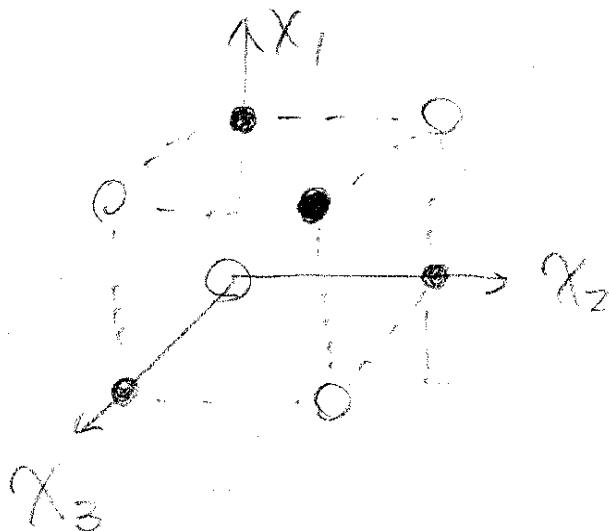


1. $0 \leq P(a) \leq 1$
 $p(\text{true}) = 1$
 $p(\text{false}) = 0$
 $p(a \text{ or } b) = p(a) + p(b) - p(a \text{ and } b)$
2. $p(m | l) = (P(l | m) * p(m)) / p(l)$
 $p(l | m) = (p(m | l) p(l)) / p(m)$
 $= .99 * 1 / (10 \text{ to the } 7) / 1(30)$
 $= .297 \times 10^{-6}$
3. Map-maximum a posteriori an approximation of Bayesian prediction where predictions are based on most probable hypothesis. That is an h_i that maximizes $p(h_i | \text{data}) * p(h_i)$
- 4.
5. $p(\text{hot and sunny}) = .4 + .1 = .5$
6. $x_1 \text{ xor } x_2 \text{ xor } x_3$
make truth table, and graph the box that is not separable by a plane
7. $w_0 = (n - .5)$
 $w_i = 1.0, 1 \leq i \leq n$
use sign function
 $h(x) = \text{sgn}(\sum w_i x_i \text{ from } i = 0 \text{ to } n - (n - .5))$ $h(x) = \text{output}$
8. $\delta_j = g'(in_j) \sum(w_{ji} \delta_i)$
 $W_{k,j} \leftarrow W_{k,j} + \alpha \times a_k \times \delta_j$
back propagation is process of adjusting the weights of the hidden layers based on how much they contributed to the final output error.
9. Mercer's theorem states any kernel of a positive definite linear operator corresponds to an inner product in some feature space. A kernel function is a function of the form $K(x_i, x_j) := F(x_i) \cdot F(x_j)$

(6) $X_1 \text{ XOR } X_2 \text{ XOR } X_3$

X_1	X_2	X_3	$X_1 \text{ XOR } X_2 \text{ XOR } X_3$
0	0	0	0
0	0	1	1
0	1	0	1
0	1	1	0
1	0	0	1
1	0	1	0
1	1	0	0
1	1	1	1



No plane can separate group of
On's and Off's

③

MAP - maximum a posteriori

- an approximation of Bayesian prediction where predictions are based on the most probable hypothesis that is h_i that maximizes $P(d|h_i)P(h_i)$

so

$$h_{MAP} = h_i \text{ s.t. } P(d|h_i)P(h_i) \text{ is largest}$$

④

Ockham's Razor - The simplest solution is the best

using h_{MAP} - maximizing

$$P(d|h_i)P(h_i) \text{ is equivalent to}$$

$$\text{minimizing } \frac{1}{P(d|h_i)P(h_i)}$$

which is equivalent to minimizing

$$\log_2 \left(\frac{1}{P(d|h_i)P(h_i)} \right) = -\log_2 P(d|h_i)P(h_i)$$

$$= -\log_2 P(d|h_i) - \log_2 P(h_i)$$

= # of bits needed to describe data d with hypothesis h_i + # of bits to describe h_i

so MAP is equivalent to choosing shortest string to explain data