

# FINAL: 90 Minutes

Last Name: Solutions  
First Name:   
RIN:   
Section: 4100 / 6100 (circle one)

Answer **ALL** questions.

**NO COLLABORATION** or electronic devices. Any violations result in an F.

**NO questions** allowed during the test. Interpret and do the best you can.

**ALWAYS** show your work and justify each answer.

## GOOD LUCK!

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	TOTAL

- You do not have time to waffle. So, don't waffle.
- Keep your answers precise and concise.
- Each question is worth 1 point.

1. What is the two step approach to learning and why do we do it that way?

①  $E_{in} \approx E_{out}$

②  $E_{in} \approx 0$

Need to do this way because we don't know  $E_{out}$

→ Can only get  $E_{in} \approx 0$

→ Learning means  $E_{out} \approx 0$

→ Need to link  $E_{in}$  to  $E_{out}$  in step 1.

2. You're a Facebook Troll. All your FB-friends have exciting lives based on their posts. In comparison, your life seems mundane and you get depressed. What learning from data trap could you be in?

Sampling bias or wrong bin

Bin 1 → Highlight Reel → estimate  $P[\text{Exciting}]$

Bin 2 → Normal Reel → estimate  $P[\text{Exciting}]$

Can't use  $P[\text{Exciting}]$  from bin 1 for bin 2.

3. Why is it a good idea to preprocess the data with input normalization?

To correct for arbitrary choices made during data collection such as units to measure (say) income.

4. Derive the growth function for the 1-dimensional positive intervals. Explain how to get the VC-dimension, and give the VC-dimension.

... n points  $\binom{N+1}{2}$  ways to place the two endpoints of the interval in two different spots.

Each gives a different dichotomy with at least 1 positive point.

Put both endpoints in same interval  $\rightarrow$  all negative.

$$m_H(N) = 1 + \binom{N+1}{2}$$

N	1	2	3
$m_H(N)$	2	4	7

$\swarrow$   $d_{VC}$        $\swarrow$  cheapest break point

VC-dimension is the max # points  $\mathcal{H}$  can shatter  
 $\rightarrow$  Cheapest break point - 1.

$$d_{VC} = 2$$

5. Based on your VC-dimension for positive intervals, what is the error bar linking  $E_{in}$  to  $E_{out}$ ?

$$E_{out} \leq E_{in} + O\left(\sqrt{\frac{d_{VC} \ln N}{N}}\right) \quad d_{VC} = 2$$

$\underbrace{\hspace{10em}}_{\text{error bar}} \quad O\left(\sqrt{\frac{2 \ln N}{N}}\right)$

More  
Exact.

$$E_{out} \leq E_{in} + \sqrt{\frac{8}{N} \ln 4 m_H(2n)}$$

$$\leq \sqrt{\frac{8}{N} d_{VC} \ln 2N + O(1)}$$

6. You decided to use 8th order polynomials to fit the data, to give your learning lots of flexibility. You suspect you might overfit the data. What does that mean?

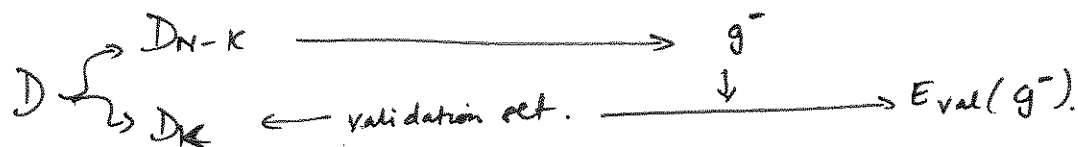
Model may fit noise in data (stochastic or deterministic)  
and thereby produce an inferior hypothesis with bad  $E_{out}$ .

7. In the previous problem, what could you do to help the learning in case of overfitting?

Regularize.

8. What is a validation set. Why do we use it?

Validation set gives peek at  $E_{out}$  with  $E_{val}$ .



Validation can be used to (for example) pick regularization parameter  $\lambda$ .

9. What are the tradeoffs in choosing the validation set? Why should it be small? Why large?

$K$  large  $\rightarrow$   $E_{val}$  is more accurate

$K$  small  $\rightarrow g^- \approx g \therefore E_{val}(g)$  is a better estimate of  $E_{out}(g)$

10. Define the leave one out cross-validation error,  $E_{cv}$ ?

$$E_{cv} = \frac{1}{N} \sum_{i=1}^N e_i(g_i^-, y_i).$$

where  $g_i^-$  is the deficient hypothesis obtained from data minus  $(x_i, y_i)$ .

$e_i$  is your error function

$e_i(g_i^-(x_i), y_i)$  is the error between your deficient predictor on the left out point.

11. Prove that  $E_{cv}$  based on  $N$  data points is an unbiased estimate of your expected out-of-sample error when you learn on  $N-1$  data points?

$$\begin{aligned} E_D[E_{cv}] &= E_D \left[ \frac{1}{N} \sum_{i=1}^N e_i(g_i^-(x_i), y_i) \right] \\ &= \frac{1}{N} \sum_{i=1}^N E_D e_i(g_i^-(x_i), y_i) \\ &= \frac{1}{N} \sum_{i=1}^N E_{D_i^-} \left[ \underbrace{E_{x_i, y_i}[e_i]}_{E_{out}(g_i^-)} \right] \\ &= \frac{1}{N} \sum_{i=1}^N E^{(N-1)} = \underline{\underline{E^{(N-1)}}}. \end{aligned}$$

12. The nearest neighbor algorithm was at most a factor of 2 away from optimal. Prove that the 3-nearest neighbor algorithm is near-optimal. That is, the out of sample error at a test point  $x$  is bounded by

$$E_{\text{out}}(3\text{-NN}) \leq E_{\text{out}}^*(1 + 3E_{\text{out}}^*).$$

If optimal performance is 1% error, how bad can 3-NN be? (In your proof, state any assumptions.)

Assume.  $\pi(x) = P[y=1|x]$  is smooth

$$x_{[1]} x_{[2]} x_{[3]} \rightarrow x \text{ as } N \rightarrow \infty.$$

$$\pi(x_{[i]}) \rightarrow \pi(x)$$

$$P[\text{error}] = \frac{P[\text{error} | y=1] \pi(x)}{3(1-\pi)^2 \pi + (1-\pi)^3} + \frac{P[\text{error} | y=-1] (1-\pi(x))}{3\pi^2 (1-\pi) + \pi^3}.$$

↖ symmetric in  $\pi, 1-\pi$

$$E_{\text{out}}^* = \min(\pi, 1-\pi)$$

∴ we can replace  $\pi$  by  $E_{\text{out}}^*$ .

$$\begin{aligned} P[\text{error}] &= E_{\text{out}}^* \pi(1-\pi) [3(1-\pi)\pi + (1-\pi)^2] + \pi(1-\pi) [3\pi(1-\pi) + \pi^2] \\ &= \pi(1-\pi) [6\pi(1-\pi) + (1-\pi)^2 + \pi^2] \\ &= E_{\text{out}}^* (4\pi(1-\pi) + 1) \cdot \pi(1-\pi) \end{aligned}$$

$$= E_{\text{out}}^* (1 - E_{\text{out}}^*) [1 + 4E_{\text{out}}^* (1 - E_{\text{out}}^*)]$$

$$= E_{\text{out}}^* (1 + 3E_{\text{out}}^* - (E_{\text{out}}^*)^2 + 4(E_{\text{out}}^*)^3)$$

$$E_{\text{out}}^* \leq \frac{1}{2}$$

$$\leq E_{\text{out}}^* (1 + 3E_{\text{out}}^*)$$

$$E_{\text{out}}^* = 0.01 \rightarrow P[\text{error}] \leq 0.01 \cdot (1.03) = 0.0103 = 1.03\%$$

3-NN asymptotically has  $E_{\text{out}} \leq 1.03\%$ ; pretty close to optimal.

3-NN is enough!

13. Explain why the 1-hidden layer neural network is more powerful than the K-RBF network. What are the pros and cons of this power?

K-RBF Network:  $g(x) = \text{sign} \left( w_0 + \sum_{i=1}^K w_i \phi \left( \frac{x - \mu_i}{\sigma_i} \right) \right)$  ↖ similarity features fixed via unsupervised learning

NN:  $g(x) = \text{sign} \left( w_0 + \sum_{i=1}^m w_i \phi(v_i^T x) \right)$  ↖ fully tunable features,  $v_i$  and  $w_i$  are jointly learned.

KRBF network:  $\mu_j \rightarrow$  unsupervised  $\rightarrow$  linear model given  $\mu_j$

NN: fully tuned features  $\rightarrow$  more powerful.

14. Explain why the optimal hyperplane with maximum margin performs well even in very high dimensional spaces where a random perceptron won't.

$E_{CV} \approx E_{out}(N-1)$  is controlled by # support vectors which can be small even in high dimensions.

15. Why is it computationally tractable to run the optimal hyperplane using a feature transform into very high, even infinite, dimensions.

In dual space it is an Inner Product Algorithm  
 $\rightarrow$  can work in original  $X$ -space if given the kernel that computes the inner product.

↑  
 Never have to go to the infinite feature space.