

1. (a) The cardinality of \mathcal{H} is d
 \Rightarrow the num of mistakes $< \log_2 d$.
 Let $d=2^n$, and suppose h_d is the true learning rule.
 Let $x_i = (-1, \dots, -1, 1, \dots, 1)$ with $n/2$ -1 's and $n/2$ 1 's. Then $h_{\text{majority}}(x_i) = \pm 1$
 WLOG, $h_{\text{majority}}(x_i) = -1$. Then we get rid of predictors h_1, \dots, h_n .
 Let $x_{n+1} = (-1, \dots, -1, -1, \dots, -1, 1, \dots, 1)$ with $n/2$ -1 's, $n/2$ 1 's, and one -1 .
 WLOG, $h_{\text{majority}}(x_{n+1}) = -1$ and 2^{n-k} are killed.
 We then require $n-1$ mistakes before we are left w/ h_d .

$2^{n-1} + 2^{n-2} + \dots + 2^{n-(n-1)} = 2^n - 1$
 Thus, The worst-case num of mistakes is $\lceil \log_2 d \rceil - 1$.
 (b). For every $d=2^m$, we may make $2^2 + 2^2 + \dots + 2^2 + 2 + 2$ mistakes.

Base Case: Let $x_0^{(0)}, x_1^{(0)}$ be $(+1, \dots, +1), (-1, \dots, -1)$
 Let $x_0^{(1)} = (+1, \dots, +1, -1, \dots, -1)$
 $x_1^{(1)} = (-1, \dots, -1, +1, \dots, +1)$

Inductive Step: Sps we have $x^{(0)} = (x_0^{(0)}, x_1^{(0)})$, $x^{(1)} = (x_0^{(1)}, x_1^{(1)})$,
 \dots , $x^{(n)} = (x_0^{(n)}, x_1^{(n)}, \dots, x_{2^{n-1}}^{(n)})$
 define $x^{(n+1)}$ as follows: divide interval $[0,1]$ into 2^{n+1} subintervals. If an interval is colored, the corresponding coordinate is 1. If it is uncolored, it is -1.

Pair the sub-intervals. A sub-interval is on if it is colored on then off $[\text{---} | \text{---}]$, off if off then on: $[\text{---} | \text{---}]$. Interpret the elements of $x^{(0)}, \dots, x^{(n)}$ as binary, and generate $|x^{(0)}| + \dots + |x^{(n)}|$ elements. By turning on the corresponding sub-interval if that element is on.

We have generated the binary sequences of len 2^n whose half are on. Now we generate the rest.
 For each element we have just generated, take its binary sequence of switches. Flip the last switch, add these in sequence to our list. Then do the same but for the second to last. Continue until we've done this to all.
 We have now generated the elements whose switches are $2^n + 1$ on, and the elements whose $2^n - 1$ are on.
 Taking this new sublist, we generate the $2^n + 2$ on elements by turning on an off switch for each $2^n + 1$ element and generate the $2^n - 2$ on elements by turning off an on switch in the $2^n - 1$ on elements.

Continuing in this manner, we have successfully defined $x^{(n+1)}$, of which there are
 $\binom{2^n}{2^{n-1}} + \binom{2^n}{2^{n-1}-1} + \binom{2^n}{2^{n-1}+1} + \dots = 2^{2^n}$
 elements.

Now, we must see how the nearest neighbor algorithm makes a mistake on each new element added.
 The recursive elements are trivial to see this for, since the distance at every element in $x^{(n+1)}$ to one other previous is always $\sqrt{2^n}$, and the distance to each element added is a constant times the distance when the corresponding element was added. Therefore, there will be a tie between the correct and incorrect neighbors.

For the non-recursive ones, each added also ties and may make a mistake.

(c) Let \mathcal{D} be the dist over \mathcal{D} of the spherical Gaussian. Labels are $y = \text{sgn}(x[i])$ for a fixed i .
 Let h_m be a nearest neighbor predictor based on $S \sim \mathcal{D}^m$ of m iid samples.

Show that $\exists c \in (0,1)$ s.t. for $m=2^{cd}$
 $E(L(h_m)) > 0.4$.

In fact, $L(h_m) > 0.4$ with high probability.

Proof: Recall: If $x \sim \mathcal{D}$ spherical gaussian on \mathbb{R}^d mean 0, unit variance. Then for all $\beta \leq \sqrt{d}$
 $P(\sqrt{d} - \beta \leq \|x\| \leq \sqrt{d} + \beta) \geq 1 - 3e^{-\beta^2/4}$

For some $c > 0$.

Thus for a, b drawn from a spherical gaussian,
 $\|a\| \approx \sqrt{d} \pm O(1)$
 $\|b\| \approx \sqrt{d} \pm O(1)$

Additionally points drawn at random are nearly orthogonal, so, by Pythagoras
 $\|a-b\|^2 \approx 2(\sqrt{d} \pm O(1))^2 = 2d + O(\sqrt{d})$.

Thus vectors picked spherical gaussian are also approximately orthogonal.

2. Parzen Window Predictor.

Let $S = \{(x_i, y_i)\}_{i=1}^m$. The Parzen Window (kernel) Density Estimate $\hat{f}(x|Y=y)$ of $f(x|Y=y)$, for each $y \in \mathcal{Y}$ is.

$$\hat{f}(x|y) := Z_y \sum_{i: y_i=y} K(x, x_i)$$

where $Z_y = (\int_{\mathcal{X}} \sum_{i: y_i=y} K(x, x_i) dx)^{-1}$, and

$$K(x, x_i) = \exp(-p(x, x_i)^2 / \sigma^2)$$

where p is a translation invariant metric on \mathcal{X} , σ is a hyper-parameter.

The Parzen Window estimator $\hat{D}(X, y)$ is given by $\hat{f}(x|y)$ combined with $\hat{p}(y) := \frac{1}{m} |\{i | y_i=y\}|$, or the count at num of diff labels among the m samples.

The Parzen Predictor is the Bayes' Optimal Predictor for $\hat{D}(X, y)$. It is specified by choice of $p(x, x')$ and σ .

(a) Show that for $\mathcal{Y} = \{-1, +1\}$, the Parzen Predictor is given by:

$$h(x) := \text{sign}\left(\sum_{i=1}^m y_i K(x, x_i)\right)$$

proof: Recall:

$$h_{\text{BayesCD}}(x) = \text{sign}(\eta_{\hat{D}}(x) - \frac{1}{2}).$$

and the Parzen Predictor is given by

$$h_{\text{BayesCD}}(x) = \text{sign}(\eta_{\hat{D}}(x) - \frac{1}{2}).$$

We have, by def'n of the estimator of the joint density.

$$\hat{f}(x, y) := \hat{f}(x|y) \cdot \hat{p}(y).$$

Recall:

$$\eta_{\hat{D}}(x) := P(Y=1 | X=x).$$

$$\therefore \eta_{\hat{D}}(x) = P(Y=1 | X=x) = \frac{\hat{f}(x, 1)}{\hat{f}(x, 1) + \hat{f}(x, -1)}$$

$$E_y[\hat{f}(x|y)] = \hat{f}_x(x).$$

$$E_y[\hat{f}(x|y)] |_{X=x}.$$

$$E_y[\hat{f}(x|y)] = \hat{f}(x|+1) \cdot \underbrace{\hat{p}(+1)}_{=p} + \hat{f}(x|-1) \hat{p}(-1)$$

$$= p Z_+ \sum_{i: y_i=1} K(x, x_i) + (1-p) Z_- \sum_{i: y_i=-1} K(x, x_i)$$

Recall: p shift-invariant implies

$$\int_{\mathcal{X}} K(x, x_i) dx = C \quad \forall i.$$

$$m \cdot \hat{p}(y).$$

$$\therefore \int_{\mathcal{X}} Z_y \sum_{i: y_i=y} K(x, x_i) dx = Z_y \cdot C \cdot |\{i: y_i=y\}| = 1.$$

$$\Rightarrow Z_y = \frac{1}{m \cdot \hat{p}(y) \cdot C}.$$

Thus.

$$\hat{f}_x(x) = \frac{p}{m \cdot p \cdot C} \sum_{i: y_i=1} K(x, x_i) + \frac{(1-p)}{m \cdot (1-p) \cdot C} \sum_{i: y_i=-1} K(x, x_i).$$

$$= \frac{1}{m \cdot C} \sum_{i=1}^m K(x, x_i)$$

$$\therefore \eta_{\hat{D}}(x) = \frac{\frac{p}{m \cdot p \cdot C} \sum_{i: y_i=1} K(x, x_i)}{\frac{1}{m \cdot C} \sum_{i=1}^m K(x, x_i)} = \frac{\sum_{i: y_i=1} K(x, x_i)}{\sum_{i=1}^m K(x, x_i)}$$

$$\therefore \eta_{\hat{D}}(x) - \frac{1}{2} = \frac{2 \sum_{i: y_i=1} y_i K(x, x_i) - \sum_{i=1}^m K(x, x_i)}{2 \sum_{i=1}^m K(x, x_i)}$$

$$= \frac{\sum_{i: y_i=1} y_i K(x, x_i) + \sum_{i: y_i=-1} y_i K(x, x_i)}{2 \sum_{i=1}^m K(x, x_i)}$$

$$= \frac{\sum_{i=1}^m y_i K(x, x_i)}{2 \sum_{i=1}^m K(x, x_i)} \geq 0$$

$$\Rightarrow h(x) = \text{sign}(\eta_{\hat{D}}(x) - \frac{1}{2}) = \text{sign}\left(\sum_{i=1}^m y_i K(x, x_i)\right)$$

(b) Using $K(x, x') = \exp(-p(x, x')^2 / \sigma^2)$, how does the Parzen predictor h behave in the limit as $\sigma \rightarrow \infty$?

Soln: as $\sigma \rightarrow \infty$, we have

$$\exp(-p(x, x')^2 / \sigma^2) \rightarrow \exp(0) = 1, \text{ so for every } x$$

$$h(x) = \text{sign}\left(\sum_{i=1}^m y_i K(x, x_i)\right) \rightarrow \text{sign}\left(\sum_{i=1}^m y_i\right).$$

i.e. Constant majority predictor.

(c) how does h behave as $\sigma \rightarrow 0$?

Soln: $\exp(-p(x, x')^2 / \sigma^2) \rightarrow \exp(-\infty) = 0$, if $x \neq x'$

if $x = x'$, $p(x, x') = 0$, so $\exp(-p(x, x')^2 / \sigma^2) \rightarrow \exp(0) = 1$.

$$\text{Thus } h(x) \rightarrow \text{sign}\left(\sum_{i=1}^m y_i \mathbb{1}\{x_i = x\}\right)$$

Memorization.

(d) When σ is small, $p(x, x')^2 / \sigma^2$ is large, except when x' is near to x . Thus, when predicting on x ,

we bias the labels s.t. x_i is near to x .

When σ grows larger, this effect is less pronounced.

In effect, this is a weighted m -NN predictor with weights $K(x, x_i)$.

Whenever one ties, we pick arbitrarily, in accordance with

$$\text{sign}(0) = 1.$$

3. Nearest Neighbor in the Statistical Setting.

(a) $X \sim \text{Unif}[-1, 1]$, $Y = \pm 1$, $P_D(Y = +1 | X) = 0.5 + 0.3 \text{sgn}(X)$

What is the Bayes Optimal Predictor h^* and the Bayes Error
Sol'n: $P_D(X)$ is precisely our $\eta(x)$, our posterior dist. $L_D(h^*)$?
 Therefore

$$\begin{aligned} h^*(x) &= \text{sgn}(P_D(X) - \frac{1}{2}) = \text{sgn}(0.3 \text{sgn}(X)) = \text{sgn}(X) \\ L_D(h^*) &= P_D(h^*(X) \neq Y^*) \\ &= P(Y \neq +1 | X \geq 0) \cdot P(X \geq 0) + P(Y \neq -1 | X < 0) \cdot P(X < 0) \\ &= P(Y = -1 | X \geq 0) \cdot \frac{1}{2} + P(Y = +1 | X < 0) \cdot \frac{1}{2} \\ &= \frac{1}{2}(1 - 0.5 - 0.3) + \frac{1}{2}(0.5 - 0.3) \\ &= \frac{1}{2} \cdot 0.8 + \frac{1}{2} \cdot 0.2 = \frac{1}{2}. \end{aligned}$$

(b) How does the error of the nearest neighbor predictor behave, $L(h_m)$ as $m \rightarrow \infty$?

Sol'n: Given $x \in [-1, 1]$, Let $Z = |X - x|$, $X \sim \text{Unif}(-1, 1)$.

The CDF of Z is:

$$\begin{aligned} P(|X - x| \leq t) &= \frac{\mu([-1, 1] \cap [x-t, x+t])}{\mu([-1, 1])} = \frac{\mu([\max\{x-t, -1\}, \min\{x+t, 1\}])}{\mu([-1, 1])} \\ &= \frac{\min\{1, x+t\} - \max\{-1, x-t\}}{2} \\ &= \frac{1 + (x+t) - 1 - (x-t)}{2} + \frac{1 - (x-t) + (x-t+1)}{2} \\ &= \frac{1}{2} \left(t+1 - \frac{|1 - (x+t)| + |x-t+1|}{2} \right) \end{aligned}$$

Let $Z_1, \dots, Z_m \stackrel{\text{iid}}{\sim} Z$, $Z_{(1)} := \min_i Z_i$

Then, the CDF of $Z_{(1)}$ is:

$$\begin{aligned} P(Z_{(1)} \leq t) &= 1 - [1 - P(Z \leq t)]^m \\ &= 1 - \left[1 - \frac{1}{2} \left(t+1 - \frac{|1 - (x+t)| + |x-t+1|}{2} \right) \right]^m \\ &= 1 - \left[\frac{1}{2} - \frac{1}{2}t - \frac{|1 - (x+t)| + |x-t+1|}{4} \right]^m \\ &= 1 - 2^{-m} \left[\frac{1}{2} - \frac{1}{2}t - \frac{|1 - (x+t)| + |x-t+1|}{4} \right]^m \\ &\approx 1 - (1-t)^m \end{aligned}$$

$$\Rightarrow P(Z_{(1)} \geq t) \approx (1-t)^m$$

Finally,

$$P(h_m(x) \neq \text{sgn}(x)) \leq P(Z_{(1)} \geq |x|) \approx (1-|x|)^m$$

\therefore as $m \rightarrow \infty$, $h_m(x)$ converges

to $h_{\text{Bayes}}(x)$ almost surely

since the set of points on which h_m does not converge to $\text{sgn}(x)$ is \emptyset .