_i$   
 - operation of the linear combination of the kernels (typically denoted with  $\lambda_i$ )

$$\hat{y}(\hat{\rho}) = \lambda_0 + \sum_{i=1}^C \lambda_i e^{-z^2/\epsilon^2}$$

$\rightarrow$  write the derivative for back-prop

$\rightarrow$  can be demonstrated that this is a universal approximator of the radial kernel has infinite support  
 - Similar to kernel density estimation (KDE) that is used to model pdf  
 - gaussian kernel will leads to vanishing gradients (when we are far from the center)  
 - I could embed RBF within an AE



- Ockem's Razor: when two hypotheses have the same explanatory power we should select the one with fewest assumptions
- Minimum Description Length (MDL): select the model with shortest encoding (model + data given) if some metric  $\rightarrow$  more interpretable and generalizable
- Greedy degradation principle: when something goes wrong fall back to minimal functionality instead of crashing
- The progressive enhancements principle: bottom-up improvement of a system  $\rightarrow$  we models that can be expanded

# EVOLUTIONARY ALGORITHMS

→ finding the global optimum of non-convex (hows not-linear) and non-differentiable is NP-hard

· If domain is discrete → exhaustive search (not feasible)

· If continuous I could discretize between min/max → grid search (curse of dimensionality)

Random Search is highly inefficient since it doesn't exploit informations. Stochastic Search is better

## ANNEALING

· In physics heats a metal and cool it slowly  
 - softens the material, restore original properties ...

· Tempering heats it at a lower temp to make it stronger

Simulated Annealing → probabilistic optimization algorithm, used to find near optimum solutions

· It moves through neighbor solutions, accepting worse solutions proportional to a temperature parameter  
 - as temp decreases I will accept worse solutions more difficultly

- I could have linear, log or exp cooling
- EXPLORATION vs EXPLOITATION tradeoff
- could be used for HT in ML
- proven that sim. annealing asymptotically find the global optimum

- 1) Initialize temp, solution, schedule
- 2) generate a neighbor solution
- 3) if  $f(x_{t+1}) < f(x_t)$  →  $x_t$  new solution  
 else accept  $x_t$  according to a criterion as Metropolis
- 4) adjust temp and if not finished go to 2)
- 5) return the best solution found

- 1) Evolution Strategies → simplest EA that uses mutation only
- 2) Genetic Algorithms → solutions as strings (typically binary), uses mutation/cross-over
- 3) Genetic Programming → solutions are in the form of computer programs (typically trees) and are evaluated to solve a target problem
- 4) NeuroEvolution → similar to above but solutions represents NN (structure or weights)

} numerical optimization

one offspring generated per epoch  
 (1+1) ES  
 maintains a single individual

using some heuristic

- 1) Initialization: parent solution  $x_0, t = 0$
- 2) Mutation → random/gaussian
- 3) Evaluation
- 4) Selection (if it is better)
- 5) Termination or goto 2)

· Typically used as baseline for new models  
 · Memory less strategy

$(\mu + \lambda)$  ES  
 maintains a population

· we have a pool of best solutions → memory  
 · still **memoryless** since it depends only on  $P_{t-1}$   
 - allows to use cross-over  
 - can be //

GENETIC ALGORITHMS → are meta-heuristic algorithms (higher level frameworks to solve difficult optimization problems)  
 ↓ designed to navigate large search space efficiently

- used for problems that are difficult, large or lacks concise math formulation → no assumptions
- mimics Darwin Evolution
- 2 Population individual is encoded via a chromosome (string)
- offspring selection is normally carried out via fitness
- Offspring population substitutes partially or completely

• Fitness function can be equivalent to the objective function or be a modification (eg. to handle constraints)

```

function [best_x, best_y] = ES_MU_plus_LAMBDA(mu, lambda, ngenes)
% Initialize population P with mu random individuals
P = initializePopulation(mu, ngenes);

% Evaluate fitness of each individual in P
fit_P = evaluateFitness(P);

% Set the current generation
g = 0;

% Repeat until termination condition is met
while terminationConditionNotMet()
% Create an empty offspring population Q
Q = [];

% Repeat until Q is of size lambda
while size(Q, 1) < lambda
% Select a parent from P using any selection method (random, binary tou
parent = selectionMethod(P);

% Perform mutation on the parent to create an offspring
offspring = mutation(parent);

% Add the offspring to Q
Q = [Q; offspring];
end

% Evaluate fitness of each individual in Q
fit_Q = evaluateFitness(Q);

% Combine P and Q to obtain a combined population R
R = [P; Q]; fit_R = [fit_P; fit_Q];

% Sort population R in descending order based on fitness
[~, indices] = sort(fit_R, 'descend');

% Select the top mu individuals from R to form the next generation P
P = R(indices(1:mu), :); fit_P = fit_R(indices(1:mu));

% Increment generation counter g by 1
g = g + 1;
end

% Return the best individual in the final population P as the result
[~, best_index] = min(fit_P);
best_x = P(best_index, :); best_y = fit_P(best_index);
end

```

```

function [x_opt, f_opt] = GA_BE_1D(f, a, b, generations, chrsize, popsize, pmutallgenes, displayAndPlot)
% Generate the initial population
P = BE_initpop(chrsize, popsize);

% Computing the fitness of each individual in the initial population
fit_P = evaluatefitness(f, P, a, b);

for t = 2:generations
% 1. Create the Mating Pool (MP) using binary tournament selection
MP = binarytournament(P, fit_P, popsize);
% - if both violates → random cv(x) - if both feasible the one with better f

% 2. Reshuffle the Mating Pool
MP = MP(randperm(popsize), :);

% 3. Apply the crossover operator to each pair in the MP, to obtain offsprings Q
for i=1:2:popsize
Q([i, i+1], :) = BE_XOVER_singlepoint(MP(i, :), MP(i+1, :));
end

% 4. Apply mutation operator to obtain the mutated version of Q
for i=1:popsize
BE_MUT_unif_all(Q(i, :), pmutallgenes);
end

% 5. Evaluate the fitness of Q
fit_Q = evaluatefitness(f, Q, a, b);

% 6. Compute R as the union of P and Q
R = [P; Q]; fit_R = [fit_P; fit_Q];

% 7. Compute the new population P, by taking the best individuals from R
[~, ind] = sort(fit_R, "ascend");
new_P = R(ind(1:popsize), :); fit_new_P = fit_R(ind(1:popsize));

% 8. Let the new population new_P become the current population and do the same for the fitness
P = new_P; fit_P = fit_new_P;
end

[f_opt, ind] = min(fit_P);
x_opt = logical2real(P(ind, :), a, b);

```

**BINARY ENCODING:** each gene is a bit

**MUTATION**

- 1) **single point** mutation: randomly select a gene and negate it
- 2) **uniform**: define a probability to mutate a gene and do it on all

if too high → destroy genetic information of the offspring  
if too low → premature convergence

**CROSS-OVER**

- 1) **one point CO**: select a point and swaps the genes that comes after

$$O_{2i} = \begin{cases} P_{1i} & \text{if } i \text{ is crossover-point} \\ P_{2i} & \text{otherwise} \end{cases} \quad O_{1i} = \begin{cases} P_{2i} & \text{if } i \text{ is crossover-point} \\ P_{1i} & \text{otherwise} \end{cases}$$

- 2) **two point cross-over**  $O_{1i} = \begin{cases} P_{1i} & \text{if } i \text{ is } \text{co-point}_1 \vee i \text{ is } \text{co-point}_2 \\ P_{2i} & \text{otherwise} \end{cases}$

**REAL ENCODING** → each gene is a real number

**MUTATION**

- 1) **Gaussian Mutation**: mutated gene  $H_i = G_i + \sigma N(0, 1)$
- 2) **Uniform Mutation**:  $H_i = G_i + \Delta_i U(-1, +1)$   
typically  $\frac{1}{10}$  or  $\frac{1}{20}$  the len of the domain of  $i$  variable

**CROSSOVER**

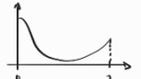
- 1) **discrete crossover**  $O_{2i} = \begin{cases} P_{1i} & \text{if } \text{rand}() \leq P_c \\ P_{2i} & \text{otherwise} \end{cases}$
- 2) **Arithmetic crossover**  $O_{2i} = \alpha_i P_{1i} + (1-\alpha_i) P_{2i}$   $\alpha_i \in \text{random in } [0, 1]$
- 3) **complex** → Arithmetic CO with  $\alpha_i \in [0, 1]$  fixed, typically 0.5

- 4) **simulated Binary CO (SBX)**: inspired by single point CO

$$O_{1i} = \frac{1}{2} [(1+\beta_q) P_{1i} + (1-\beta_q) P_{2i}]$$

$$O_{2i} = \frac{1}{2} [(1-\beta_q) P_{1i} + (1+\beta_q) P_{2i}]$$

controls distribution spread



optional check  $\alpha_b, \alpha_b$

$$\beta_q = \begin{cases} (2u)^{\frac{1}{\eta+1}} & u \leq 0.5 \\ \left(\frac{1}{2(1-u)}\right)^{\frac{1}{\eta+1}} & \text{otherwise} \end{cases}$$

with  $u \sim U(0, 1)$   
 $\eta$  set to values small (2, ..., 5) or large (20, ..., 100)  
from parents (expectation)      expectation

- Amore advanced version:
- 1)  $x_{\min} = \min(P_{1i}, P_{2i})$       2)  $u_i \sim U(0, 1)$  open interval
  - $x_{\max} = \max(P_{1i}, P_{2i})$
  - 3)  $\beta_i = 1 + 2 \frac{x_{\max} - x_{\min}}{x_{\max} + x_{\min}}$
  - 4)  $\alpha_i = 2 - \beta_i^{-\eta}$

- 5) compute  $\beta_{q2} = \begin{cases} (u_i \alpha_i)^{\frac{1}{\eta+1}} & \text{if } u_i \leq \frac{1}{\alpha_i} \\ \left(\frac{1}{2 - u_i \alpha_i}\right)^{\frac{1}{\eta+1}} & \text{if } u_i > \frac{1}{\alpha_i} \end{cases}$
- 6) compute  $c_1, c_2$

7) from  $c_1$  and  $c_2$  get  $O_1, O_2$  by selecting some genes from  $c_1$  and the other from  $c_2$  and viceversa

**Number of Generation** typically in  $\frac{10-50 K}{N_{pp}}$  → # of fitness eval - sometimes computing no many times  $f$  could be prohibitive → **surrogate model** (e.g. for CFD or ML hyperparametering)

GAs don't guarantee convergence to global optimum in all cases but:

- Schema Theorem:**  $E[\# \text{ of individuals in a schema}] \propto e^{\Delta}$  if the schema has above avg fitness and C.O. and mutation are not disruptive  
 $\Delta$  = net of individuals with similar characteristics
- Building Block Hypothesis:** small low order schemas with high fitness are called building blocks  
 - GA are good at preserving and recombining building blocks
- Markov Chain Analysis:** GA can be modelled as MC where  $x_t = P_t$   
 - under certain conditions GA are stationary  $\Rightarrow$  converge to global optimum and if it is concentrated around the optimum

**HEMATIC ALGORITHM**  $\rightarrow$  combine metaheuristic with local search strategy  
 - GA (crossover and mutation) + gradient descent to optimize the offspring in a deterministic way

can be used in Travelling Salesman Problem (TSP)

- over all pairs of edges, check if swapping them the result improves  
 $(i,j), (k,e) \rightarrow (i,k), (j,e)$
- must start from a feasible solution

$\rightarrow$  until no improvements || too many iterations

**2-OPT ALGORITHM**  
 (local search operator)

$\rightarrow$  reduces complexity

**SURROGATES**

- 1) DNN  $\rightarrow$  complex non-linearities
- 2) Gaussian Process Regression  $\rightarrow$  provides uncertainty estimates
- 3) Radial Basis Function  $\rightarrow$  simple
- 4) Polynomial Regression  $\rightarrow$  fast but may lack accuracy

**procedure:**

- 1) generate  $P_0$  and evaluate with  $f$
  - 2) train surrogate model  $g$  on the evaluation
- while  $gen < T$ :
- 3) select, C.O., mutate
  - 4) evaluate using  $g$
  - 5) update  $g$  with new evals periodically

**NO FREE LUNCH THEOREM (NFL)**

No optimization algorithm outperforms the others when averaged over all possible problems

- universe optimality is impossible
- algorithm is problem dependent  $\rightarrow$  important to have domain knowledge
- sparked critical analysis among opt. algorithms to test in relevant domains

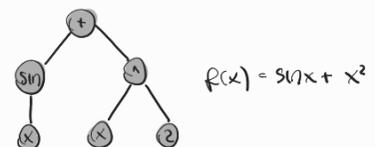
**Limitations:**

- 1) the theorem assumes  $f$  to optimize to be uniformly probable  $\rightarrow$  we do not encounter random arbitrary problems
- 2) NFL considers all opt. problems as black boxes  $\rightarrow$  in reality we have some exploitable patterns as continuity, symmetry, ...
- 3) typically we are not interested of performance among all functions

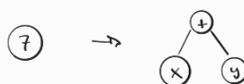
**GENETIC PROGRAMMING** evolves computer programs to solve a problem

- problems as tree structures (syntax trees)

- 1) functions in internal nodes
- 2) terminal in leaves (var/const)



**Cross over** will swap subtrees, while **mutation** randomly altering a node/subtree



Can be used for **Symbolic Regression** (find a function that approximates a Dataset)

- fitness typically MSE
- doesn't assume any model form





positive feedback loop

to remove an edge  $T_{ij} = 0$

- $\alpha$ : exploitation
- $\beta$ : exploration

$\alpha$  and  $\beta$  balances the two behaviours

- $\alpha \uparrow \Rightarrow$  more exploitation
- $\beta \uparrow \Rightarrow$  less exploration (greedy)
- $\rho \uparrow \Rightarrow$  fast forgetting, more exploration

### VARIANTS

- EAS** (elitist Ant System): gives extra push to the best-so-far solution speeding-up convergence <sup>new pheromone "elitist"</sup>
- MMAS** (max-min ant system): impose explicit upper and lower bound on pheromones to encourage continued exploration
- ACS** (Ant colony system): only allows the best so far ant to deposit pheromones

ACO encodes the state of the search in a pheromone matrix. This allows a **warmstart** of the problem (not feasible in SA, EA, PSO)

↳ non-random start

SA  $\rightarrow$  single solution based  
- probabilistic acceptance of worse solutions to escape local minima

EA is population based  
- CO, selection and mutation  
- discrete or continuous probs.

ACO  $\rightarrow$  population based  
- indirect communication via pheromones  
- can only solve discrete problems over graph  
 $\rightarrow$  particularly suited for **combinatory problems** like TSP

- $\rightarrow$  there is actually an extension on the **continuum** with probability densities
- $\rightarrow$  I could add **nettie** by storing the best solution per each ant. The social push comes from P or the best ant only
- $\rightarrow$  it's not Markov Algorithm since it is based on the history of all past nodes visited, in order to understand where to move next  $\rightarrow$  more powerful than RL

### NEURAL ARCHITECTURE SEARCH (NAS)

- using **GA**: I could select network hyperparameters  $\rightarrow$  evaluation is time consuming and encoding difficult and time consuming  
- slow convergence
- using **RL**: hyperparameters space is the environment. An agent (controller) proposes an architecture, that is trained and evaluated. The performance eg. accuracy is used as a **reward signal**  
 $\rightarrow$  may require many iterations to converge

### DARTS differentiable architecture search

- use a continuous search space and  $\nabla$  to optimize
- significantly reduces time w.r.t GA and RL
- the supernet is memory intensive and we may stuck in local optima
- it encounters also approximation challenges

- 1) Define a **supernet** that contains all possible operations in the computational graph with different edges (eg. linear vs. conv2d)  
- use different weights on edges and softmax to combine the operations on one node  $\rightarrow$  **DIFFERENTIABLE**
- 2) optimize the edges weights on **validation set**
- 3) select the architecture with operations with highest weights

### KNOWLEDGE DISTILLATION $\rightarrow$ transfer knowledge from a large **teacher** model to a smaller (more efficient) **student** model

Problems with large models 1) latency 2) resource limits (eg. edge devices) 3) memory footprint

- KD has better performance than a direct training + converge more quickly
- $\rightarrow$  student learns from imitating internal states of teacher and/or labels
- ideas from 90s and formalized in 2015

### $\rightarrow$ FOR CLASSIFICATION (original paper) [multiclass]

- teacher uses softmax as final layer  $\rightarrow$  the full distribution serves as **soft target** for the student
- $\rightarrow$  this will improve the model performance (eg. the relative probabilities of the wrong classes helps the model to learn)  
 $\rightarrow$  a car is more similar to a truck than a concert

**PROBLEM:** The CE loss will not change much since the probabilities are so close to zero

⇒ define a **temperature** parameter for the smoothness of the softmax

$$P_i = \frac{e^{z_i/T}}{\sum_j e^{z_j/T}}$$

• higher means less sharp (smooth) → more easy to learn

• In earlier implementation I have the model to match directly the **logits** (not so small, but more unstable) → almost the same if I zero mean the logits with high temp softmax

The knowledge of a teacher is passed with **distillation signals** (e.g. softmax)

• loss for classification:  $\mathcal{L} = D_{KL}(P^t \parallel P^s) = \sum P_i^t \log \frac{P_i^t}{P_i^s}$

• loss for regression:  $\mathcal{L} = \|y^t - y^s\|_2^2$  (MSE) → High Precision Loss, refined/denoised estimate of ground truth

Intermediate **feature matching**

↗ to project features with correct shapes

$$\mathcal{L} = \|f^t(x) - f_w^s(f^s(x))\|_2^2$$

↳ for a specific layer

- in addition or in substitution

**Gradient Matching**

$$\mathcal{L}_{GM} = \|\nabla_x \mathcal{L}^t - \nabla_x \mathcal{L}^s\|_2^2$$

same loss (e.g. CE)

→ reduce distance between gradients w.r.t.  $x$

→ the loss is added to the student overall loss

→ learns teacher inputs sensitivity → but requires an additional backward pass

- could suffer from vanishing gradients problem especially in deep architectures

**Adversarial Gradient Matching**

• a new module (discriminator  $D$ ) learns to distinguish between teacher and student gradients  $\nabla_x \mathcal{L}^s$  vs.  $\nabla_x \mathcal{L}^t$

The student learns to fool  $D$  and create gradients with same distribution as teacher

$$\mathcal{L} = -\log D(\nabla_x \mathcal{L}^s(x)) \rightarrow \text{to be minimized}$$

↳ = 1 when student fools  $D$

→ gradients w.r.t. intermediate layers (intermediate feature map)  $\nabla_f \mathcal{L}$

eg. here



• match how the teacher will correct the intermediate features to lower the loss

→ since copying gradient path w.r.t. the weight is infeasible

1) using the teacher checkpoints for the distillation signals → CURRICULUM LEARNING

2) preserving feature space similarity

$$\mathcal{L}_{\text{total}} = \alpha \mathcal{L}_{\text{task}} + \sum_k \lambda_k \mathcal{L}_{\text{distill}}^k$$

↳ w.r.t. HARD

↳ distill signals

eg for classification

$$\mathcal{L}_T = (1-\lambda) \mathcal{L}_{\text{task}} + \lambda \mathcal{L}_{\text{soft}}$$

**DISTILBERT** → using temperature probs and internal hidden states during training