

## PERCEPTRON LEARNING ALGORITHM

LEARN  $(\vec{x}^i, y^i)$ ,  $y \in \{-1, 1\}$  BY FINDING  $w$  ST.  $y^i = \Theta(\vec{w} \cdot \vec{x}^i)$   
WORKS ON LINEARLY SEPARABLE DATA  
USE  $\begin{cases} w_0 = -1 \\ x_0^i = 1 \end{cases} \mid \Theta(\cdot) = I(x \geq 0)$

(ASSUMING SOLUTION EXISTS)

1. INITIALISE  $\vec{w}_0 = 0$ ,  $\eta$  (LEARNING RATE)
2. REPEAT UNTIL CORRECT:
  - A. PICK A DATAPoint  $\mu^i$
  - B. IF  $\vec{w}_t$  MISCLASSIFIES  $y^i$ 

$$\vec{w}_{t+1} = \vec{w}_t + y^i \vec{x}^i$$

(ALIGNS THE PLANE'S NORM TO THE POINT)

## PROOF OF CONVERGENCE:

FIND UPPER AND LOWER BOUNDS FOR  $w^{k+1}$  AND SOLVE FOR A BOUND ON  $k$

TERMINATES IF  $\oplus$  POINTS ARE SEPARATED FROM  $\ominus$

ITERATIONS TO CONVERGE  $k \leq R^2 / \gamma^2$

### LOWER BOUND

$w_1 = 0$ , IF  $(k \geq 1)$   $x^i$  IS MISCLASSIFIED  $w^{k+1} \cdot w^k = (w^k + y^i x^i) \cdot w^k$

USING  $w^{k+1} \cdot w^k \leq \|w^{k+1}\| \cdot \|w^k\|$

(INDUCTION)  $\|w^{k+1}\| \geq k\gamma$

### UPPER BOUND

$\|w^{k+1}\|^2 = \|w^k + y^i x^i\|^2 = \|w^k\|^2 + \|x^i\|^2 + 2(w^k \cdot x^i)y^i \leq \|w^k\|^2 + \|x^i\|^2$

(INDUCTION)  $\|w^{k+1}\|^2 \leq kR^2$

### COMBINATION

$$k^2 \gamma^2 \leq \|w^{k+1}\|^2 \leq kR^2 \rightarrow k < \frac{R^2}{\gamma^2}$$

$$\gamma = \min_i [y^i (w^* \cdot x^i)]$$

← LEAST ALIGNMENT BETWEEN DATA AND POINTS

$y$  MUST BE MISCLASSIFIED

$$\leq \|w^k\|^2 + R^2$$

## SVM

### PRIMAL FORMULATION

PERCEPTRON ALG. WITH WIDEST MARGIN (STABLE WITH ERRORS)

$$\max_w \frac{2}{\|w\|} \text{ (MARGIN)} \\ \text{ST. } y_n (w^T x_n + b) \geq 1$$

$$\min_{w, b, \xi} \frac{1}{2} \|w\|^2 + C \sum_{n=1}^N \xi_n \leftarrow \text{(SOFT MARGIN)} \\ \text{ST. } y_n (w \cdot x_n + b) \geq 1 - \xi_n \\ \text{(SHALL C. SOFT MARGIN)} \quad \xi_n \geq 0$$

LEARNING: USE HINGE LOSS  $\xi \mapsto \max\{0, \xi_n\}$

CONVEX PROBLEM OBJECTIVE IS SUM OF TWO CONVEX FUNCTIONS

DERIVATION OF GRADIENT DESCENT:

$$w_{t+1} = w_t - \eta_t \nabla_w C(w_t) \quad \left( \text{USING } \lambda = \frac{1}{NC} \right) \\ C(w) = \frac{1}{N} \sum_{n=1}^N \left( \frac{\lambda}{2} \|w\|^2 + \max(0, 1 - y_n (w \cdot x_n + b)) \right)$$

### DUAL

EXPLANATION:

NUMBER OF PARAMETERS INCREASES WITH THE NUMBER OF SAMPLES

USE THE DUAL IF YOU HAVE FEW DATA-POINTS, WITH MANY FEATURES

$$\min_{\alpha} \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N y_i y_j \alpha_i \alpha_j x_i \cdot x_j - \sum_{i=1}^N \alpha_i \\ \text{ST. } \sum_{i=1}^N y_i \alpha_i = 0 \quad \forall i: 0 \leq \alpha_i \leq C$$

REPRESENTER THEOREM DERIVATION:

WRITE LAGRANGIAN OF SOFT MARGIN ( $\alpha_n, \xi_n \geq 0$ )

$$\mathcal{L}(w, b, \xi, \alpha, \gamma) = \frac{1}{2} \|w\|^2 + C \sum_{n=1}^N \xi_n - \sum_{n=1}^N \alpha_n (y_n (w \cdot x_n + b) - 1 + \xi_n) - \sum_{n=1}^N \gamma_n \xi_n$$

$$\frac{\partial \mathcal{L}}{\partial w} = w^T - \sum_{n=1}^N \alpha_n y_n x_n^T = 0 \quad \left. \begin{array}{l} \text{OPTIMAL VECTOR} \\ \text{IS IN THE SPAN OF THE DATA} \end{array} \right\} \\ w = \sum_{n=1}^N \alpha_n y_n x_n$$

$$\bullet \text{ POLY DEG } \leq d: (x \cdot y + 1)^d \\ \bullet \text{ RBF} = \exp\left(-\frac{\|x - y\|^2}{2\sigma^2}\right)$$

NON-LINEAR/KERNEL TRICK

NON-LINEARLY SEPARABLE DATA CAN FIRST BE MAPPED TO MAKE IT SEPARABLE

KERNEL FUNCTION: PRODUCES THE DOTTED RESULT WITHOUT THE INTERMEDIATE STEPS OF COMPUTATION  $k(x, y) = \phi(x) \cdot \phi(y)$

SYMMETRY ( $k(x, y) = k(y, x)$ )

! POS-DEF:  $K_{ij} = k(x^{(i)}, x^{(j)})$  IS POSITIVE SEMIDEF

ENSURES CONVEXITY

← REPRESENTS WHAT DATA THE ALGORITHM SEES

**RNN ENERGY FUNCTION (SYMMETRIC NETWORKS)**

SYMMETRIC CONNECTION MATRIX,  $J_{ij} = 0$  | SYMMETRIC  $\Rightarrow$  ENERGY F.  $\oplus$  PARALLEL IN RANDOM TIME TO LOWER E.

$E(S_1 \dots S_N) = -\frac{1}{2} \sum_{i \neq j} J_{ij} S_i S_j$  } DECREASES MONOTONICALLY TO A MINIMUM

$S_k(t+dt) = \text{SIGN} \left( \sum_j J_{kj} S_j(t) \right)$

ASYNC. DYNAMIC  $\rightarrow$  PICK  $k$  RANDOM NEURON TO UPDATE | WEIGHTS ON THE EDGES

**HOPFIELD:**  $N$  NEURONS ( $\pm 1$ )

$P$  PATTERNS TO MEMORISE:  $S_i^{\mu} = \pm 1$  (IN THE WHOLE SYSTEM)  $\rightarrow \mu$ : PATTERN NO.,  $i$ : VERTEX NO.

**HEBBIAN ( $\approx$  CORRELATION) MATRIX:**  $J_{ij} = \frac{1}{N} \sum_{\mu} S_i^{\mu} S_j^{\mu}$  (CALCULATED OVER ALL THE MEMORIES)

CONDITIONS FOR FIXED POINT ATTRACTORS:

- A CORRELATED STATE CONVERGES TO THAT COMBINATION OF STATES
- CONVERGES TO PATTERN IN ONE UPDATE IF THERE IS A FINITE OVERLAP  $\rightarrow$  **PATTERN COMPLETION**

$\rightarrow$  NOISE ANALYSIS:  $S_i(t) = \text{SIGN} \left( S_i^{\mu} + \frac{1}{N} \sum_{j \neq i} S_j^{\mu} S_j(t) \right)$

NOISE IS CAUSED BY OTHER PATTERNS INTERFERING

SYSTEM IS INITIALIZED IN A STORED PATTERN

MEAN/VARIANCE  $\rightarrow (V = \frac{1-1}{N})$  FOR NUMBER OF MEMORIES

$\rightarrow$  CAPACITY IS LINEAR WITH # NEURONS AND LIMITED BY STATES (LOCAL MIN) MERGING AND CREATING SPURIOUS MEMORIES

$\downarrow$

NEARBY/CORRELATED PATTERNS MERGE INTO A NEW MINIMUM

$\rightarrow$  NOISE ALLOWS TO CROSS BARRIERS

$\downarrow$

SIMULATED ANNEALING PRIORITISES TRANSITION TO LOWER PROBABILITIES

$P = 0.14 \cdot N \Rightarrow$  OTHERWISE THERE IS TOO MUCH INTERFERENCE

**RBM**  $\rightarrow$  ENCODING/DECODING

UL-DIVERGENCE | SYSTEM REACHES THERMAL EQUILIBRIUM IN ONE STEP  $\rightarrow$  DATA VECTORS ARE BINARY

1. LEARNING: SHOWN DATA VECTORS THAT IT MUST GENERATE WITH HIGH IP ACCURACY

2. GENERATION: FIX CONNECTION WEIGHTS, USE ENERGY, SAMPLE LOW E STATES

HIDDEN

VISIBLE

(BIPARTITE GRAPH)

$z_i = b_i + \sum_j w_{ij} S_j$  (STATE UPDATE)

**COST FUNCTION:**  $IP(S_i = 1) = \frac{1}{1 + e^{-z_i}}$  (STATE PROBABILITY)

$P(S) = \frac{1}{Z} \exp(-E(S))$  (RANDOM UPDATING  $\rightarrow$  BOLTZMANN DISTRIBUTION)

$Z = \sum_S \exp(-E(S))$  (PARTITION FUNCTION)

$E(S) = -\sum_{i,j} w_{ij} S_i S_j - \sum_i b_i S_i$

**GRADIENT DESCENT:**  $\max C = -\frac{1}{N} \sum_{\mu} \log P(x^{\mu})$  (LOG-LIKELIHOOD OVER THE DATA)

$\rightarrow$  BOLTZMANN DISTRIB.

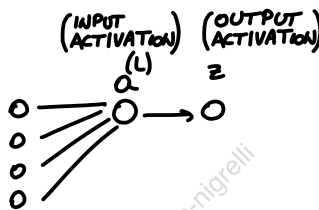
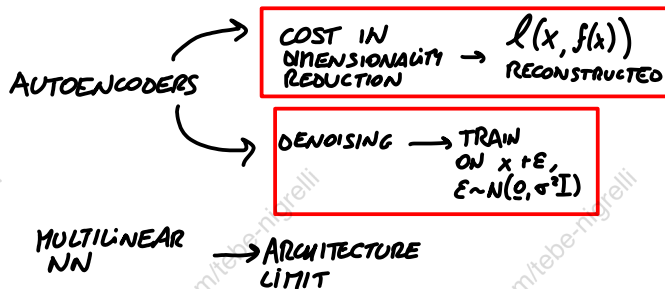
**DERIVATION OF LEARNING ALGORITHM:**  $\frac{\partial C}{\partial \theta} = \frac{1}{N} \sum_{\mu} E_h \left( \frac{\partial E(x^{\mu}, h)}{\partial \theta} \right) - E_{x,h} \left( \frac{\partial E(x, h)}{\partial \theta} \right)$

$\rightarrow$  MONTE CARLO SIMULATION/METHODS ARE USED TO APPROXIMATE THE  $E(\dots)$

**CONTRASTIVE DIVERGENCE:** REPLACE AVERAGE OF NEGATIVE TERM WITH AN ESTIMATE (GIBBS SAMPLING) IN  $k$  STEPS

$\oplus$

BATCHES OF SIZE  $k$  ARE USED TO COMPUTE EACH ESTIMATE



$$z_j^L = \sigma(a_j^L + b) = \sigma(a_j^{(L)} + w_{j0}^{(L)}) = \sigma\left(\sum_{i=1}^d w_{ji}^{(L)} z_i^{(L)} + w_{j0}^{(L)}\right)$$

SURROGATE LOSS  $\rightarrow$  (NP-HARD) 0-1 LOSS  $\downarrow$  USE A RELAXATION }  $\ell_2$  LOSS FOR CONVEXITY

**ACTIVATION FUNCTIONS**

RELU =  $\max(0, z)$  LOGISTIC =  $\frac{1}{1+e^{-z}}$

SOFT PLUS =  $\log(1+e^z)$  HYPERB. =  $\frac{e^z - e^{-z}}{e^z + e^{-z}}$

HARD THRESH =  $I(z > 0)$  TANG =  $\frac{e^z - e^{-z}}{e^z + e^{-z}}$

**KORNIK THM STATEMENT**  $\rightarrow \sigma: \mathbb{R} \rightarrow \mathbb{R}$

NONCONST BOUNDED NONDECREASING CONTINUOUS } 4

PROOF EXISTENCE } 2

NO IDEA ON ARCHITECTURE } 2

$S \subseteq \mathbb{R}^d$  CLOSED BOUNDED

$f: S \rightarrow \mathbb{R}, \forall \epsilon > 0$

$\exists$  NN WITH ONE HIDDEN LAYER, FINITE HIDDEN LAYERS

$h(x) = \sum_{j=1}^J w_j^{out} \sigma(w_j^T x + b_j)$  ST.  $|h(x) - f(x)| < \epsilon \forall x \in S$

**XOR EXAMPLE**

TRAIN NN TO CLASSIFY:  $\phi\left(\begin{matrix} x_1 \\ x_2 \end{matrix}\right) \rightarrow$

x	xor	$\phi(x)$
0 0	0	0 0 0
0 1	1	0 1 0
1 0	1	1 0 0
1 1	0	1 1 1

$\rightarrow$  USE:  $w = \begin{pmatrix} -0.2 \\ 0.5 \\ 0.5 \\ -1.5 \end{pmatrix}, x = \begin{pmatrix} 1 \\ x_1 \\ x_2 \\ x_3 \end{pmatrix}$

**BACKPROPAGATION ALGORITHM GENERAL SCHEME**

FIRST LAYER  $\rightarrow$  FILTERS EXTRACT FEATURE LEARNING  $\rightarrow$  BASIC FEATURES AT THE FIRST LAYERS

**EFFICIENT SCHEME (DERIVATION)**

- FORWARD PASS: COMPUTE AND STORE WEIGHTS  $x \mapsto f(x)$
- BACKWARD PASS: COMPUTE DERIVATIVES

BEST ESTIMATE

1. INITIALISE  $\delta_j^{(L+1)} = -(y - \hat{y}) = -(y - a^{(L+1)})$  ERROR

2. COMPUTE THE DERIVATIVES  $\frac{\partial \ell(y, \hat{y})}{\partial w_{ji}^{(L+1)}} = -(y - \hat{y}) z_j^{(L)}$

CONES FROM  $\nabla \ell^{(L+1)}$  EQUAL TO IT UP TO A NORMALISATION

3. COMPUTE  $\delta_j^{(L)} = \sigma'(a_j^{(L)}) \sum_{k=1}^d w_{kj}^{(L+1)} \delta_k^{(L+1)}$  ACCOUNTS FOR ALL INFLUENCE ON

YOU COMPUTE THE  $\Delta$  BACKWARDS, USING THE KNOWLEDGE FROM THE LAST LAYER TO GET ALL PREVIOUS ONES

$\frac{\partial \ell(y, \hat{y})}{\partial w_{ji}^{(L)}} = \delta_j^{(L)} z_i^{(L-1)}$

**SGD (FULL BATCH)**  $\rightarrow$  CONVEX:  $w_0 = \vec{0}$  ! NOT CONVEX:  $w_0 \neq \text{CONST}$  (OTHERWISE UPDATES ARE TOO CORRELATED)

$\vec{w}_{t+1} = \vec{w}_t - \eta \nabla \ell(w)$

(GD)  $w^{(k)} = w^{(k-1)} + \eta^{(k)} \mathbb{E} \left[ (y_n - w^{(k-1)} \cdot x_n) x_n \right]$  (USING SQUARE LOSS)

(SGD) SAME AS GO BUT ON SMALLER BATCH  $\rightarrow$  m BATCH SIZE  $\rightarrow$  LOWER VARIANCE VS MORE COMPUTE OF DATA

**HEURISTICS** - DECAYING STEP SIZE - UNSTABLE FOR BIG  $\eta$

**MULTI-CLASS CLASSIFICATION**  $\rightarrow$  PRESERVES SMOOTH FUNCTION MODEL CONVEXITY

(SOFT-MAX)

GENERALISES BINARY, PROVIDES "SMOOTHING"

$p(y = \ell | x, w^{(1)} \dots w^{(L)}) = \frac{\exp(w^{(L)} \cdot x)}{\sum_{i=1}^K \exp(w^{(L)} \cdot x_i)}$