# 7 Perceptron

- McCulloch-Pitts Neuron $y = \Theta(\underline{w} \cdot \underline{x} - T)$

- Feedforward/Recurrent
  (Pattern Recognition) (Short-term Memory)
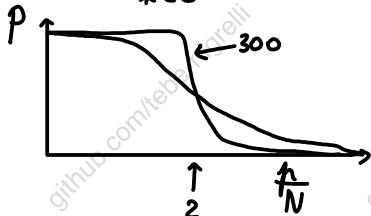
- Perceptron Learning
  ↳ Algorithm:
  For wrong classifications of $\vec{x}^\mu$
  $$\vec{w}_{t+1} = \vec{w}_t + \eta \vec{x}^\mu$$
  → Convergence in $t < \dfrac{k^2}{\epsilon^2}$  (Proof)

- Cover Formula $\begin{pmatrix} P \text{ finding } w \\ 100\% \text{ correct} \end{pmatrix}$
  $$P = \frac{1}{2^{P-1}} \sum_{k=0}^{N-1} \binom{P-1}{k}, \quad \binom{P-1}{k} = 0 \text{ if } k > P-1$$
  (Proof)



- Cover Generalisations
  → Sign constrained ⟹ Capacity $P = N$

- Cover Limitations
  → Slower Learning if data is unbiased
  → No way to quantify robustness
    ↳ Learning should be robust against noise

# 8 Support Vector Machines

- Linear Separability Assumption

- Maximise distance between points:
  Pick $w$ st.
  $$\begin{cases} w^T x_+ + b = 1 \\ w^T x_- + b = -1 \end{cases} \text{ for } x_+, x_- \text{ support vectors}$$

- Optimisation Problem:
  $$\max_w \frac{2}{\|w\|} \text{ (margin)}$$
  st. $\begin{cases} w^T x_n \geq 1 & \text{if } y_n = +1 \\ w^T x_n \leq -1 & \text{if } y_n = -1 \end{cases}$

  OR
  $$\min \|w\|^2 \text{ st. } y_n(w^T x_n + b) \geq 1 \quad \begin{pmatrix} \text{use} \\ \text{grad.} \\ \text{descent} \end{pmatrix}$$

  Tradeoff: Margin/# mistakes

- Soft Margin $\begin{pmatrix} \text{allows for} \\ \text{some mistakes} \end{pmatrix}$
  $$\min_{w,b,\xi} \frac{1}{2}\|w\|^2 + C\sum_{n=1}^N \xi_n \leftarrow \text{slack per variable}$$
  st. $y_n(\underline{w} \cdot \underline{x}_n + b) \geq 1 - \xi_n$
  $\xi_n \geq 0$

  $C$ reg. parameter: larger $C$ means more severe penalisation, so narrower margin

  OR $f(x) = \underline{w} \cdot \underline{x} + b$
  $\begin{cases} y_n f(x_n) \geq 1 - \xi_n \\ \xi_n \geq 0 \end{cases}$

  $\xi_n = \max(0, 1 - y_n f(x_n))$
  Just minimise:
  $$\min_w \frac{1}{2}\|w\|^2 + C\sum_{n=1}^N \max(0, 1 - y_n f(x_n))$$
  Regularisation term ‖ Hinge loss

  The problem is convex
  $(\sum \text{conv} = \text{conv})$

- Hinge loss
  $$\ell(t) = \max\{0, 1-t\} \quad \text{only if wrong}$$
  $t = y f(x) = y(\underline{w} \cdot \underline{x} + b)$

  Zero-one loss → cannot be differentiated, difficult to minimise

- Gradient Descent:
  $$w_{t+1} = w_t - \eta_t \nabla_w \underbrace{C(w_t)}_{g(w) \text{ to be minimised}}$$
  $$= w_t - \eta \frac{1}{N}\sum_{n=1}^N \left(\lambda \underline{w}_t + \nabla_w \ell(x_n, y_n, \underline{w}_t)\right)$$

  SGD uses mini batches instead of all the data

  Single loss over data, model

  Due to hinge loss, only points strictly violating the margin contribute to the gradient

# 9 Support Vector Machines II

**Primal SVM** ⟺ **Dual SVM**

Parameters increase with #features

Parameters increase with #data

Use if more data than features

Use if more features than data

Typical application of SVM in real life

---

Using convex duality,
$$w^* = \sum_{n=1}^{N} \alpha_n y_n x_n$$

→ **Representer Theorem:** solution to optimisation is in span of data

---

Nonlinear SVM → more flexible

Discriminant function:
$$f(x) = \begin{cases} 1 & \hat{w} \cdot \phi(x) + b > 0 \\ -1 & \text{else} \end{cases}$$

"linearity" is kept

Shortcut to apply the nonlinearity while skipping some computations and only obtaining the $f(x) \cdot f(y) = ?$ instead of eval. $f$ twice

→ ! more flexible
! less compute

**Kernel function:**
$$k : X \times X \to \mathbb{R}$$
- $k(x,y) = k(y,x)$
- ! pos-definit.

**Kernel trick:** apply non-lin to data, then apply linear SVM

⟹ ensures convexity

---

$\phi(x)$ is quicker to compute because it skips intermediate high-dim steps

$$k(x_1, x_2) = \phi(x_1) \cdot \phi(x_2)$$

want to describe ↑ rearrange in terms of $\phi$

**Kernel property is kept.**
- $\lambda \geq 0 \to \lambda k_1$
- $k_1 + k_2$
- $k_1^n$, $\prod k_i$

---

**Sample kernels**
- Polyn, deg $= d \to K(x,y) = (x \cdot y)^d$
  deg $\leq d \to (1 + x \cdot y)^d$
- Radial Basis Functions $\to \exp\left(-\frac{\|x-y\|^2}{2\sigma^2}\right)$

---

# 10 Recurrent NN

Allow for memory operations

Types of neurons → $\pm 1$, $\{0,1\}$, $1 \ldots N$

$$S_i(t+1) = \phi\left(\sum_j w_{ij} S_j(t) + I_i(t)\right)$$

Offset depends on $i, t$. $W$ matrix of weights for cross-behaviour

Parallel / asynchronous updates

Discrete / continuous time → $r_i(t+1) = \ldots$
$$\tau \frac{dr_i}{dt} = -r_i + \ldots$$

$\phi$ Activation function (sigmoid, ReLU)

Symmetric network ($J_{ij} = J_{ji}$) → Energy function (Lyapunov): $-\frac{1}{2}\sum_{j \neq i} J_{ij} S_i S_j = E(S_1 \ldots S_N)$

⟱ As "temp. decreases" AKA convergence

$$P(S_1 \ldots S_h) = \frac{1}{2}\exp(-\beta E(S_1 \ldots S_N))$$

---

# 11 Hopfield Model

**Description:** N binary ($\pm 1$) neurons,
$$S_i(t+1) = \text{sgn}\left(\sum_j J_{ij} S_j(t)\right)$$

Network connects inputs and outputs, must embed some logic in the limit state

P patterns to memorise ($\xi_i^\mu = \pm 1$)

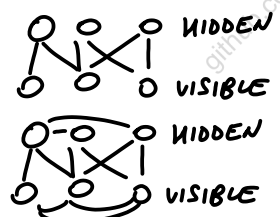**Hebbian matrix:** $J_{ij} = \frac{1}{N}\sum_\mu \xi_i^\mu \xi_j^\mu$

Noise, cross-interactions

**Pattern completion:** if there is an overlap with a stored pattern, the network converges to it in one step

---

# 12 Restricted Boltzmann Machines

- Train machines to learn the distribution of the data features

- Restricted → connections only between visible ↔ non visible
- Boltzmann machine → same type neurons are connected

Phases:
- **Learn:** visible = x, must learn to generate x with high prob.
- **Sample:** fixed weights, define energy f. natural convergence to lowest value


Hidden / Visible


Hidden / Visible