

RECITATION 6

LEARNING THEORY, PROBABILISTIC LEARNING, PRECISION AND RECALL

10-301/10-601: INTRODUCTION TO MACHINE LEARNING
10/25/24

1 Learning Theory

1.1 PAC Learning

Some Important Definitions

1. Basic notation:

- Probability distribution (unknown): $X \sim p^*$
- **True function** (unknown): $c^* : X \rightarrow Y$
- **Hypothesis space** \mathcal{H} and **hypothesis** $h \in \mathcal{H} : X \rightarrow Y$
- Training dataset $\mathcal{D} = \{x^{(1)}, \dots, x^{(N)}\}$

2. **True Error (expected risk)**

$$R(h) = P_{x \sim p^*(x)}(c^*(x) \neq h(x))$$

3. **Train Error (empirical risk)**

$$\begin{aligned}\hat{R}(h) &= P_{x \sim \mathcal{D}}(c^*(x) \neq h(x)) \\ &= \frac{1}{N} \sum_{i=1}^N \mathbb{1}(c^*(x^{(i)}) \neq h(x^{(i)})) \\ &= \frac{1}{N} \sum_{i=1}^N \mathbb{1}(y^{(i)} \neq h(x^{(i)}))\end{aligned}$$

The **PAC criterion** is that we produce a high accuracy hypothesis with high probability. More formally,

$$P(\forall h \in \mathcal{H}, \text{_____} \leq \text{_____}) \geq \text{_____}$$

$$P(\forall h \in \mathcal{H}, |R(h) - \hat{R}(h)| \leq \epsilon) \geq 1 - \delta$$

Sample Complexity is the minimum number of training examples N such that the PAC criterion is satisfied for a given ϵ and δ

Sample Complexity for 4 Cases: See Figure 1. Note that

- **Realizable** means $c^* \in \mathcal{H}$
- **Agnostic** means c^* may or may not be in \mathcal{H}

	Realizable	Agnostic
Finite $ \mathcal{H} $	Thm. 1 $N \geq \frac{1}{\epsilon} [\log(\mathcal{H}) + \log(\frac{1}{\delta})]$ labeled examples are sufficient so that with probability $(1 - \delta)$ all $h \in \mathcal{H}$ with $\hat{R}(h) = 0$ have $R(h) \leq \epsilon$.	Thm. 2 $N \geq \frac{1}{2\epsilon^2} [\log(\mathcal{H}) + \log(\frac{2}{\delta})]$ labeled examples are sufficient so that with probability $(1 - \delta)$ for all $h \in \mathcal{H}$ we have that $ R(h) - \hat{R}(h) \leq \epsilon$.
Infinite $ \mathcal{H} $	Thm. 3 $N = O(\frac{1}{\epsilon} [\text{VC}(\mathcal{H}) \log(\frac{1}{\epsilon}) + \log(\frac{1}{\delta})])$ labeled examples are sufficient so that with probability $(1 - \delta)$ all $h \in \mathcal{H}$ with $\hat{R}(h) = 0$ have $R(h) \leq \epsilon$.	Thm. 4 $N = O(\frac{1}{\epsilon^2} [\text{VC}(\mathcal{H}) + \log(\frac{1}{\delta})])$ labeled examples are sufficient so that with probability $(1 - \delta)$ for all $h \in \mathcal{H}$ we have that $ R(h) - \hat{R}(h) \leq \epsilon$.

12

Figure 1: Sample Complexity for 4 Cases

The **VC dimension** of a hypothesis space \mathcal{H} , denoted $\text{VC}(\mathcal{H})$ or $d_{VC}(\mathcal{H})$, is the maximum number of points such that there exists at least one arrangement of these points and a hypothesis $h \in \mathcal{H}$ that is consistent with any labelling of this arrangement of points.

To show that $\text{VC}(\mathcal{H}) = n$:

- Show there exists a set of points of size n that \mathcal{H} can shatter
- Show \mathcal{H} cannot shatter any set of points of size $n + 1$

Questions

- For the following examples, write whether or not there exists a dataset with the given properties that can be shattered by a linear classifier.
 - 2 points in 1D
 - 3 points in 1D
 - 3 points in 2D
 - 4 points in 2D

How many points can a linear boundary (with bias) classify exactly for d-Dimensions?

- Yes
- No
- Yes
- No

$$d + 1$$

2. Consider a rectangle classifier (i.e. the classifier is uniquely defined 3 points $x_1, x_2, x_3 \in \mathbb{R}^2$ that specify 3 out of the four corners), where all points within the rectangle must equal 1 and all points outside must equal -1

(a) Which of the configurations of 4 points in figure 2 can a rectangle shatter?

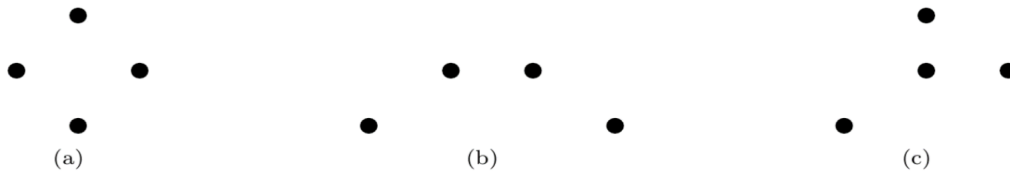


Figure 2

(a), (b), since the rectangle can be scaled and rotated it can always perfectly classify the points. (c) is not perfectly classifiable in the case that all the exterior points are positive and the interior point is negative.

(b) What about the configurations of 5 points in figure 3?



Figure 3

None of the above. For (d), consider (from left to right) the labeling 1, 1 -1, -1, 1. For (e), same issue as (c).

3. In the below table, state in which case the sample complexity of the hypothesis falls under.

Problem	Hypothesis Space	Realizable/ Agnostic	Finite/ Infinite																				
A binary classification problem, where the data points are linearly separable	Set of all linear classifiers																						
Predict whether it will rain or not based on the following dataset: <table border="1" style="margin: 10px auto;"> <thead> <tr> <th>Temp</th> <th>Humid</th> <th>Wind</th> <th>Rain?</th> </tr> </thead> <tbody> <tr> <td>High</td> <td>Yes</td> <td>Yes</td> <td>Yes</td> </tr> <tr> <td>Low</td> <td>Yes</td> <td>No</td> <td>No</td> </tr> <tr> <td>Low</td> <td>No</td> <td>Yes</td> <td>Yes</td> </tr> <tr> <td>High</td> <td>No</td> <td>No</td> <td>Yes</td> </tr> </tbody> </table>	Temp	Humid	Wind	Rain?	High	Yes	Yes	Yes	Low	Yes	No	No	Low	No	Yes	Yes	High	No	No	Yes	A decision tree with max depth 2, where each node can only split on one feature, and the features cannot be repeated along a branch		
Temp	Humid	Wind	Rain?																				
High	Yes	Yes	Yes																				
Low	Yes	No	No																				
Low	No	Yes	Yes																				
High	No	No	Yes																				
Classifying a set of real-valued points where the underlying data distribution is unknown	Set of all linear classifiers																						
A binary classification problem on a given set of data points, where the data is not linearly separable	K-nearest neighbour classifier with Euclidean distance as distance metric																						

	Realizable/ Agnostic	Finite/ Infinite
1	Realizable	Infinite (All possible linear classifiers)
2	Realizable (We can split the given data using a depth 2 decision tree)	Finite (There are only a finite set of decision trees that can be formed with the given constraints)
3	Agnostic (The data may or may not be linearly separable)	Infinite
4	Agnostic (The KNN classifier may or not be able to perfectly classify each point)	Finite (The hypothesis space is the set of all possible partitions of the input space into k-nearest regions - which is finite for all possible values of k)

4. Let x_1, x_2, \dots, x_n be n random variables that represent binary literals ($x \in \{0, 1\}^n$). Let the hypothesis class \mathcal{H}_n denote the conjunctions of no more than n literals in which each variable occurs at most once. Assume that $c^* \in \mathcal{H}_n$.

Example: For $n = 4$, $(x_1 \wedge x_2 \wedge x_4), (x_1 \wedge \neg x_3) \in \mathcal{H}_4$

Find the minimum number of examples required to learn $h \in \mathcal{H}_{10}$ which guarantees at least 99% accuracy with at least 98% confidence.

$$|\mathcal{H}_n| = 3^n$$

$$|\mathcal{H}_{10}| = 3^{10}, \epsilon = 0.01, \delta = 0.02$$

$$N(\mathcal{H}_{10}, \epsilon, \delta) \geq \lceil \frac{1}{\epsilon} [\ln |\mathcal{H}_{10}| + \ln \frac{1}{\delta}] \rceil = \lceil 1489.81 \rceil = 1490$$

2 Probabilistic Learning

In probabilistic learning, we are trying to learn a target probability distribution as opposed to a target function. We'll review two ways of estimating the parameters of a probability distribution, as well as one family of probabilistic models: Naive Bayes classifiers.

2.1 MLE/MAP

As a reminder, in MLE, we have

$$\begin{aligned}\hat{\theta}_{MLE} &= \arg \max_{\theta} p(\mathcal{D}|\theta) \\ &= \arg \min_{\theta} -\log(p(\mathcal{D}|\theta))\end{aligned}$$

For MAP, we have

$$\begin{aligned}\hat{\theta}_{MAP} &= \arg \max_{\theta} p(\theta|\mathcal{D}) \\ &= \arg \max_{\theta} \frac{p(\mathcal{D}|\theta)p(\theta)}{\text{Normalizing Constant}} \\ &= \arg \max_{\theta} p(\mathcal{D}|\theta)p(\theta) \\ &= \arg \min_{\theta} -\log(p(\mathcal{D}|\theta)p(\theta))\end{aligned}$$

-
1. Imagine you are a data scientist working for an advertising company. The advertising company has recently run an ad and wants you to estimate its performance.

The ad was shown to N people. Let $Y^{(i)} = 1$ if person i clicked on the ad and 0 otherwise. Thus $\sum_i^N y^{(i)} = k$ people decided to click on the ad. Assume that the probability that the i -th person clicks on the ad is θ and the probability that the i -th person does not click on the ad is $1 - \theta$.

(a) Note that

$$p(\mathcal{D}|\theta) = p((Y^{(1)}, Y^{(2)}, \dots, Y^{(N)}|\theta) = \theta^k(1 - \theta)^{N-k}$$

Calculate $\hat{\theta}_{MLE}$.

$$\begin{aligned}\hat{\theta}_{MLE} &= \arg \min_{\theta} -\log(p(\mathcal{D}|\theta)) \\ &= \arg \min_{\theta} -\log(\theta^k(1 - \theta)^{N-k}) \\ &= \arg \min_{\theta} -k * \log(\theta) - (N - k) \log(1 - \theta)\end{aligned}$$

Setting the derivative equal to zero yields

$$\begin{aligned}0 &= \frac{-k}{\theta} + \frac{(N - K)}{1 - \theta} \\ \implies \hat{\theta}_{MLE} &= \frac{k}{N}\end{aligned}$$

(b) Suppose $N = 100$ and $k = 10$. Calculate $\hat{\theta}_{MLE}$.

$$\hat{\theta}_{MLE} = \frac{k}{N} = 0.10$$

(c) Your coworker tells you that $\theta \sim \text{Beta}(\alpha, \beta)$. That is:

$$p(\theta) = \frac{\theta^{\alpha-1}(1 - \theta)^{\beta-1}}{B(\alpha, \beta)}$$

Recall from lecture that $\hat{\theta}_{MAP}$ for a Bernoulli random variable with a Beta prior is given by:

$$\hat{\theta}_{MAP} = \frac{k + \alpha - 1}{N + \alpha + \beta - 2}$$

Suppose $N = 100$ and $k = 10$. Furthermore, you believe that in general people click on ads about 6 percent of the time, so you, somewhat naively, decide to set $\alpha = 6 + 1 = 7$, and $\beta = 100 - 6 + 1 = 95$. Calculate $\hat{\theta}_{MAP}$.

$$\hat{\theta}_{MAP} = \frac{k + \alpha - 1}{N + \alpha + \beta - 2} = \frac{10 + 7 - 1}{100 + 102 - 2} = \frac{16}{200} = 0.08$$

(d) How do $\hat{\theta}_{MLE}$ and $\hat{\theta}_{MAP}$ differ in this scenario? Argue which estimate you think is better.

Both estimates are reasonable given the available information. Note that $\hat{\theta}_{MAP}$ has lower variance than $\hat{\theta}_{MLE}$, but $\hat{\theta}_{MAP}$ is more biased. If you believe that this advertisement is similar to advertisements with a 6 percent click rate, then $\hat{\theta}_{MAP}$ may be a superior estimate, but if the circumstances under which the advertisement was shown were different from the usual, then $\hat{\theta}_{MLE}$ might be a better choice.

2. Suppose you are an avid Neural and Markov fan who monitors the @neuralthenarwhal Instagram account each day. Suppose you wish to find the probability that Neural or Markov will post at any time of day. Over three days you look on Instagram and find the following number of new posts: $x = [3, 4, 1]$

A fellow fan tells you that this comes from a Poisson distribution:

$$p(x|\theta) = \frac{e^{-\theta}\theta^x}{x!}$$

Also, you are told that $\theta \sim \text{Gamma}(2, 2)$ — that is, its pdf is:

$$p(\theta) = \frac{1}{4}\theta e^{-\frac{\theta}{2}}, \theta > 0$$

Calculate $\hat{\theta}_{MAP}$.

(See also https://en.wikipedia.org/wiki/Conjugate_prior)

Note:

$$p(\mathcal{D}|\theta) = \frac{e^{-\theta}\theta^3}{3!} \frac{e^{-\theta}\theta^4}{4!} \frac{e^{-\theta}\theta^1}{1!}$$

$$\begin{aligned} \hat{\theta}_{MAP} &= \arg \min_{\theta} -\log(p(\mathcal{D}|\theta)p(\theta)) \\ &= \arg \min_{\theta} -\log\left(\frac{e^{-\theta}\theta^3}{3!} \frac{e^{-\theta}\theta^4}{4!} \frac{e^{-\theta}\theta^1}{1!} \times \frac{1}{4}\theta e^{-\frac{\theta}{2}}\right) \\ &= \arg \min_{\theta} -\log\left(\frac{e^{-3\theta-\frac{\theta}{2}}\theta^9}{3! \times 4!}\right) \\ &= \arg \min_{\theta} -\left(\left(-3\theta - \frac{\theta}{2}\right) \log e + 9 \log \theta - \log(3! \times 4!)\right) \\ &= \arg \min_{\theta} \left(3\theta + \frac{\theta}{2}\right) - 9 \log \theta + \log(3! \times 4!) \end{aligned}$$

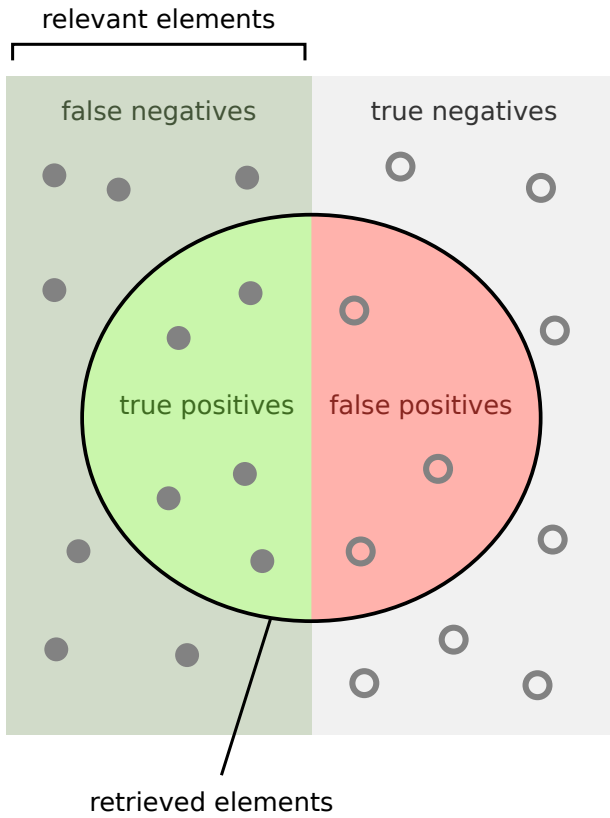
Taking the derivative gives us

$$\frac{d}{d\theta} \left(3\theta + \frac{\theta}{2}\right) - 9 \log \theta + \log(3! \times 4!) = \left(3 + \frac{1}{2}\right) - \frac{9}{\theta}$$

Setting the derivative equal to zero yields

$$\begin{aligned} 0 &= \left(3 + \frac{1}{2}\right) - \frac{9}{\theta} \\ \implies \theta_{MAP} &= \frac{9}{3 + \frac{1}{2}} = 2.57142857143 \end{aligned}$$

3 Precision and Recall



How many retrieved items are relevant?

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

How many relevant items are retrieved?

$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

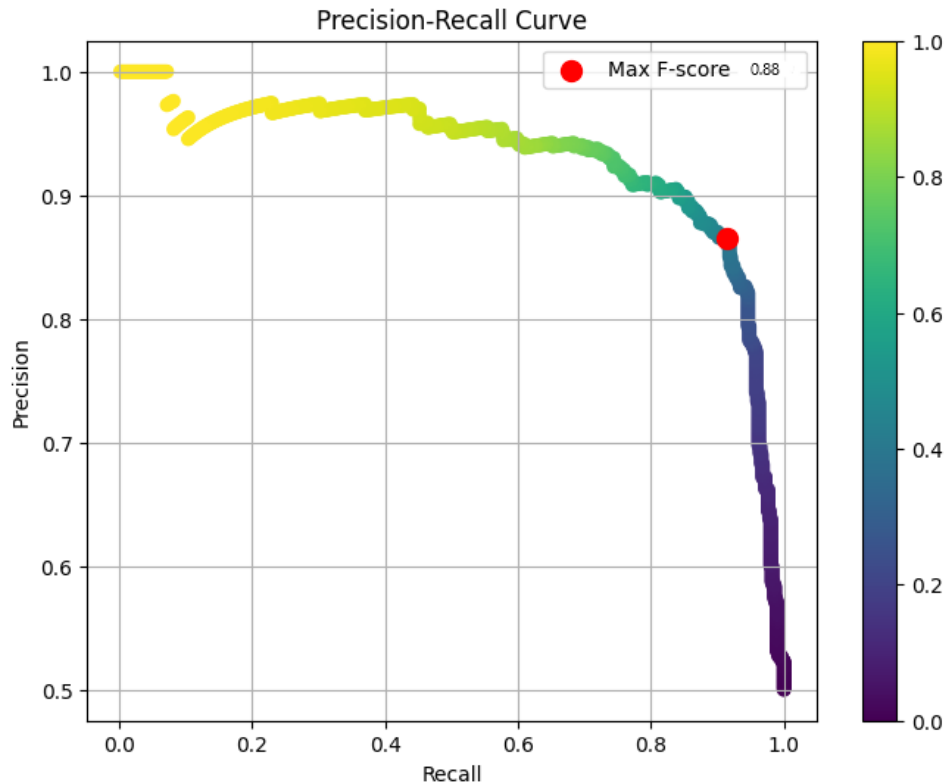
The following chart is known as a *confusion matrix* and helps formalize the concepts displayed above. There are 4 categories in the chart:

- *True positives*: items that are predicted positive and have actual label positive
- *False positives*: items that are predicted positive but have actual label negative
- *True negatives*: items that are predicted negative and have actual label negative
- *False negatives*: items that are predicted negative but have actual label positive

		Actual Values	
		Positive (1)	Negative (0)
Predicted Values	Positive (1)	TP	FP
	Negative (0)	FN	TN

- *Type I error*: occurs when we predict a false positive (erroneously predict a positive label when the true label is negative)
 - *Type II error*: occurs when we predict a false negative (erroneously predict a negative label when the true label is positive)
1. What is the formula for precision in terms of the values in the confusion matrix? What about recall? $\text{Precision} = \text{TP}/(\text{TP} + \text{FP})$, $\text{Recall} = \text{TP}/(\text{TP} + \text{FN})$
 2. The *base rate* is the proportion of items that have true label positive. What is the formula for the base rate in terms of the confusion matrix? $\text{base rate} = (\text{TP} + \text{FN}) / (\text{TP} + \text{FP} + \text{FN} + \text{TN})$
 3. Suppose we predict every item to be positive. What is the precision? What is the recall? $\text{precision} = \text{base rate}$, $\text{recall} = 1$

4. The F_1 score is defined as the harmonic mean of the precision and recall: $F_1 = \frac{2}{1/P+1/R}$. The following image shows an example curve of precision and recall for a classifier when varying the threshold between the positive and negative classes. The point on the curve with highest F_1 score is marked.



Draw an example precision-recall curve for a “better” classifier than the one shown. Mark the point with the optimal F_1 score.

Draw an example precision-recall curve for a “worse” classifier than the one shown. Mark the point with the optimal F_1 score.

