

(e) Express the generalization error (in terms of the training set error) as a function of the training set size.

4. [20 points] Properties of VC dimension

In this problem, we investigate a few properties of the Vapnik-Chervonenkis dimension, mostly relating to how $\text{VC}(H)$ increases as the set H increases. For each part of this problem, you should state whether the given statement is true, and justify your answer with either a formal proof or a counter-example.

- (a) Let two hypothesis classes H_1 and H_2 satisfy $H_1 \subseteq H_2$. Prove or disprove: $\text{VC}(H_1) \leq \text{VC}(H_2)$.
- (b) Let $H_1 = H_2 \cup \{h_1, \dots, h_k\}$. (I.e., H_1 is the union of H_2 and some set of k additional hypotheses.) Prove or disprove: $\text{VC}(H_1) \leq \text{VC}(H_2) + k$. [Hint: You might want to start by considering the case of $k = 1$.]
- (c) Let $H_1 = H_2 \cup H_3$. Prove or disprove: $\text{VC}(H_1) \leq \text{VC}(H_2) + \text{VC}(H_3)$.

5. [20 points] Training and testing on different distributions

In the discussion in class about learning theory, a key assumption was that we trained and tested our learning algorithms on the same distribution \mathcal{D} . In this problem, we'll investigate one special case of training and testing on different distributions. Specifically, we will consider what happens when the training labels are *noisy*, but the test labels are not.

Consider a binary classification problem with labels $y \in \{0, 1\}$, and let \mathcal{D} be a distribution over (x, y) , that we'll think of as the original, "clean" or "uncorrupted" distribution. Define \mathcal{D}_τ to be a "corrupted" distribution over (x, y) which is the same as \mathcal{D} , except that the labels y have some probability $0 \leq \tau < 0.5$ of being flipped. Thus, to sample from \mathcal{D}_τ , we would first sample (x, y) from \mathcal{D} , and then with probability τ (independently of the observed x and y) replace y with $1 - y$. Note that $\mathcal{D}_0 = \mathcal{D}$.

The distribution \mathcal{D}_τ models a setting in which an unreliable human (or other source) is labeling your training data for you, and on each example he/she has a probability τ of mislabeling it. Even though our training data is corrupted, we are still interested in evaluating our hypotheses with respect to the original, uncorrupted distribution \mathcal{D} .

We define the generalization error *with respect to* \mathcal{D}_τ to be

$$\varepsilon_\tau(h) = P_{(x,y) \sim \mathcal{D}_\tau}[h(x) \neq y].$$

Note that $\varepsilon_0(h)$ is the generalization error with respect to the "clean" distribution; it is with respect to ε_0 that we wish to evaluate our hypotheses.