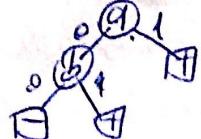
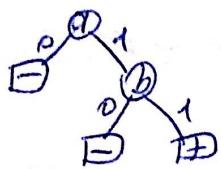


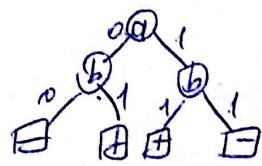
$$\text{d} f(A, B) = A \vee B$$



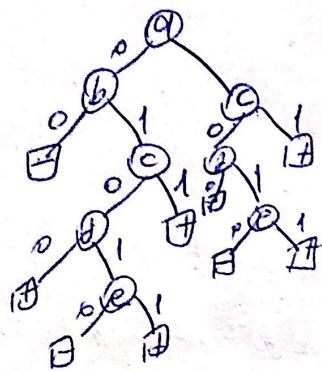
$$f(A, B) = A \wedge B$$



$$f(A, B) = A \oplus B$$



$$f(A, B, C, D, E) = (A \vee B) \wedge (C \vee \neg D \vee E)$$



d. TDIDT algo ?

$$\text{Gain}(S, \text{outlet}) = H\left(\frac{6}{11}, \frac{5}{11}\right) - \left(\underbrace{\frac{6}{11}H\left(\frac{2}{3}, \frac{2}{3}\right)}_{\text{outlet}} + \underbrace{\frac{5}{11}H\left(\frac{1}{4}, \frac{3}{4}\right)}_{\text{non-outlet}} \right)$$

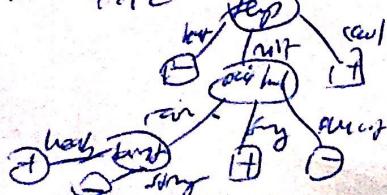
$$\text{Gain}(S, \text{outlet}) = 0.629$$

$$\text{Gain}(S, \text{work}) = 0.028$$

calculator gain for temp & mH?

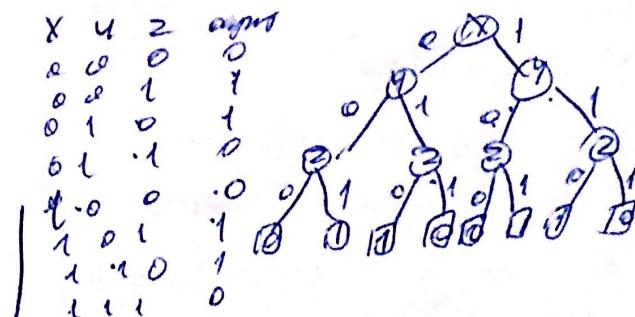
$$\text{Gain}(S, \text{outlet}) = H\left(\frac{2}{3}, \frac{1}{3}\right) - \left[\frac{1}{3}H\left(\frac{1}{2}, \frac{1}{2}\right) + \frac{2}{3}H\left(\frac{1}{2}, \frac{1}{2}\right) \right] = 0.87$$

$$\text{Gain}(S, \text{work}) = -0.012$$



d. why TDIDT algo does not produce smallest decision tree, binary trees?

X	A	B	outlet
0	0	0	0
0	0	1	1
0	1	0	1
0	1	1	0
1	0	0	0
1	0	1	0
1	1	0	1
1	1	1	0



d. X nodes and depth n wrt $p_1 \rightarrow p_n$ wrt $\sum p_i = 1$.
Pure entropy $H(X) = H(p_1 - p_n) = \sum p_i \log p_i \leq \log n$

By def, $\log_2 p = 0$ for $p=0$

function $2 \mapsto \log_2 X$ is concave as second derivative is -ve $\Rightarrow E(f(X)) \leq f(E(X))$

$H(x) = E_p(x) \cdot [-\log_2 P(x)]$ where $P(x)$ is prob. mass function.

$$= \sum_{i=1}^n (-\log_2 p_i) \cdot p_i$$

$$= \sum_{i=1}^n (-\log_2 P_X(v_i)) \cdot P_X(v_i)$$

$$= E_{P_X}[-\log_2 P_X(v_i)]$$

By JT & CS Theory

$$= E_{P_X} \left[\log_2 \frac{1}{P_X(v_i)} \right] \leq \log_2 \left(E_{P_X} \left[\frac{1}{P_X(v_i)} \right] \right)$$

$$= \log_2 \sum_{i=1}^n P_X(v_i) \cdot \frac{1}{P_X(v_i)}$$

$$= \log_2 n$$

d. pure information gain \Rightarrow always non-negative?

clearly $\text{Gain}(X, Y) \geq 0 \Leftrightarrow -\text{Gain}(X, Y) \leq 0$

$$\text{Gain}(X, Y) = \sum_i P(X=i) \cdot \log P(X=i)$$

$$- \sum_v P(Y=v) \cdot \sum_i P(X=i | Y=v) \cdot \log P(X=i | Y=v)$$

$$= \sum_i \sum_v P(X=i | Y=v) \cdot \log P(X=i | Y=v)$$

$$- \sum_i \sum_v P(X=i | Y=v) \cdot P(Y=v) \cdot \log P(X=i | Y=v)$$

$$= \sum_i \sum_v P(X=i | Y=v) \cdot P(Y=v) \left(\log \frac{P(X=i | Y=v)}{P(X=i | Y=v)} \right)$$

By JT & CS Theory

$$\leq \sum_v P(Y=v) \cdot \log \frac{\sum_i P(X=i | Y=v) \cdot P(X=i)}{P(X=i | Y=v)}$$

$$= 0$$

11

$\text{d) } h_1 \wedge h_2 \text{ iff } (\forall x \in \mathcal{X}) [h_2(x) = 1 \Rightarrow h_1(x) = 1]$
 Show \leq is paralleler ordered basis over \mathcal{X} ?
 \rightarrow reflexivity: $\forall h \in H = h \leq h$

Let $h \in H, x \in \mathcal{X}$ with $h(x) = 1 \Rightarrow h_1(x) = 1$
 $h \leq h_1$, it is reflexive

\rightarrow antisymmetry $h_1, h_2 \in H$; if $h_1 \leq h_2$ and $h_2 \leq h_1$
 let $h_1(x) = 1 \wedge h_2(x) = 1$

$\Rightarrow h_2(x) = 1 \Rightarrow h_1(x) = 1$

$\Rightarrow h_1(x) = 1 \Rightarrow h_2(x) = 1$

$\Rightarrow h_1(x) = 1 \wedge h_2(x) = 1$ and $h_1 = h_2$ and unique

\rightarrow transitivity of $h_1 \leq h_2$ and $h_2 \leq h_3$

Let $h_1, h_2, h_3 \in H, x \in \mathcal{X}$ with $h_1 \leq h_2$ and $h_2 \leq h_3$

if $h_2(x) = 1 \Rightarrow h_1(x) = 1$

if $h_3(x) = 1 \Rightarrow h_2(x) = 1$

$\Rightarrow h_3(x) = 1 \Rightarrow h_1(x) = 1$ and unique

$\text{d) } \mathcal{X} \text{ be } \{0,1\}^n \text{ with } n \in \mathbb{N}, \text{ then do } x_i \rightarrow x_i$
 for $h \in H = \bigvee_{l \in L} \text{disj. } L \subseteq \{x_1 \rightarrow x_1, x_2 \rightarrow x_2, \dots\}$

Reparatur: take input E & consider $\mathcal{E} \subseteq \mathcal{X}$ and
 what happens when we do $H(E)$?

Let E^+ front set of edges $c(\vec{x}_i) = 1$
 E^- back set of edges $c(\vec{x}_i) = 0$

Input $E \rightarrow \{x_1 \rightarrow x_2\}$ ~~is~~

$e_i = (\vec{x}_i, c(\vec{x}_i))$ with $x_i \in \{0,1\}$

$h = \vec{x}_1 \vee \vec{x}_2 \vee \dots \vee \vec{x}_n \vee \vec{x}_1 \wedge \dots \wedge \vec{x}_n$
 for $e = (\vec{x}_i, c(\vec{x}_i)) \in E^-$ \rightarrow ~~it goes away~~

$\text{if } h(\vec{x}) = 1$
 remove all ~~front~~ ~~back~~ edges in \mathcal{X} from h

for $e = (\vec{x}_i, c(\vec{x}_i)) \in E^+$ \rightarrow ~~goes away~~
 $\text{if } h(\vec{x}) > 0$
 then $\vec{x} \in E^+$

$\text{d) } \mathcal{X} \text{ be } \{0,1\}^6, H_0 \text{ best captures over } x_1 \text{ to } x_6$
 for all $h \in H$, $h = \bigwedge_{l \in L} (x_1 \rightarrow x_2 \wedge x_3 \rightarrow x_4 \wedge \dots \wedge x_5 \rightarrow x_6)$
 $E = E^+ \cup E^-$ with $E^+ = \{(1,1,1,1,1,1), (1,1,0,1,1,1), (1,1,0,1,0,0)\}$, $E^- = \{(0,0,0,1,1,1,0)\}$

d) Using cond.-elimination \rightarrow S contains V_{true} ?

$G_0 = \{\top\}$

$S_0 = \{x_1 \wedge x_2 \wedge \dots \wedge x_6 \rightarrow x_6\}$

$E_1 = \{(1,1,1,1,1,1)\}$

$G_1 = \{\top\}$

$S_1 = \{x_1 \wedge x_2 \wedge \dots \wedge x_6 \rightarrow x_6\}$

$E_2 = \{(1,1,0,1,1,1)\}$

$G_2 = \{\top\}$

$S_2 = \{x_1 \wedge x_2 \wedge x_3 \wedge x_5 \wedge x_6\}$

$E_3 = \{(1,1,0,1,0,0)\}$

$G_3 = \{\top\}$

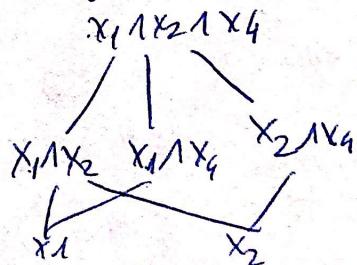
$S_3 = \{x_1 \wedge x_2 \wedge x_4\}$

$E_4 = \{(0,0,0,1,1,0)\}$

$G_4 = \{x_1, x_2\}$

d) Generate all hyperedges \rightarrow V_{true} ?

$V_{\text{true}} = \{h \in H \mid \exists g \in G, S \in S : g \subseteq h \subseteq S\}$



$V_{\text{true}} = \{x_1 \wedge x_2, x_1 \wedge x_4, x_1 \wedge x_2 \wedge x_4, x_2 \wedge x_4, x_1 \wedge x_2 \wedge x_4\}$

\mathcal{L} Pre $H \subseteq 2^X > X, E^+, E^- \subseteq X, \subset H$,
can do also \exists same, \mathcal{E} not generalizable
 $\mathcal{L}(S, E, E^-)$?

Let $E = \{e_1 \rightarrow e_m\}$

$E_k \Rightarrow \{e_1 \rightarrow e_k\} \rightarrow$ elements after step

S_k, G_k denote boundaries after seen E_k
 S_k^*, G_k^* denote actual boundaries.

Show $S_k \supseteq S_k^*$, $G_k \supseteq G_k^*$ for $k \leq m$

observation = $1 \leq k \leq m-1$, if h is compatible E_{k+1} ,
it is consistent with E_k .

Notation on k

if $k=0$, $E_0 = \emptyset$, normally holds as
 G_0, S_0 are initialized.

if $k=1$, $G_k = G_k^*$, $S_k = S_k^*$

assume E_{k+1} is positive,

$G_{k+1}^* = G_{k+1}$,

claim = $G_{k+1}^* \subseteq G_{k+1}$

Let $g \in G_{k+1}^*$, by def of general boundary.

g^* is consistent with E_{k+1} ,

any $g \in g^*$ is inconsistent with E_{k+1} ,

E_{k+1} is positive so any $g \in g^*$ is inconsistent with E_k .
by observation g^* is consistent with E_k

then $g^* \in G_k^*$.

by induction $G_k^* = G_k$, then $g^* \in G_k$

$g^* \in G_{k+1}^*$ covers E_{k+1} then $g^* \in G_{k+1}$

claim = $G_{k+1} \subseteq G_{k+1}^*$

Let $g \in G_{k+1}$, by induction E_{k+1} is positive

then $g \in G_k = G_k^*$ and

g is consistent with E_{k+1} ,

any $g' \in g$ is inconsistent with E_k

by observation g' is inconsistent with E_{k+1}

then $g \in G_{k+1}^*$

It implies $G_{k+1} \subseteq G_{k+1}^*$

\mathcal{L} Pre or the theory of cognitive errors done
at most n prediction mistakes?

Let have $x_1 \dots x_n$, set E of m examples
and so H also, start with $h = x_1 \rightarrow x_2 \wedge x_3 \rightarrow x_4$

if $m \leq n$, the claim is true

if $m > n$,

if first prediction mistake is made, then delete
all wrong literals in example from h
so only n literals remain

for each wrongly classified example, where is
at least one literal gets deleted from h
so before at least one literal that
there may be at most n more mistakes
then in total at most $n+1$ mistakes.

\mathcal{L} for E and H , if draw $m \geq \frac{1}{\delta} (\ln |H| + \ln \frac{1}{\delta})$ samples

and consistent hypo from H ; can we get $\Pr[\text{error}(h)] \leq \delta$

claim = if draw $m \geq \frac{1}{\delta} (\ln |H| + \ln \frac{1}{\delta})$ samples,
we can find consistent hypothesis,
we can PTC len C has error by ~~at most δ~~
with $\Pr[\text{error}(h)] \leq \delta \geq 1 - \delta$

$\Rightarrow \Pr[\text{error}(h) \geq \epsilon] \leq \delta$

let $H_B = \{h_{i,j} \rightarrow h_k\} \subseteq H$. B graph hypo
 h_i with $\text{error}(h_i) \geq \epsilon$ and it can be output
if it is consistent with all examples

Let $E = \{e_1 \rightarrow e_m\}$ some sequence,
 $h \in H_B$

then $\Pr[h \text{ is consistent with } E] \leq 1 - \delta$

$\Rightarrow \Pr[h \text{ is consistent}] \leq (1 - \delta)^m$

$\Pr[\neg(h \in H_B \text{ consistent})] \leq |H_B| \cdot (1 - \delta)^m \leq |H| \cdot (1 - \delta)^m$
 $(1 - \delta) \leq e^{-\delta}$ for $0 \leq \delta \leq 1$

$\Rightarrow |H| \cdot e^{-\delta m} \leq \delta$

probability of bad hypothesis being consistent with m
examples are bounded by $|H| \cdot e^{-\delta m}$ and
should be about δ

$\Rightarrow -\delta m \leq \ln(\frac{\delta}{|H|})$

$m \geq \frac{1}{\delta} (\ln |H| - \ln \delta)$

$m \geq \frac{1}{\delta} (\ln |H| + \ln \frac{1}{\delta})$

101

of size k -CNP_n set of all CNFs of each clause contains at most k literals from $\{x_1, \bar{x}_1; \dots, x_n, \bar{x}_n\}$

Then k -CNP_n is efficiently PAC-learnable?

claim = it is PAC-learnable

Basis |H|, each clause has at most k literals and total there are $2n$ literals, $\binom{2n}{k}$ clauses with k literals.

total num of different clauses are bounded by

$$\sum_{i=1}^k \binom{2n}{i} \leq \binom{2n}{n}, \frac{2n-(k-1)}{2n-(2k-1)}$$

$$\leq (2n)^k \cdot \frac{2n-(k-1)}{2n-(2k-1)} = \text{poly}(n)$$

each clause is part of certain k-CNF or not,
total number of k-CNF is bounded by $|H| \leq 2^{\text{poly}(n)}$
And consider hypothesis

$$h_i = \frac{1}{\sum} (\ln |H| + \ln \frac{1}{\sum})$$

$$= \frac{1}{\sum} (\ln 2^{\text{poly}(n)} + \ln \frac{1}{\sum})$$

$$= \frac{1}{\sum} (\text{poly}(n) \ln 2 + \ln \frac{1}{\sum}) \text{ examples}$$

thus $\Delta \text{poly}(n) \ln \frac{1}{\sum} / \ln 2$.

so there are $\leq \text{poly}(n)$ different clauses.

~~General hypothesis~~, let broken variables x'_1, \dots, x'_m where $m \leq \text{poly}(n)$ is num of different clauses.

$x'_i \neq 1$ iff clause is satisfied

any k-CNF formula over x_1, \dots, x_n is equivalent to

conjunction over $x_1 \wedge \dots \wedge x_m$,

consensus hypothesis can be found by FALS also.
by applying also for cognitive concepts,

$\delta \gamma = R^2$, compare Colv λ to set of any parallel aligned rectangles, Is Colv λ learnable?

R	LR ₂	z
R ₁	+	R ₃
+	+	R ₄
R ₅	+	

(if first consensus hypothesis), take smaller rectangle contains all positive examples called R_E

area of possible misclassification $\rightarrow R \setminus R_E$

probability example drawn from $R \setminus R_E$ is error(R_E)
Guarantee $\Pr[\text{error}(R_E) \leq \epsilon] \geq 1 - \delta$

If $P(R) \leq \epsilon$, then $P(R \setminus R_E) \leq P(R) \leq \epsilon$

If $P(R) > \epsilon$, then consider subregions $R_1 \cup R_4 \cup R \setminus R_E$
s.t. $P(R_i) = \frac{\epsilon}{4}$

choose enough examples with higher prob. prob
there is a point in each R_i , then R_E misclassifies
all regions and error $\geq \epsilon$ atmost ϵ .

$$R \setminus R_E \subseteq \bigcup_{i=1}^4 R_i$$

$$P(R \setminus R_E) \leq P\left(\bigcup_{i=1}^4 R_i\right) \leq \sum_{i=1}^4 P(R_i)$$

$$\leq \sum_{i=1}^4 \frac{\epsilon}{4}$$

$$\leq \epsilon$$

find lower bound,

$$P(R_i) \text{ is prob}, i=1 \text{ to } 4$$

$$\Leftrightarrow P\left(\bigcup_{i=1}^4 R_i\right) \text{ is prob}$$

$$\leq \sum_{i=1}^4 P(R_i) \text{ prob}$$

$$= \sum_{i=1}^4 (1 - P(R_i))^m$$

$$> \sum_{i=1}^4 (1 - \frac{\epsilon}{4})^m$$

$$= \sum_{i=1}^4 4 \cdot (1 - \frac{\epsilon}{4})^m$$

$$\leq 4 \cdot e^{-\frac{\epsilon m}{4}} \quad \epsilon \leq \delta$$

$$\text{where } 1 - \frac{\epsilon}{4} \leq e^{-\frac{\epsilon}{4}}$$

$$\Leftrightarrow \frac{4}{\epsilon} m \ln 4 \leq \ln \delta$$

$$m \geq \frac{4}{\epsilon} \ln \left(\frac{4}{\delta} \right)$$

If median at least m examples, we output
hypothesis with error atmost ϵ ,

with probability ~~at least~~ at least $1 - \delta$

α Conv X , $C \subseteq \mathbb{R}^n$, C shades $\subseteq X$ for $A \subseteq S$,
there is $C \subseteq C$, C covers every $A \in \text{VC}(C)$ since layers
from subset of X can be shaded by C .

a) Circles in \mathbb{R}^2 ?

$$\text{claim} = \text{VC}(C) = 3$$

case 0: \rightarrow draw 3 points, set of points is shaded by C
why not 4?
case 1: some points x in interior of convex hull of
others, we can't use circle to cover exactly
points in convex hull

case 2: every point in convex hull, we have 2 points such
that any circle contains A also contains B or
any circle contains B also contains points in A .

c) convex layers?

$$\text{claim} = \text{VC}(C) = 2d+1$$

take $2d+1$ points on circle that can be shaded by
d half spaces.

induction $d=1$, 3 points by 1 half space ~~is~~

$d \geq 3$; intersection of these halfspaces is convex-gms.
why not $2d+2$?

case 1: some points are in interior of convex hull
of others, we can't separate convex hull from
these points with any convex shape.

case 2: all points in convex-hull, subset contains
every other points in convex hull, ~~we~~
we can't separate subset from its complement
with d-gms because we need different d+1 half

b) by C it is $\cdot 7$ for $\text{VC}(\text{halfspaces})$ spaces

α a) $|C| \leq n$ show $\text{VCdm}(C) \leq \log_2 |C|$?

By def, $\text{VCdm}(C)$ is largest subset of
 C is shaded by C

then $\forall B \subseteq A: \exists C \subseteq C: C \cap A = B$

$$2 \text{VCdm}(C) = 2^{|A|} \leq |C|$$

$$\Leftrightarrow \text{VCdm}(C) \leq \log_2 |C|$$

b) $C_1, C_2 \subseteq 2^X$ is ranapocher over X ,
show $\text{VCdm}(C_1) \leq \text{VCdm}(C_2)$?

A, B layer ~~sets~~ set shaded by C_1 and for each
 $C \in C_1 = C \subseteq C_2$ where $C_1 \subseteq C_2$,
 $\text{VCdm}(C_1) = |A|$, then A is also shaded by C_2
then $\text{VCdm}(C_2) \geq |A| = \text{VCdm}(C_1)$

c) $C \subseteq 2^X$ over $X: \bar{C} = \{X \setminus C : C \subseteq C\}$,
show $\text{VCdm}(C) = \text{VCdm}(\bar{C})$?

Yes, let $A \subseteq X$ st. $\text{VCdm}(C) = |A|$,
 $\text{claim}_2: A$ is also shaded by \bar{C}
 A is shaded by C ,
foreach $B \subseteq A$ where $\exists C \in C$ st. $C \cap A = B$
then ~~$A \setminus C \neq \emptyset$~~

$$\frac{X \setminus C \in \bar{C}}{A \cap (X \setminus C)} = A \setminus C = A \setminus (A \cap C) = A \setminus (A \setminus B) = B$$

so A is shaded by \bar{C}

$$\text{then } \text{VCdm}(\bar{C}) \geq \text{VCdm}(C)$$

α borden $f = \{0, 1\}^n \rightarrow \{0, 1\}$ is symmetric if
 $f(x_1, \dots, x_n) = f(x_{\pi(1)}, \dots, x_{\pi(n)})$ for $\pi \in S_{n, 1}^n$,
can \exists 1 ishow C_n is poly. PAC-learnable?
 $\text{claim} = \text{VCdm}(C) > n+1$

separating $A \subseteq \{0, 1\}^n$ st. only after y of A is
given by symmetric func $\forall x \in X (x \in A \Leftrightarrow f_y(x) = 1)$
if it satisfies then elements of A must be transversal
st. $\exists x \in X$ and have $5 \neq 5'$

consider equivalence classes $S \sim S'$ iff s and s' one
permutation of each other, then any A there is
shaded by C must contain ~~at least~~ at most one
representative of each equivalence class,

there are $n!$ such classes, then largest A is $n!$
then $\text{VCdm}(C) = n! = \text{poly}(n)$
we can find $h \in C$ consistent with S in linear manner
 $S \subseteq X$ and $C \subseteq C$

then $\exists h \in C$ is poly. PAC-learnable.

α Show C consisting of all axis parallel rectangle
 β poly. PAC learnable by using Corollary 1?

Consistent hypothesis can be found in poly. time by
 partitioning smaller rectangle containing all positive
 $\text{cl}_{\text{dm}} = V(\text{dm}(C)) = 4$ examples.

Consider all subsets can be refined by $\square \rightarrow \square \times \square$

some axis-parallel rectangle such
 set of A with 5 points. If 3 others are $\oplus \ominus \oplus$
 on line, this cannot be shattered.

Assume none are collinear, let $x_{\max}, x_{\min}, y_{\max}, y_{\min}$
 denote max/min x/y-coordinates of points in A ,
 consider bounding box $BB(A)$ is given by points

$(x_{\max}, y_{\max}), (x_{\max}, y_{\min}), (x_{\min}, y_{\max}), (x_{\min}, y_{\min})$.
 if at least one point lies inside $BB(A)$, then it cannot

if all lie on bounding box, then there exists
 labeling that 4-positive and 1-negative that
 cannot be refined by C

then $V(\text{dm}(C)) \leq 5$ so $V(\text{dm}(C)) = 4$.

α Consider test results giving 98% false alarm rate,
 correct rejection rate 97% to detect real cases,
 0.008 of entire pop. have cancer

$$P(C|+) = \frac{P(+|C) \cdot P(C)}{P(+)} = \frac{0.0038}{P(+)} = 0.0038$$

$$P(TC|+) = \frac{P(+|TC) \cdot P(TC)}{P(+)} = \frac{0.0028}{P(+)} = 0.0028$$

$$\text{we know } P(C|+) + P(TC|+) = 1$$

$$P(C|+) = \frac{0.0038}{0.0038 + 0.0028} = 0.21$$

Suppose we decide to take second test for some reason,
 suppose second test response as well. What are
 posterior prob. of cancer and no cancer following these
 independent tests?

prob. of two tests given

$$\text{prior has been } P(+, +|C) = P(+|C) \cdot P(+|C)$$

$$P(+|C) = 0.0038, P(C) = 0.0008$$

$$P(-|TC) = 0.97, P(TC) = 0.0022$$

$$P(+|TC) = 0.03 \text{ and } P(C|+) = 0.001$$

$$P(C|+ \cap +) = \frac{P(+|C) \cdot P(C)}{P(+|C) + P(+|TC) \cdot P(TC)} = \frac{0.0038 \cdot 0.0008}{0.0038 + 0.0022} = 0.0007$$

$$P(+, +|C) = P(C) + P(+|TC) \cdot P(TC) = 0.0007$$

α unfair coin, estimate sequence of n consecutive unknown
 prob. θ of head. What is the likelihood estimator $\hat{\theta}$ of
 θ from sequences of n coin flips?

We observe n outcomes of heads during n tosses,
 then prob. of seeing T is $L(\theta) = \theta^h \cdot (1-\theta)^{n-h}$

n heads $n-h$ tails

(log likelihood)

$$L(\theta) = h \cdot \log \theta + (n-h) \cdot \log(1-\theta)$$

so maximise L , take first derivative and set it to 0

$$\frac{d L(\theta)}{d \theta} = \frac{h}{\theta} + \frac{(n-h)}{1-\theta} = 0$$

$$\Leftrightarrow \frac{h}{\theta} = \frac{n-h}{1-\theta}$$

$$\Leftrightarrow \theta = \frac{h}{n} \text{ so } \hat{\theta} = \frac{h}{n} \text{ maximizes likelihood}$$

α $\mathcal{L} \subseteq \{01\}^n \times \{01\}$ of having exactly m 0's = $\{(1,1,1,1,1),$
 $(1,1,0,1,0), (0,1,1,0,0), (0,1,0,1,0)\}$

$X = (1,1,1,0)$ and $H = \{h_1, h_2, h_3\}$ with

if $a_1=1$ then h_1 returns $\frac{1}{4}$ and $\frac{1}{4}$ if $a_2=1$ and $\frac{1}{4}$ if $a_3=1$

if $a_2=1$, h_2 returns $\frac{1}{4}$ with $\frac{1}{4}$ and $\frac{1}{4}$ if $a_1=1$ and $\frac{1}{4}$ if $a_3=1$

if $a_3=1$ then h_3 returns $\frac{1}{4}$ with $\frac{1}{2}$, and $\frac{1}{2}$ if $a_1=1$ and $\frac{1}{2}$ if $a_2=1$

assume hypotheses h_1, h_2, h_3 are uniformly functioning.

a) What is resp. hypo. against $P(h|D)$?

$$H_{\text{hyp}} = \arg \max_{h \in H} \frac{P(D|h) \cdot P(h)}{P(D)}$$

$$= \arg \max_{h \in H} P(D|h) \cdot P(h)$$

$$\text{by def } P(h_1) = \frac{1}{3}, P(h_2) = \frac{1}{3}, P(h_3) = \frac{1}{3}$$

$$= \arg \max_{h \in H} P(D|h)$$

$$h_1 \geq P(D|h_1) = \frac{4}{5} \cdot \frac{1}{5} \cdot \frac{4}{5} \cdot \frac{4}{5} = 0.11024$$

$$h_2 \geq P(D|h_2) = \frac{4}{5} \cdot \frac{4}{5} \cdot \frac{1}{5} \cdot \frac{4}{5} = 0.11024$$

$$h_3 \geq P(D|h_3) = \frac{2}{3} \cdot \frac{2}{3} \cdot \frac{2}{3} \cdot \frac{2}{3} = 0.11025 \text{ always}$$



Q) What's bayes optimal classification of x ?

approx $\sum_{h \in H} P(h_j | h_i) \cdot P(h_i | D)$, we denote $P(h_i | D)$

$$P(h_1 | D) = \frac{0.1024}{0.1024 + 0.0256 + 0.9752} \approx 0.25$$

$$P(h_2 | D) = 0.0256$$

$$P(h_3 | D) = \frac{0.9752}{\dots} \approx 0.99$$

$$X = (h_1, h_2)$$

$$\text{for } D' \Rightarrow \sum_{h \in H} P(h_1 | h_i) \cdot P(h_i | D)$$

$$= P(O|h_1) \cdot P(h_1 | D) + P(O|h_2) \cdot P(h_2 | D) + P(O|h_3) \cdot P(h_3 | D)$$

$$= 0.2848 + \frac{1}{3} \cdot 0.0256 + \frac{2}{3} \cdot 0.9752$$

$$\text{For } D' \Rightarrow \frac{4}{3} \cdot 0.0256 + \frac{1}{3} \cdot 0.9752 \approx 0.32713$$

c) What's naive bayes classifier of v ?

$$V_{10} = \text{approx } P(X_j) \cdot \prod_{i \neq j} P(X_i | V_{10})$$

$$\Rightarrow P(V=1) \cdot \prod_{i=1}^3 P(X_i | V=1) = \frac{1}{4} \cdot 1 \cdot 1 \cdot 0.20$$

$$P(V=0) \cdot \prod_{i=1}^3 P(X_i | V=0) = \frac{3}{4} \cdot \frac{1}{3} \cdot \frac{1}{3} \cdot \frac{1}{3} = \frac{1}{36}$$

naive bayes classifier is '0'

d) $P \subseteq D^2$, $|P| = n$ (each all contours are perfect),

\emptyset be overall random & half $\emptyset(P_{(i)}) \subset \emptyset(P_{(j)})$ for all

a) $n \geq 3$ exist $P \subset R^2$ of n non overlapping

sq. all cells of Voronoi diagram of P are unbounded?

take n points in circle,

Voronoi diagrams edges correspond to perpendicular bisectors of segments of neighbour. points

then no bounded cell in diagram

b) $P \subset R^2$ with $|P| = n$ for $n \geq 3, 5$ how many vertices of Voronoi diagram of P is atmost $2n-5$ and no of AB edges is atmost $3n-6$ where points in plane B vertex of Voronoi diagram, point in plane lies on edge?

Contours $G = (V(G), E(G))$ as $V(G)$ has all vertices of diagram $\emptyset + X$, $E(G)$ has all edges + bounded edges
 V denote num of vertices, e denote num of edges
 observations (G is planar, $|V(G)| = v+1$, $|E(G)| = e$,
 num of faces of $G = n$ (each face corresponds to Voronoi cell)
 any vertex of Voronoi diagram must atleast 3 edges

$$|\delta(v)| \geq 3 \Rightarrow V_{\text{vert}}(G)$$

$$2|E(G)| = \sum_{v \in V(G)} |\delta(v)| \geq 3|V(G)|$$

$$\Leftrightarrow v+1 \leq \frac{2}{3}e$$

$$\text{from euler, } |V(G)| - |E(G)| + f = 2$$

$$(v+1) - (e) + n = 2$$

$$2e - e + n \geq 2$$

$$e \leq 3n - 6$$

$$\text{and } (v+1) - (e) + n = 2$$

$$v+1 - (3n-6) + n \leq 2$$

$$v \leq 3n-5$$

d) 100 rand. examples, PR classifier 87 correctly,
 a) what's 95% confidence interval for true error of it?

$$\text{errors}(t) \pm 2.082 \cdot \sqrt{\frac{\text{errors}(t) \cdot (1 - \text{errors}(t))}{n}}$$

$$= 0.13 \pm 1.96 \cdot \sqrt{\frac{0.13 \cdot 0.87}{100}}$$

interval is $[0.0641, 0.1959]$

b) What's the meaning of this confidence interval?

with at least 95% confidence probability -
 this interval contains value of true error.

d) Define equivalent bin F_B-size and

$$F_d = \frac{1}{\alpha \cdot \frac{1}{P} + (1-\alpha) \cdot \frac{1}{R}} \quad \text{for details}$$

$$\text{above } F_B = (1+\beta^2) \frac{P \cdot R}{\beta^2 \cdot P + R} ?$$

~~$F_d = F_B$~~

observation for a weighted mean & harmonic mean of P and R ; $2 \cdot \frac{P \cdot R}{P+R}$

claim = if we choose $\alpha = \frac{1}{\beta^2 + 1}$, then $F_d = F_B$

$$F_d = \frac{1}{\frac{1}{\beta^2} \cdot \frac{1}{P} + (1 - \frac{1}{\beta^2}) \cdot \frac{1}{R}}$$

$$= \frac{1}{1 + \beta^2} \left(\frac{1}{P} + \frac{\beta^2}{R} \right)$$

$$= (1 + \beta^2) \frac{P \cdot R}{\beta^2 P + R} = F_B$$

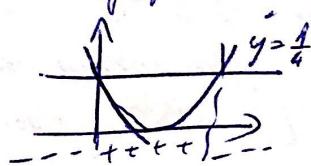
II

d) $S^+ \subset C(0,1)$ and $S^- \subset R \setminus [0,1]$ be finite sets of positive and negative reals. Define $\Phi: R \rightarrow R$ as $\forall x \in S^+$, $\Phi(x) = x + \epsilon^+$ and $\{\Phi(x) : x \in S^-\}$ be one linearly separable?

$$\Phi: R \rightarrow R^2, N=2$$

$$X \rightarrow (x, (x - \frac{1}{2})^2)$$

see we define the points



2) $S^+ \subset B$ and $S^- \subset R \setminus B$ be finite sets with

$$B = \{(x, y) : x^2 + y^2 \leq 1\}, \Phi: R^2 \rightarrow R$$
 as $\forall x \in S^+$, $\{\Phi(x) : x \in S^+\}$ and $\{\Phi(x) : x \in S^-\}$ be one linearly separable?

$$\Phi: R^2 \rightarrow R$$

$$(x, y) \rightarrow x^2 + y^2 + 2$$

observe if $(x, y) \in S^+ \subseteq B$, $x^2 + y^2 + 2 \leq 1$

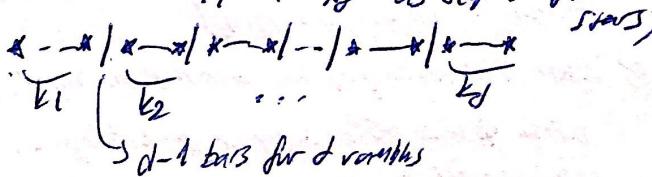
if $(x, y) \in S^- \subseteq R \setminus B$, $x^2 + y^2 + 2 > 1$

so point 1 is separating hyperplane

c) X_1, \dots, X_d real valued variables, moment of degree k over X_1, \dots, X_d is $B = \prod_{i=1}^d X_i^{k_i}$ with nonnegative k_1, \dots, k_d . $\sum k_i = l$

d) Let D be num. of monomials of degree k over X_1, \dots, X_d
show that $D = \binom{d+k-1}{d-1} = \binom{d+k-1}{k}$?

encode monomials $X_1^{k_1} \dots X_d^{k_d}$ to $X_j^{k_j}$ as sequence of bits and



$$\frac{(d+k-1)!}{(d-1)! \cdot k!} = \binom{d+k-1}{d-1} = \binom{d+k-1}{k}$$

b) $\Phi: R^d \rightarrow R^D$ with $\Phi = \Phi = (d_1, \dots, d_D) \mapsto (M_1(\Phi), M_2(\Phi), \dots)$

for $\Phi \in R^d$, $M_1 \mapsto M_2 \mapsto \dots$ here follow number of degree k over

X_1, \dots, X_d . $S^+ \subseteq R^d$ with $|S^+| = n$, $S^- = \{\Phi(S) : S \in S^+\}$, $(S^{(1)}, S^{(-)})$ be dichotomy of S^+ characteristically about.

Give prob. $P = P(n, d, k)$ of event that $(S^{(1)}, S^{(-)})$ is linearly separable?

probability that dichotomy random selected uniformly from general positions of R^N . Is linearly separable by $P = \frac{1}{2^n} \left(\sum_{k=0}^n \binom{n-1}{k} \right)$

there are ~~$\frac{1}{2^n}$~~ $D = \binom{d+k-1}{d-1}$ monomials of degree k choosing $N = D$,

$$P = \frac{1}{2^n} \left(\sum_{k=0}^D \binom{n-1}{k} \right)$$

d) Prove $-S^+ \Phi(S^+) \cup (S^- \Phi(S^-)) \subseteq R^d \times R$,
and vector $w \in R^d$ s.t. $f(w) = (x^2, y^2)$ has no plate S
meaning minimaes empirical error other w.r.t. square loss,
minimize $\text{Emp}(f) = \frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2$?

claim = empirical risk w.r.t. square loss function is

$$\text{Emp}(f) = \frac{1}{n} \|xw - y\|^2$$

pref $\vec{y} = (y_1, \dots, y_n)$ and X be $n \times d$ matrix
that rows are vectors x_1, \dots, x_n .

$$x = \begin{pmatrix} x_{11} & \dots & x_{1d} \\ x_{21} & \dots & x_{2d} \\ \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{nd} \end{pmatrix} \text{ } n \times d$$

$$\text{Emp}(f) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2$$

$$= \frac{1}{n} \left(\sum_{i=1}^n (x_i w - y_i)^2 \right)$$

$$\stackrel{\text{definition}}{=} \frac{1}{n} \|xw - \vec{y}\|^2 \text{ where}$$

$$\|a\| = (a_1^2 + \dots + a_d^2)^{\frac{1}{2}}, a \in R^d$$

2) Prove convex min prob.

$$\min R_2(w, s) = \min_w \|xw - y\|^2 + 2\|w\|^2 / R^2,$$

claim = $\nabla_w R_2(w, s) = 2x^T xw - 2x^T y + 2Rw$

$$\nabla^2 f = \nabla_w (\langle xw - y, xw - y \rangle + 2\langle w, w \rangle)$$

$$= \nabla_w (\underbrace{x^T x w - 2x^T y w + y^T y}_{\text{by } \|w\|^2 = \langle w, w \rangle = \langle w, w \rangle} + 2\langle w, w \rangle)$$

$$= \nabla_w (w^T x^T x w) = \nabla_w (w^T x^T) \cdot x w + w x^T \cdot \nabla_w (x w)$$

$$= x^T x w + w^T x^T x = 2x^T x w$$

$$2 = \nabla_w (-2y^T x w) = -2 \nabla_w ((y^T x) w + y^T x w)$$

$$= -2x^T y$$

$$3 = \nabla_w (w^T w) = \nabla_w (w_1^2 + \dots + w_d^2) = 2(2w_1 + \dots + 2w_d) = 2Rw$$

$$= 2x^T x w - 2x^T y w + 2Rw$$

$\mathcal{L} \subseteq \mathbb{R}^d \times \mathbb{R}^{d-1}$ the set of d -dim vectors
at least one positive, one negative entries
Let $C = \frac{1}{\|x\|_2} \sum_{i=1}^d x_i$ and $C^- = \frac{1}{\|x\|_2} \sum_{i=1}^d x_i$ be mean of

+ out-exp class assign vector $x \in \mathbb{R}^d$ to class whose mean
is closer to x , i.e. correctly classified. Can implement x
easily through inner products?

assign each non point to class which minimizes
distance to its mean

$$f(x) = \text{sgn}(\|C - x\|^2 - \|C^+ - x\|^2)$$

$$= \text{sgn}(\langle C - x, C - x \rangle - \langle C^+ - x, C^+ - x \rangle)$$

$$\text{we have } f(x) = \begin{cases} 1 & \text{if } x \text{ is closer to } C \\ -1 & \text{if } x \text{ is closer to } C^+ \end{cases}$$

α

1) $k: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ be kernel function, is it true that underlying
feature space embedding func. Φ (ex for which $k(x, y) = \langle \Phi(x), \Phi(y) \rangle$)
 $\langle \Phi(x), \Phi(y) \rangle$ is unique, why or not?

No, Φ be embedding such $k(x, y) = \langle \Phi(x), \Phi(y) \rangle$,
then we have $\langle -\Phi(x), -\Phi(y) \rangle = \langle \Phi(-x), \Phi(-y) \rangle$
 $= k(x, y)$
so $-\Phi$ is also feature space embedding

2) $k: \mathbb{R}^2 \times \mathbb{R}^2 \rightarrow \mathbb{R}$ be first defined by $k(x, y) = \langle x, y \rangle + 1$
for all $x, y \in \mathbb{R}^2$, is k kernel function, if not
give first Φ for which $k(x, y) = \langle \Phi(x), \Phi(y) \rangle$?

$$\begin{aligned} \text{Yes, } k(x, y) &= \langle x, y \rangle + 1 \\ &= \langle x, y \rangle + 2\langle x, y \rangle + 1 \\ &= (x_1 y_1 + x_2 y_2)^2 + 2(x_1 y_1 + x_2 y_2) + 1 \\ &= x_1^2 y_1^2 + 2x_1 y_1 x_2 y_2 + x_2^2 y_2^2 + 2x_1 y_1 + 2x_2 y_2 + 1 \end{aligned}$$

$$\text{choose } \Phi(x, y) = (x_1^2, \sqrt{2}x_1 x_2, x_2^2, \sqrt{2}x_1 \sqrt{2}x_2, 1)$$

for all $x = (x_1, x_2)$ and $y = (y_1, y_2)$

We have $k(x, y) = \langle \Phi(x), \Phi(y) \rangle$

3) $k: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ be defined $k(x, y) = \sum_{i=1}^d (x_i + y_i)$ for all
 $x, y \in \mathbb{R}^d$, is k kernel?

No, inner product of vector with i -th value 0
consider $x = (-1, 0, \dots) \in \mathbb{R}^d$ nonnegative,
 $y = (1, 0, \dots) \in \mathbb{R}^d$

$$\Rightarrow k(x, y) = -2$$

so there cannot be any Φ s.t. $k(x, y) = \langle \Phi(x), \Phi(y) \rangle$

4) $k: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ defined by $k(x, y) = \prod_{i=1}^d \left(\frac{x_i + y_i}{b} \right)$
forall $x, y \in \mathbb{R}^d$ and either with $b \neq 0$, is k kernel function?

Yes, $k(x, y) = f(x) \cdot f(y)$ is kernel over X for all functions

$$f: X \rightarrow \mathbb{R} \text{ set } f(x) = \prod_{i=1}^d \left(\frac{x_i + y_i}{b} \right)$$

then $k(x, y) = f(x) \cdot f(y)$, k is kernel

5) $k: \mathbb{R}^d \setminus \{0\} \times \mathbb{R}^d \setminus \{0\} \rightarrow \mathbb{R}$ defined by

$$k(x, y) = -\frac{\langle x, y \rangle}{\|x\|_2 \|y\|_2}, \text{ is } k \text{ kernel function?}$$

No, consider $d=1$, $k(x, y) = -\frac{x \cdot y}{\|x\|_2 \|y\|_2} = -1 \leq 0$

but $\langle \Phi(x), \Phi(y) \rangle \geq 0$ should be for any inner product
so k is not kernel

6) V be finite set and $k: 2^V \times 2^V \rightarrow \mathbb{R}$ by $k(S_1, S_2) = |S_1 \cap S_2|$
forall $S_1, S_2 \subseteq V$, is k kernel function?

Yes, let $V = \{v_1, \dots, v_m\}$

define $\Phi: 2^V \rightarrow \{0, 1\}^m$

$$S \mapsto x^S \text{ where } x_i^S = \begin{cases} 1 & \text{if } v_i \in S \\ 0 & \text{otherwise} \end{cases}$$

for $S_1, S_2 \subseteq V$,

$$\langle \Phi(S_1), \Phi(S_2) \rangle = \sum_{i=1}^m x_i^{S_1} \cdot x_i^{S_2}$$

$\underbrace{1}_{\text{iff } v_i \in S_1 \cap S_2}$

$$= |S_1 \cap S_2|$$

7) $k: X \times X \rightarrow \mathbb{R}$ our $X = \{x_1 \rightarrow x_2\}$ for $n \in N$ and $k(x, y)$
Gram (or kernel) matrix of X wrt k .

1) for any nonnegative real number c , define $\kappa_{x,y}$ and κ_x^c

$\kappa_{x,y}^c$ by $(\kappa_{x,y}^c)_{ij} = (\kappa_{x,y})_{ij} + c$ for all $i, j = 1, \dots, n$,

Show $\kappa_{x,y}^c$ is Gram matrix of X wrt k' for some kernel

$k': X \times X \rightarrow \mathbb{R}$ over X ?

k' is kernel and where Φ a feature map $\Phi: X \rightarrow F$
with $\kappa(x, y) = \langle \Phi(x), \Phi(y) \rangle$

define Φ' , $X \rightarrow (F \times \mathbb{R})$

$$x \mapsto (\Phi(x), \sqrt{c})$$

define $\langle (x, \sqrt{c}), (y, \sqrt{c}) \rangle_{F \times \mathbb{R}} = \langle x, y \rangle + \sqrt{c}$

for $x, y \in F$, $x, y \in \mathbb{R}$

define $K'(x, y) = \kappa(x, y) + c$

$$\langle \Phi(x), \Phi(y) \rangle = \langle (\Phi(x), \sqrt{c}), (\Phi(y), \sqrt{c}) \rangle_{F \times \mathbb{R}}$$

$$= \langle \Phi(x), \Phi(y) \rangle + c$$

$$= \kappa(x, y) + c$$

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

so k' is kernel with $\kappa_{x,y}$ as its gram matrix

51

2) for any nonnegative real number c , define a norm matrix $K_{X \times E}$

$$K_{X \times E}^c (x_i, y_j) = \begin{cases} (K_{X \times E})_{ij} & \text{if } i \neq j \\ (K_{X \times E})_{ii} + c & \text{otherwise} \end{cases}$$

Show that $K_{X \times E}^c$ is Gram matrix of X w.r.t. E for some kernel $K' : X \times X \rightarrow \mathbb{R}$ over X ?

K is kernel that there is feature map $\Phi : X \rightarrow \mathbb{R}^n$ with

$$K(x_i, y_j) = \langle \Phi(x_i), \Phi(y_j) \rangle$$

define $\Phi' : X \times X \rightarrow \mathbb{R}^n$

$$x \mapsto (\Phi(x), \sqrt{c}, \mathbf{e}_n)$$

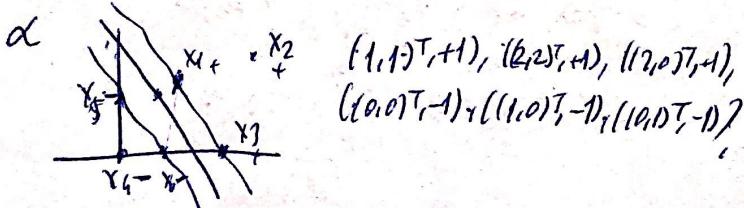
$$x \in \mathbb{R}^n, (\Phi(x))_i = \begin{cases} 1 & \text{if } x = x_i \\ 0 & \text{otherwise} \end{cases}$$

define $\langle (x, \tilde{x}), (y, \tilde{y}) \rangle_{F_X \times F_X} = \langle x, y \rangle + \langle \tilde{x}, \tilde{y} \rangle_{\mathbb{R}^n}$

define $K'(x_i, y_j) = K(x_i, y_j) + c \langle \mathbf{e}_n, \mathbf{e}_n \rangle_{\mathbb{R}^n}$

$$\begin{aligned} \langle \Phi(x_i), \Phi(y_j) \rangle &= \langle (\Phi(x_i), \sqrt{c}, \mathbf{e}_n), (\Phi(y_j), \sqrt{c}, \mathbf{e}_n) \rangle \\ &= \langle \Phi(x_i), \Phi(y_j) \rangle + c \langle \mathbf{e}_n, \mathbf{e}_n \rangle_{\mathbb{R}^n} \\ &= K(x_i, y_j) + c \langle \mathbf{e}_n, \mathbf{e}_n \rangle \\ &= K'(x_i, y_j) \end{aligned}$$

so K' is kernel with $K_{X \times E}$ as its gram matrix



1) Give primal optimization problem for hard margin SVM?

$$\min \frac{1}{2} \|\vec{w}\|^2 \text{ s.t. } y_i (\langle \vec{w}, \vec{x}_i \rangle + b) \geq 1 \quad \forall i \in \{1, 2, \dots, n\}$$

for each training data point $i = 1 \dots n$

$\vec{w} = (w_1, w_2)$ and plugging in data points into inequality

$$\min \frac{1}{2} (w_1^2 + w_2^2) \text{ subject to}$$

$$-1((w_1 + w_2) + b) + 1 \leq 0$$

$$-1((2w_1 + 2w_2) + b) + 1 \leq 0$$

$$-1((2w_1) + b) + 1 \leq 0$$

$$b + 1 \leq 0$$

wrong Lagrangian,

$$\min_{\vec{w}, b} \max_{d \geq 0} \left(-\frac{1}{2} \|\vec{w}\|^2 + \sum_{i=1}^n d_i y_i (\langle \vec{w}, \vec{x}_i \rangle + b) \right)$$

$$\min_{\vec{w}, b} \max_{d \geq 0} \left(\frac{w_1^2 + w_2^2}{2} + \sum_{i=1}^n d_i y_i (\langle \vec{w}, \vec{x}_i \rangle + b) \right)$$

where d_i 's are LHS of inequalities

2) Give dual optimization problem for hard margin SVM?

$$L_d(\vec{a}) = \max_{\vec{a} \geq 0} \left(\sum_{i=1}^n a_i - \frac{1}{2} \sum_{i,j=1}^n a_i a_j y_i y_j (\vec{x}_i \cdot \vec{x}_j) \right)$$

$$\text{s.t. } \sum_{i=1}^n a_i y_i = a_1 + a_2 + a_3 - a_4 - a_5 - a_6 = 0$$

3) What are support vectors?

$$(0, 1)^T, (1, 0)^T, (1, 1)^T, (2, 0)^T$$

4) How large is margin margin?

$$\gamma = \frac{2}{\|\vec{w}\|} = \frac{2}{\sqrt{w_1^2 + w_2^2}} = \frac{\sqrt{2}}{2}$$

5) Give dualization of dual optimization problem for soft margin SVM?

primal optimization for $C > 0$,

$$\min_{\vec{w}, b} C \sum_{i=1}^n \varepsilon_i + \frac{1}{2} \|\vec{w}\|^2 \text{ s.t. } y_i (\langle \vec{w}, \vec{x}_i \rangle + b) \geq 1 - \varepsilon_i$$

$$\text{where } \varepsilon_i \geq 0, i = 1 \dots n$$

Lagrangian multiplier a_i, ε_i

$$L(w, b, \varepsilon, a, \lambda) = C \sum_{i=1}^n \varepsilon_i + \frac{1}{2} \|\vec{w}\|^2 - \sum_{i=1}^n a_i y_i (\langle \vec{w}, \vec{x}_i \rangle + b) - (1 - \varepsilon_i) - \sum_{i=1}^n \lambda_i \varepsilon_i$$

optimization problem will be,

$$\max_{\vec{w}, b} \min_{\varepsilon, a} L(w, b, \varepsilon, a, \lambda) \text{ s.t. } a_i, \lambda_i \geq 0.$$

compute derivative and set it to 0,

$$\frac{\partial L}{\partial \vec{w}} = \vec{w} - \sum_{i=1}^n a_i y_i \vec{x}_i = 0$$

$$\vec{w} = \sum_{i=1}^n a_i y_i \vec{x}_i$$

$$\frac{\partial L}{\partial b} = - \sum_{i=1}^n a_i y_i = 0$$

$$\sum_{i=1}^n a_i y_i = 0$$

$$\frac{\partial L}{\partial a_i} = C - a_i - x_i = 0$$

$$C = a_i + x_i \text{ where } 0 \leq a_i \leq C$$

$$\max_{\vec{w}, b} \min_{\varepsilon, a} L(w, b, \varepsilon, a, \lambda)$$

$$\max_{\vec{w}, b} \min_{\varepsilon, a} C \sum_{i=1}^n \varepsilon_i + \sum_{i=1}^n a_i - \frac{1}{2} \sum_{i,j=1}^n a_i a_j y_i y_j (\vec{x}_i \cdot \vec{x}_j)$$

$$\rightarrow \sum_{i=1}^n a_i \varepsilon_i - \sum_{i=1}^n a_i \varepsilon_i$$

$$\max_{\vec{w}, b} \min_{\varepsilon, a} \sum_{i=1}^n (C - a_i - x_i) \varepsilon_i + \sum_{i=1}^n a_i - \frac{1}{2} \sum_{i,j=1}^n a_i a_j y_i y_j (\vec{x}_i \cdot \vec{x}_j)$$

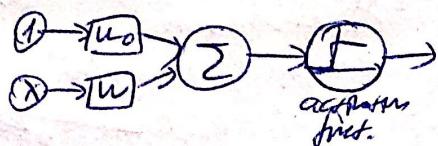
$$\Rightarrow \max_{\vec{w}, b} \min_{\varepsilon, a} \sum_{i=1}^n a_i - \frac{1}{2} \sum_{i,j=1}^n a_i a_j y_i y_j (\vec{x}_i \cdot \vec{x}_j) \text{ s.t. } 0 \leq a_i \leq C$$

$$\sum_{i=1}^n a_i \varepsilon_i = 0$$

α Constant perceptors

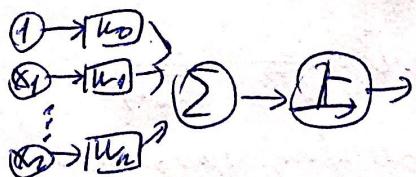
a) $f(x) = \gamma x$?

$$\text{sgn}(z) = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{else} \end{cases}$$



$$w_0 = 1, w_1 = -1, \text{ output is } 1 \text{ iff } n = 0$$

b) $f(x_1, \dots, x_n) = \bigvee_{i=1}^n x_i$?

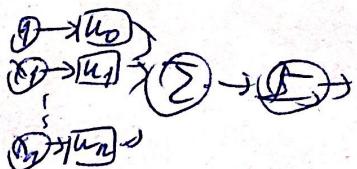


$$w_0 = -n+1, w_i = 1 \text{ for } i = 1 \dots n$$

$$\sum_{i=0}^n x_i w_i = \underbrace{x_0 w_0}_{-n+1} + \sum_{i=1}^n x_i w_i$$

output is 1 iff all inputs are 1

c) $f(x_1, \dots, x_n) = \bigvee_{i=1}^n x_i$?



$$w_0 = -n+1, w_i = 1 \text{ for } i = 1 \dots n$$

$$\sum_{i=0}^n x_i w_i = \underbrace{x_0 w_0}_{-n+1} + \sum_{i=1}^n x_i w_i$$

output is 1 iff at least one input is 1

α Can you give perceptron for $(x_1 \wedge x_2 \wedge x_3) \vee (x_1 \wedge x_2 \wedge \neg x_3)$?

~~assume there is no perceptron~~

claim = it is not linearly separable.

proof = assume there is a perceptron satisfies this,

then for x_3 s.t. $\sum_{i=1}^3 w_i x_i > 0$ iff formula is satisfied

input $(0,0,0)$ satisfies, $w_0 > 0$

input $(1,0,0), (0,1,0), (0,0,1)$ doesn't satisfy, $w_0 < 0, w_1 < 0$

$w_2 < 0$

$w_0 + w_1 + w_2 < 0$

input $(1,1,1)$ satisfies the formula, $w_0 + w_1 + w_2 > 0$

but this is contradiction because $w_0 > 0$ and $3w_0 + w_1 + w_2 < 0$

for w_0, w_1, w_2 should be smaller than 0.

α Show every boolean function can be represented by some 2-layer neural network?

every boolean function $f(x_1, \dots, x_n)$ can be written in CNF

$$f = \bigwedge_{i \in I} C_i \text{ where } C_i = \bigvee_{j \in J} x_j$$

first perceptron C_i for each C_i ,

$$C_i = b = 1$$

$$w_j = \begin{cases} \frac{1}{|C_i|} & \text{if } x_j \in C_i \\ -\frac{1}{|C_i|} & \text{if } \neg x_j \in C_i \\ 0 & \text{else} \end{cases}$$

so C_i implements C_i , join all C_i together in 2nd layer implementing conjunction over I

$$A = \cdot b = -1$$

$$w_i = \frac{1}{|I|} + \frac{1}{|I|} x_i^2$$

α Suppose $(x_1, y_1), \dots, (x_n, y_n)$ are training examples

with $x_i \in \mathbb{R}^d$ and $y_i \in \{-1, 1\}$ s.t. y_i 's

- have bounded Euclidean norm $\|x_i\| \leq R$ for OLR,
- are linearly separable with margin $\gamma > 0$, where $\exists w \in \mathbb{R}^d$ with $\|w\| = 1$ s.t.

$$y_i(w^T x_i) \geq \gamma \text{ for all } i = 1 \dots n.$$

Let $w^T = \beta + \beta^T x_i \rightarrow w^T \in \mathbb{R}^d$ be weights calculated

ideally by perceptron learning rule. Perceptron Blackboard

- Euclidean norm of all examples in D is bounded by R

- positive and negative examples in D are linearly separable

with finite margin $\gamma > 0$, then for the number of k of updates in loss: $W_{k+1} = W_k + y_i x_i$ (perceptron rule)

perceptron also, we have $k \leq \left(\frac{R}{\gamma}\right)^2$, provided?

perceptron algo updates weights if classification was made wrong or β is bounded by $\left(\frac{R}{\gamma}\right)^2$.

vector that was categorized after i -th mistake was made

i -th class = $x_i^T D \geq k \gamma$

let i to k iterate indices of examples when mistake was made.

$$W_k^T D = \left(\sum_{j=1}^k x_j y_j \right)^T D = \sum_{j=1}^k y_j x_j^T D \geq \sum_{j=1}^k \gamma$$

W_k comes from W_{k-1} by adding $y_k x_k$... = $L \gamma$

points are separable with margin γ

~~2nd class~~



5

$$-2nd \text{ claim: } \|\vec{w}_k\|^2 \leq \|\vec{w}_{k-1}\| + R^2$$

$$\text{Let's note } \|\vec{w}_k\|^2 = \vec{w}_k \cdot \vec{w}_k$$

$$= (\vec{w}_{k-1} + y_{ik} \vec{x}_{ik}) \cdot (\vec{w}_{k-1} + y_{ik} \vec{x}_{ik})$$

$$= \vec{w}_{k-1} \cdot \vec{w}_{k-1} + 2y_{ik} \vec{w}_{k-1} \cdot \vec{x}_{ik} + y_{ik} \vec{x}_{ik} \cdot \vec{y}_{ik} \vec{x}_{ik}$$

observe $2y_{ik} \vec{w}_{k-1} \cdot \vec{x}_{ik} \leq 0$ because mistake was made,
 $y_{ik} y_{ik} = 1$ so we have

$$\begin{aligned} \|\vec{w}_k\|^2 &\leq \vec{w}_{k-1} \cdot \vec{w}_{k-1} + \vec{x}_{ik} \cdot \vec{x}_{ik} \\ &\leq \|\vec{w}_{k-1}\|^2 + \|\vec{x}_{ik}\|^2 \\ &\leq \|\vec{w}_{k-1}\|^2 + R^2 \end{aligned}$$

now we can combine,

$$\text{by first claim, } (kj)^2 \leq (\vec{w}_k \vec{b})^2 \leq \|\vec{w}_k\|^2 \cdot \|\vec{b}\|^2$$

$$(kj)^2 \leq \|\vec{w}_k\|^2 \quad \boxed{1}$$

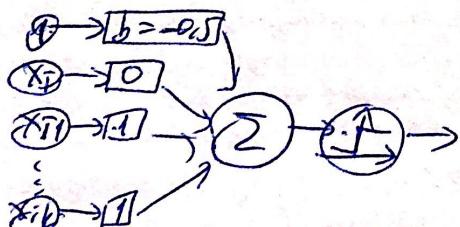
$$\text{by 2nd claim } (kj)^2 \leq \|\vec{w}_{k-1}\|^2 + R^2$$

$$\leq \|\vec{w}_{k-1}\|^2 + 2R^2$$

$$(kj)^2 \leq \|\vec{w}_0\|^2 + kR^2 = kR^2$$

solving for k , it yields $k \leq \frac{R^2}{j^2}$, i.e. is proved.

Q 1) Show that any k -class misclassification in SOL^n can be represented by perceptron?



$$w_i = \begin{cases} 1 & \text{if } j \in \{1, -1\} \\ 0 & \text{otherwise} \end{cases}$$

so decisions over labels can be realized by perceptron.

2) Give upper bound on number of prediction mistakes made by perceptron also during learning of unknown k -class misclassification in SOL^n , give your upper bound?

7

Ensure ~~that~~ to find separating hyperplane through origin.
 Let $(x_1, y_1), \dots, (x_n, y_n)$ form training set
 transform each $\vec{x}_i = (x_{i1}, \dots, x_{in})$ to $\vec{v}_i = (1, y_1, \dots, y_n)$
 so all are dimensioned the same and all points get same coordinate in this new dimension.

The points \vec{v}_1 to \vec{v}_n are separable by hyperplane through origin by $\vec{b} = (-1, v_1, \dots, v_n)$ where $v_i = \begin{cases} 2 & \text{if } j \in \{1, -1\} \\ 0 & \text{otherwise} \end{cases}$

then, $\vec{b} \cdot \vec{v}_i \geq 1$ if \vec{v}_i satisfies formula,
 $\vec{b} \cdot \vec{v}_i \leq -1$ otherwise

$$\text{normalize } \vec{b}' \text{, then } \vec{b} = \frac{\vec{b}'}{\|\vec{b}'\|} = \frac{\vec{b}'}{\sqrt{Gk+1}}$$

apply theorem, radius of any input point is bounded by $\sqrt{Gk+1}$ and margin is lower bounded by

$$j \geq \frac{1}{\sqrt{Gk+1}}$$

$$\text{we get } k \leq \frac{R^2}{j^2} = (n+1) \cdot (4k+1)$$