

1.) $H = \{ \text{sgn}(ax^2 + bx + c); a, b, c \in \mathbb{R} \}$

↳ 1 when positive
0 o/w

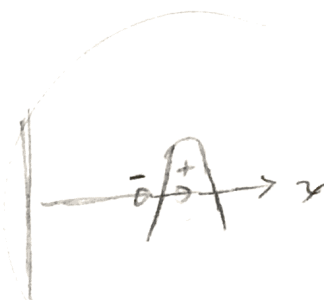
VC Dimension = 3

Proof...

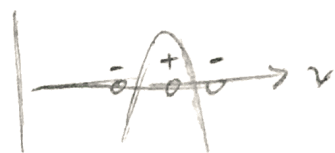
VC ≥ 1 :



VC ≥ 2 :

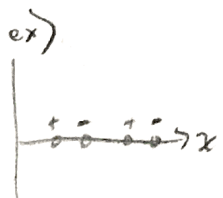


VC ≥ 3 : * all 3 same sign obvious extensions of above cases
* If 2 of same sign adjacent, simply place 3rd point c less than or greater than zero of parabola from corresponding example above
Final case: some sign points not adjacent...



Note, if choice of labels switched, simply flip parabola

VC < 4 : Consider 4 points arbitrarily placed c points along the x axis



such that adjacent points have alternating signs - note sign given to each arbitrary. If 2 adjacent points have alternating signs, they must be separated by zeros of the parabola. $\exists \leq 2$ zeros of the parabola, \therefore for any placement of the 4 points, \nexists enough zeros to separate the 3 sign switches. Therefore the given classifier can't shatter 4 points \Rightarrow VC Dimension = 3

\therefore sign switches 3 times btw points

2.) $K_B(x, z) = (1 + Bx \cdot z)^3 \mid B > 0; x, z \in \mathbb{R}^2$. Find $\phi_B(\cdot)$.

$$\hookrightarrow K_B(x, z) = \phi_B(x)^T \phi_B(z)$$

$$(1 + Bx \cdot z)^3 = (1 + Bx^T z)^3 = (1 + Bx_1 z_1 + Bx_2 z_2)^3$$

$$= B^3 x_1^3 z_1^3 + B^3 x_2^3 z_2^3 + 3B^3 x_1^2 z_1^2 z_2 + 3B^3 x_1 x_2^2 z_1 z_2^2 + 3B^2 x_1^2 z_1^2 + 3B^2 x_2^2 z_2^2 + 6B^2 x_1 x_2 z_1 z_2 + 3B x_1 z_1 + 3B x_2 z_2 + 1$$

found using Wolfram Alpha

$$= \begin{pmatrix} 1 \\ \sqrt{3B} x_1 \\ \sqrt{3B} x_2 \\ \sqrt{6B} x_1 x_2 \\ \sqrt{3B} x_1^2 \\ \sqrt{3B} x_2^2 \\ \sqrt{3B^{3/2}} x_1^2 x_2 \\ \sqrt{3B^{3/2}} x_1 x_2^2 \\ B^{3/2} x_1^3 \\ B^{3/2} x_2^3 \end{pmatrix}^T \begin{pmatrix} 1 \\ \sqrt{3B} z_1 \\ \sqrt{3B} z_2 \\ \sqrt{6B} z_1 z_2 \\ \sqrt{3B} z_1^2 \\ \sqrt{3B} z_2^2 \\ \sqrt{3B^{3/2}} z_1^2 z_2 \\ \sqrt{3B^{3/2}} z_1 z_2^2 \\ B^{3/2} z_1^3 \\ B^{3/2} z_2^3 \end{pmatrix}$$

$$\therefore \phi_B(x) =$$

$$\begin{pmatrix} 1 \\ \sqrt{3B} x_1 \\ \sqrt{3B} x_2 \\ \sqrt{6B} x_1 x_2 \\ \sqrt{3B} x_1^2 \\ \sqrt{3B} x_2^2 \\ \sqrt{3B^{3/2}} x_1^2 x_2 \\ \sqrt{3B^{3/2}} x_1 x_2^2 \\ B^{3/2} x_1^3 \\ B^{3/2} x_2^3 \end{pmatrix}$$

B adds a constant multiplier equal to $B^{1/2}$ (1/2 degree of the term) to each term of the feature vector & the kernel function of $K(x, z)$.

This allows you to weight the importance of higher order terms, in direct or indirect proportion to their order.

3.) a. $y w^T x = 1$

$$1 \Rightarrow 1(w_1, w_2) \cdot (1, 1) = 1$$

$$w_1 + w_2 = 1$$

$$2 \Rightarrow -1(w_1, w_2) \cdot (1, 0) = 1$$

$$w_1 = -1 \rightarrow w_2 = 2$$

$$\therefore w^* = \begin{bmatrix} -1 \\ 2 \end{bmatrix}$$

b. $y(w^T x + b) = 1$

$$1 \Rightarrow 1[(w_1, w_2)(1, 1) + b] = 1$$

$$w_1 + w_2 + b = 1$$

$$2 \Rightarrow -1[(w_1, w_2)(1, 0) + b] = 1$$

$$-w_1 - b = 1$$

$$w_2 = 2$$

$$\text{minimize } \|w\| \therefore w_1 = 0$$

$$\therefore b = 1$$

$$\Rightarrow w^* = \begin{bmatrix} 0 \\ 2 \end{bmatrix} \quad b^* = 1$$

Note $\|w_{b=0}^*\| = 5 < \|w_{b \neq 0}^*\| = 4$

\therefore allowing offset yields
larger margin

4.2.b) Maintaining Class Proportions in K-Fold Cross-Validation

It is important to maintain the proportions of each class across folds in K-fold validation —called stratification—because ideally, one wants each fold to be representative of the entire dataset. This is achieved by distributing the class instances to each fold proportionally to how they occur in the entire sample dataset...which hopefully is a good model for the underlying distribution.

4.2.d) Selecting the Optimal Hyperparameter

The table to the right shows the results of performing K-fold stratified cross validation on the parameter C of linear-kernel SVM classifiers, for three different classifiers—accuracy, F1-score, and area under the receiver operating characteristic curve (AUROC). We see that for judging by accuracy, the optimal $C = 10$ (Accuracy = 0.8300); by F1-score, the optimal $C = 10$ (F1-Score = 0.8839); by AUROC, the optimal $C = 1$ (AUROC = 0.8994). Note that for many of all three of the metrics, multiple C 's resulted in the same performance score. Accordingly, we chose the smallest such C with the max performance value, to get the largest margin. Typically, it is impossible to say which will improve test results the most, as this depends on the data you encounter during testing, but I found that any of the matching C 's actually yielded identical test performance.

C	accuracy	F1-score	AUROC
10^{-3}	0.7102	0.8106	0.8542
10^{-2}	0.7246	0.8170	0.8542
10^{-1}	0.8121	0.8762	0.8809
10^0	0.8211	0.8779	0.8994
10^1	0.8300	0.8839	0.8994
10^2	0.8300	0.8839	0.8994
best C			

4.3.c) Test Set Performance

Performance Metric Used	Value of C parameter	Test Performance Score
Accuracy	10	0.7391
F1_Score	10	0.4706
AUROC	1	0.7431