

① We need to find β such that: $P[w^T X - b \geq 0] \geq \beta$
i.e. β such that $P[b - w^T X > 0] \leq 1 - \beta$

$$P[b - w^T X > 0] = P\left[-\frac{2}{3}X_1 + X_2 > 0\right]$$
$$= P\left[\underbrace{-\frac{2}{3}(X_1 - 2) + (X_2 - \frac{3}{4})}_{\text{say, } Z} > \frac{7}{12}\right]$$

Our bounds answer the upper bound query; hence this trick

Our bounds work on mean zero too; hence this trick.

$$= P[e^{\lambda Z} > e^{\lambda \cdot 7/12}] \quad \forall \lambda > 0$$

$$\leq \frac{E[e^{\lambda Z}]}{e^{\lambda \cdot 7/12}} \quad \rightarrow \text{Markov Inequality}$$

$$= \frac{E\left[e^{-\frac{2\lambda}{3}(X_1 - 2)}\right] E\left[e^{\lambda(X_2 - \frac{3}{4})}\right]}{e^{\lambda \cdot 7/12}} \quad (\because X_1 \text{ is ind. of } X_2)$$

Now,
 $0.5 \leq X_1 \leq 2.5 \Rightarrow -\frac{1}{3} \leq -\frac{2}{3}(X_1 - 2) \leq 1 \Rightarrow E\left[e^{-\frac{2\lambda}{3}(X_1 - 2)}\right] \leq e^{\lambda \cdot 2/9}$

$0.5 \leq X_2 \leq 1.5 \Rightarrow -\frac{1}{4} \leq X_2 - \frac{3}{4} \leq \frac{3}{4} \Rightarrow E\left[e^{\lambda(X_2 - \frac{3}{4})}\right] \leq e^{\lambda \cdot 3/8}$

Hoeffding bound

$$\therefore \leq e^{2\lambda/9} e^{3\lambda/8} e^{-\lambda \cdot 7/12} \quad \forall \lambda > 0 \quad \left. \vphantom{\leq} \right\} \text{ Chernoff bounding technique}$$

$$\Rightarrow \leq e^{-0.245} \approx 0.7827$$

$$\Rightarrow P[w^T X - b \geq 0] \geq 1 - e^{-0.245} \approx 0.2173$$

* Given only the mean of the test data point one would label it to be a positive data point. But the prob. that it is indeed such one is can to take into account only

② Using hints given in lectures, Smola book we arrive at following claim:

claim $P\left[\sup_{f \in \mathcal{F}} \left\{ \text{Rep}^m[f] - R[f] \right\} > \epsilon\right] \rightarrow 0 \text{ as } m \rightarrow \infty \quad \forall \epsilon > 0$
 (Required ~~necessary~~ suff. cond.)

$$\sup_{f \in \mathcal{F}} \text{Rep}^m[f] \xrightarrow{P} \sup_{f \in \mathcal{F}} R[f] \quad (\text{Emp. risk. max.})$$

i.e. $P\left[\left| \sup_{f \in \mathcal{F}} \text{Rep}^m[f] - \sup_{f \in \mathcal{F}} R[f] \right| > \epsilon\right] \rightarrow 0 \text{ as } m \rightarrow \infty$

Proof:

$$P\left[\left| \sup_{f \in \mathcal{F}} \text{Rep}^m[f] - \sup_{f \in \mathcal{F}} R[f] \right| > \epsilon\right] = P\left[\sup_{f \in \mathcal{F}} \text{Rep}^m[f] - \sup_{f \in \mathcal{F}} R[f] > \epsilon\right] + P\left[\sup_{f \in \mathcal{F}} R[f] - \sup_{f \in \mathcal{F}} \text{Rep}^m[f] > \epsilon\right]$$

$$\Rightarrow \left(\text{Let } f^m = \underset{f \in \mathcal{F}}{\text{argmax}} \text{Rep}^m[f] ; f^{\text{opt}} = \underset{f \in \mathcal{F}}{\text{argmax}} R[f] \right)$$

$$\Rightarrow P\left[R^m[f^m] - R[f^{\text{opt}}] > \epsilon\right] + P\left[R[f^{\text{opt}}] - \text{Rep}^m[f^m] > \epsilon\right]$$

$$= P\left[\underbrace{R^m[f^m] - R[f^m]}_{\leq 0} + \underbrace{R[f^m] - R[f^{\text{opt}}]}_{\leq 0} > \epsilon\right] + P\left[\underbrace{R[f^{\text{opt}}] - \text{Rep}^m[f^{\text{opt}}]}_{\leq 0} + \underbrace{\text{Rep}^m[f^{\text{opt}}] - \text{Rep}^m[f^m]}_{\leq 0} > \epsilon\right]$$

$$\leq P\left[R^m[f^m] - R[f^m] > \epsilon\right] + P\left[R[f^{\text{opt}}] - \text{Rep}^m[f^{\text{opt}}] > \epsilon\right]$$

$$\leq P\left[\sup_{f \in \mathcal{F}} \left\{ \text{Rep}^m[f] - R[f] \right\} > \epsilon\right] + P\left[R[f^{\text{opt}}] - \text{Rep}^m[f^{\text{opt}}] > \epsilon\right]$$

$\downarrow m \rightarrow \infty$
 0
 by ~~necessary~~ suff. cond. given

$\downarrow m \rightarrow \infty$
 0
 by Chernoff's bound.

* In other words, the ~~two~~ ^{one}-sided uniform conv. are necessary for ER Min. & ER Max. & hence, the two-sided (regular) uniform conv. conditions are necessary for both ER Min & ER Max.

③ Refer pg. 40 (Thm. 1 & Proof) in appendix of Burger Tutorial.

④ Let $\mathcal{F}_{11} \equiv$ set of all even 11° ellipses

$\mathcal{F}_{45} \equiv$ set of all ellipses at 45°

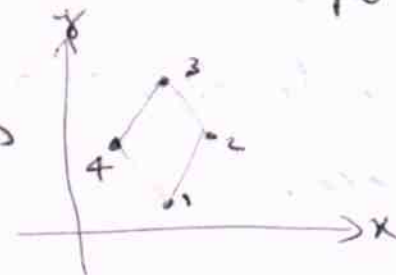
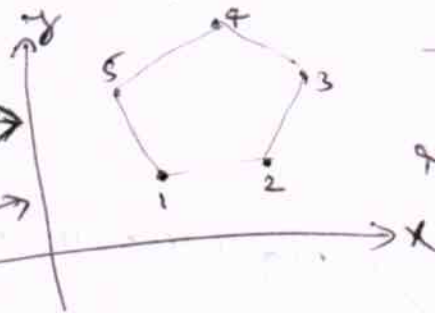
$\mathcal{F} = \mathcal{F}_{45} \cup \mathcal{F}_{11}$

(Everywhere we are drawing regular polygons)

It is easy to see that \mathcal{F}_{11} shatters

also \mathcal{F} shatters

\mathcal{F}_{11} does not shatter



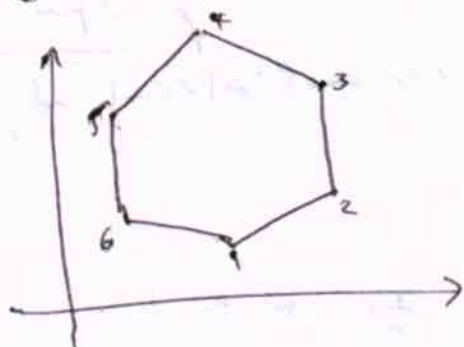
square
regular

regular pentagon

(In some sense regular polygons are optimal arrangements; we need to look at orientation with axis though for this problem)

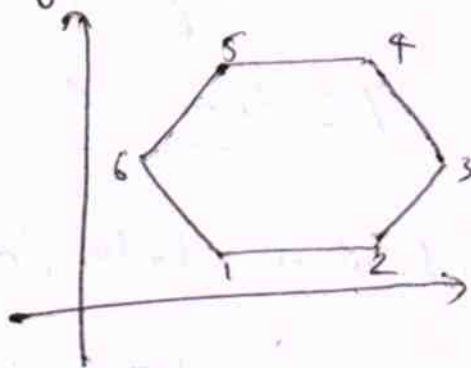
Here \mathcal{F}_{11} 's VC dim. = 4 (see also ndm. to next problem)

\mathcal{F} does not shatter



regular hexagons in both orientations

and



Here \mathcal{F} 's VC dim. = 5

* It is easy to verify all claims regarding shattering made here by imagining ellipses to be limiting cases of rectangles

⑤ Refer to paper titled "Classification by Polynomial Surfaces" by Martin Anthony

or
 simply use the fact illustrated in lecture that polynomial discriminants are essentially linear discriminants in expanded feature space.

For J_{11} :

$$f(x) = \text{sign}((x-c)^T S (x-c) - 1) = \text{sign}\left(\begin{bmatrix} x_1 - c_1 & x_2 - c_2 \end{bmatrix} \begin{bmatrix} D_1 & 0 \\ 0 & D_2 \end{bmatrix} \begin{bmatrix} x_1 - c_1 \\ x_2 - c_2 \end{bmatrix} - 1\right)$$

$$= \text{sign}(\omega_1 z_1 + \omega_2 z_2 + \omega_3 z_3 + \omega_4 z_4 + \omega_5)$$

where $\omega_1 = D_1, \omega_2 = D_2, \omega_3 = -2D_1 c_1, \omega_4 = -2D_2 c_2, \omega_5 = D_1 c_1^2 + D_2 c_2^2 - 1$

$\&$ $z_1 = x_1^2, z_2 = x_2^2, z_3 = x_1, z_4 = x_2$

\downarrow
 fixed if
 $\omega_1, \omega_2, \omega_3, \omega_4$
 are fixed.

This is essentially ~~linear~~ linear discriminants in \mathbb{R}^4 passing through fixed point

\therefore VC dim = 4 (not 5 since ~~the~~ threshold is fixed)

⑥ Bound which is to be plotted is:

$$\left\{ \begin{array}{l} 4 \exp\left\{ m \log 2 - m \epsilon^2 / 8 \right\} \quad \text{for } m \leq 3 \\ 4 \exp\left\{ 3 \left(\log \frac{m}{3} + 1 \right) - m \epsilon^2 / 8 \right\} \quad \text{for } m > 3 \end{array} \right.$$

Section II

Assignment 3 Solutions

Q2 In 3-d, the VC dim. of hyperplanes is 4, which is upper limit on that for conical hyperplanes

The margin margin achievable with 4 points lying on a sphere is when they are on a tetrahedron (regular).

height of tetrahedron = $\frac{4}{3}R$ (\because margin = $\frac{2}{\|w\|} \leq \frac{4R}{3}$)

\therefore If $\|w\| \geq \frac{3}{2R}$ then VC dim = 4

Next possibility is still higher margin where 3 pts. can be scattered \rightarrow pts. lying on an equilateral Δ with circumradius = R.

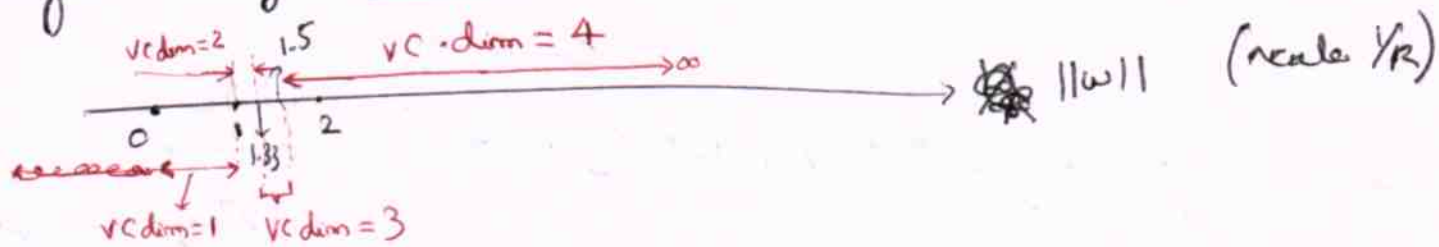
\rightarrow height of ~~the~~ this eq. Δ = $\frac{3R}{2}$ (\because margin = $\frac{2}{\|w\|} \leq \frac{3R}{2}$)

\therefore If $\frac{3}{2R} \leq \|w\| < \frac{3}{R}$, then VC dim = 3

|| by next possibility is 2 pts. & margin $\leq 2R$

\Rightarrow If $\frac{1}{R} \leq \|w\| < \frac{4}{3R}$, then VC dim = 2

of course if $\|w\| < \frac{1}{R}$, then VC dim = 1.



Q3 SVM primal with reg. loss:

min w, b, ξ $\frac{1}{2} \|w\|^2 + C \sum \xi_i$

s.t. $y_i(w^T x_i - b) \geq 1 - \xi_i, \xi_i \geq 0.$

$$\mathcal{L} = \frac{1}{2} \|\omega\|^2 + C \sum \xi_i^2 + \sum \alpha_i (y_i \omega^T x_i + y_i (b+1 - \xi_i)) - \sum \beta_i \xi_i$$

$$\nabla_{\omega} \mathcal{L} = 0 \Rightarrow \omega = \sum \alpha_i y_i x_i$$

$\alpha_i, \beta_i \geq 0$ are L. mul.

$$\frac{\partial \mathcal{L}}{\partial \xi_i} = 0 \Rightarrow \alpha_i + \beta_i = 2C \xi_i \quad \leftarrow \text{(only change compared to hinge-loss)}$$

$$\frac{\partial \mathcal{L}}{\partial b} = 0 \Rightarrow \sum \alpha_i y_i = 0$$

$$\mathcal{L} = -\frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j x_i^T x_j + \sum \alpha_i + \sum_i \frac{(\alpha_i + \beta_i)^2}{4C}$$

(eliminating ω, b, ξ_i)

Dual is $\max_{\alpha_i, \beta_i} -\frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j x_i^T x_j + \sum \alpha_i + \sum_i \frac{(\alpha_i + \beta_i)^2}{4C}$

s.t. $\sum \alpha_i y_i = 0, \alpha_i \geq 0, \beta_i \geq 0$

Same as $\max_{\alpha_i} -\frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j x_i^T x_j + \sum \alpha_i + \max_{\beta_i \geq 0} -\sum_i \frac{(\alpha_i + \beta_i)^2}{4C}$

s.t. $\sum \alpha_i y_i = 0, \alpha_i \geq 0$

Same as $\max_{\alpha_i} -\frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j x_i^T x_j - \frac{\sum \alpha_i^2}{4C} + \sum \alpha_i \Rightarrow \beta_i = 0 \forall i$

s.t. $\alpha_i \geq 0, \sum_i \alpha_i y_i = 0$

Can be merged

Same as $\max_{\alpha_i} -\frac{1}{2} \sum \alpha_i \alpha_j y_i y_j k_2(x_i, x_j) + \sum \alpha_i$

s.t. $\alpha_i \geq 0, \sum_i \alpha_i y_i = 0$

where $k_2(x_i, x_j) = \begin{cases} x_i^T x_j + \frac{1}{2C} & y_i = y_j \\ x_i^T x_j & y_i \neq y_j \end{cases}$

Realize that exactly same as hard-margin SVM dual

only difference is gram matrix is incremented by $\frac{1}{2C} I$

Dual is also equal to:
$$\min_{\alpha} \frac{1}{2} \alpha^T Q_c \alpha - 1^T \alpha$$
 (vectorial notation)
$$\text{s.t. } \alpha \geq 0, y^T \alpha = 0$$
 (D)
 vector of 0's & 1's

where y is vector of labels of tr. datapts.

γ is $\text{diag}(y)$.

$$Q_c = \gamma \left(K + \frac{1}{2c} I \right) \gamma$$

gram matrix of tr. datapts

Note that Hessian of (D) is Q_c which is strictly positive definite $\forall c$.
 \therefore the objective is strictly convex & hence solution is unique.

→ recovery of w & b :

w can be (uniquely) computed ensuring $w = \sum \alpha_i y_i x_i$.

Once we get w , look at those eq. where $\alpha_i > 0$ let them be set $S \rightarrow$ set of support vectors.

$$\forall i \in S, \quad \underbrace{\epsilon_i}_{\text{Complementary slackness cond.}} = 1 + y_i b - \underbrace{y_i w^T x_i}_{\epsilon_i \geq 0} \geq 0$$

say set S_+
 $\Rightarrow \forall i \in S, y_i = 1$, we have $b \geq w^T x_i - 1$

say set S_-
 $\forall i \in S, y_i = -1$, we have $b \leq w^T x_i + 1$

$$\max_{i \in S_+} w^T x_i - 1 \leq b \leq \min_{i \in S_-} w^T x_i + 1$$

Q4 There are many ways of showing this is an indefinite kernel.

Here is one simple way:

1) If k is positive then Cauchy-Schwarz ^{inner} must hold

i.e. if Cauchy-Schwarz eq. does not hold then k is not positive.

2) In our case $k(x,y)^2 \leq k(x,x)k(y,y)$

would imply $(x^T y)^2 \leq 0 \quad \forall x,y$ which is impossible

~~(since it is the kernel)~~

$\therefore k$ is not positive. $\|Iy - k\|$ is also not positive, \therefore It is an indefinite kernel.

Q5 We will show that $\frac{k(A_1, A_2)}{|A_1 \cap A_2|}$ is a positive kernel. Then of course $|A_1 \cap A_2|^2$ is also positive kernel.

* [This kernel is actually very close to what we have seen in lecture]

Consider the prob. space $\mathbb{P} = (\Omega, 2^\Omega, P)$ $\rightarrow P$ is "uniform distribution" or "classical probability"

$$\text{i.e. } P(A) = \frac{|A|}{|\Omega|}$$

In this \mathbb{P} , we have $P(A_1 \cap A_2) = \frac{|A_1 \cap A_2|}{|\Omega|} = k(A_1, A_2)$.

Consider the Indicator r.v.: ~~mapping~~ ~~space~~ ~~to~~ ~~sets~~

$$X_{A_1} = \mathbb{1}_{A_1} \quad (\text{then } E\{X_{A_1}\} = P(A_1) = \frac{|A_1|}{|\Omega|})$$

Here, $E\{X_{A_1} X_{A_2}\} = P(A_1 \cap A_2) = k(A_1, A_2)$

\downarrow

we already know that, $E\{xy\}$ is a valid inner product in

Note that $E\{x^2\} \geq 0 \quad \forall X$

$\& E\{x^2\} = 0 \Rightarrow X = 0$ with prob. 1

(Non-negativity)

$\& E\{XY\} = E\{YX\}$

(Symmetry)

$\& E\{\alpha_1 X_1 + \alpha_2 X_2, Y\} = \alpha_1 E\{X_1, Y\} + \alpha_2 E\{X_2, Y\}$

(Linearity)

Hence $E\{XY\}$ is indeed a valid inner product.

Hence $k(A, A_2) = \frac{|A_1 \cap A_2|}{|A_1|}$ is a positive kernel.

In fact the vector space of all k 's defined on \mathcal{P} is an inner-product space where k is an inner-product.

Another argument:

Consider the kernel $k(A, A_2) = |A_1 \cap A_2|$

Now we will see "another" feature space where this is an inner-product.

Consider a Euclidean space in $|R|$ dimensions.

(one axis for each element of R)

For each $A \in \mathcal{P}$, consider the mapping of A to a ~~vector~~ ^{binary vector} of 0,1's such that a '1' is put ~~where~~ for all elements in A .

In this space, of course the natural dot product ~~gives~~ ^{is} between vectors ~~for~~ formed from A_1 & A_2 will give $|A_1 \cap A_2|$

⑥ Refer Scholkopf, Smola, Williamson, Bartlett. New Support Vector Algorithms, Neural Computation, 12, pp. 1207-1245, 2000.

⑦ Since it is mentioned that it is a real-world regression problem we can safely assume that the loss is bounded.

$$a \leq (y - f(x))^2 \leq b \quad \text{--- (1/2 mark)}$$

Now the learning bounds we derived in lectures using R.A. do hold for all bounded loss functions. Here the problem boils down to simply upper bounding the conditional R.A. of the indexed set of loss functions:

$$R_m(\mathcal{L}) = E_{\sigma} \left[\sup_{\|w\|_{SW} \leq R} \left\{ \frac{1}{m} \sum_{i=1}^m \sigma_i (y_i - w^T x_i)^2 \right\} \right] \quad \text{--- (1/2 mark)}$$

Now bounding $R_m(\mathcal{L})$ is not easy. (1/2 mark) will be given to any attempt (even if wrong) towards it. Here one (incomplete) way of bounding based on derivations in lectures:

$$\begin{aligned} R_m(\mathcal{L}) &\leq E_{\sigma} \left[\sup_{\|w\|_{SW} \leq R} \left\{ \frac{1}{m} \sum_i \sigma_i y_i^2 \right\} \right] + E_{\sigma} \left[\sup_{\|w\|_{SW} \leq R} \left\{ \frac{2}{m} \sum_i \sigma_i w^T x_i \right\} \right] \\ &\quad + E_{\sigma} \left[\sup_{\|w\|_{SW} \leq R} \left\{ \frac{1}{m} \sum_i \sigma_i (w^T x_i)^2 \right\} \right] \\ &\leq \frac{N^2}{m} E_{\sigma} \left[\left(\text{max. eigenvalue of } \sum_{i=1}^m \sigma_i x_i x_i^T \right)^2 \right] \end{aligned}$$

(using ideas given in lecture)

$$\leq \frac{2N}{m} \sqrt{\text{trace}(YX^T)}$$

This bound is fine if $\rightarrow 0$ as $m \rightarrow \infty$. Can we ensure that it happens?

Coming to the issue of a completely correct answer to this problem, refer to Cor. 6.7 in "Local Rademacher Complexities" by Bartlett, Bourquet & Mendelson, *Annals of Statistics*, 33(4):1497-1537, 2005.

⑧ Refer to related wikipedia articles.

- Applications:
- Ⓐ generalizations of χ^2
 - Ⓑ Efron-Stein inequality (Kiefer test.)
 - Ⓒ Ueda's lectures for proving McDiarmid ineq.

⑨ $f(X) = \text{sign} \left(\overbrace{1.5(X_1-1)^2 + 1.5(X_2-1)^2 + (X_1-1)(X_2-1) - 1}^{-g(X_1, X_2)} \right)$.

We require $P[-g(X_1, X_2) > 0] \stackrel{?}{\geq}$? → there is a typo in the problem!! we need a lower bound instead of upper bound

We have $E[g] = +1 - 1.5 - 1.5 \times 7.5 - (-0.5)(-1.5) = -5$

We can get $P[g > 0] \leq ?$ using McDiarmid's inequality of Hoeffding's inequality.

McDiarmid $P[g > 0] = P[g - E[g] > 5] \leq e^{-2(5)^2 / \sum_{i=1}^2 c_i^2} \rightarrow$ b.d. property

Now $\max_{X_1, X_1'} |g(X_1, X_2) - g(X_1', X_2)| = \max_{X_1, X_1'} |(X_1' - X_1)[1.5(X_1' + X_1) + X_2 - 4]|$

$\leq \begin{cases} 44 \rightarrow \text{obvious bound} \\ 21 \rightarrow \text{through mat lab.} \end{cases} = c_1$

||| by $c_2 = 12$

$\therefore P[f = 1] \geq 1 - e^{-2(5)^2 / (c_1^2 + c_2^2)} \approx 0.08192$

Hoeffding's Ineq.

$$\begin{aligned} \text{again } P\{g > 0\} &= P\{\overbrace{g - E(g)}^Z > 5\} \\ &= P\{e^{\lambda Z} > e^{\lambda 5}\} \leq \frac{E\{e^{\lambda Z}\}}{e^{\lambda 5}} \end{aligned}$$

We can show that

$$-19.5 \leq Z \leq 6$$

$$\leq e^{\frac{5\lambda}{8} (e^{5\lambda} - 5\lambda)}$$

$\forall \lambda > 0$

$$\therefore P\{B=1\} \geq 1 - e^{-5 \cdot (25.5)^{-1}} = 0.07401$$

again, since McDiarmid inequality employs information that the r.v.s X_1, X_2 are independent, it wins!