

CM 146 — Machine Learning: Midterm

Fall 2017

Name: Zhouyang Xue

UID: 104629708

Instructions:

1. This exam is CLOSED BOOK and CLOSED NOTES.
2. You may use scratch paper if needed.
3. The time limit for the exam is 1 hour, 45 minutes.
4. Mark your answers ON THE EXAM ITSELF. If you make a mess, clearly indicate your final answer (box it).
5. For true/false questions, CIRCLE True OR False and provide a brief justification for full credit.
6. Unless otherwise instructed, for multiple-choice questions, CIRCLE ALL CORRECT CHOICES (in some cases, there may be more than one) and provide a brief justification if the question asks for one.
7. If you think something about a question is open to interpretation, feel free to ask the instructor or make a note on the exam.

Q	Problem	Points	Score
1	ML basics	6	
2	Application	4	
3	True/False	12	
4	Multiple choice	7	
5	Maximum likelihood	5	
6	Decision Trees	10	
7	Regression	16	
Total		60	

1. (6 pts) Machine Learning Basics

(a) (2 pts) Consider supervised and unsupervised learning. What is the main difference in the inputs and the goals?

- ① Supervised learning takes labels as input. Its goal is to label an object into ^{based on its knowledge} certain categories.
- ② Unsupervised learning doesn't have labels. Its goal is to form clusters of the given data.

(b) (2 pts) What is the main difference between classification and regression?

- ① The result/data of classification is discrete.
- ② While that of regression is real numbers (continuous).

(c) (2 pts) What is the motivation to separate the available data into training and test data?

Training data and test data don't intersect. The separation of training and test data allows a model to learn on a part of the data base, and test on unfamiliar, undisturbed data. This makes the model much more reliable in real world, because normally the data to be predicted is totally new towards the model.

2. (4 pts) **Application** Suppose you are given a dataset of cellular images from patients with and without cancer.

(a) (2 pts) Consider the models that we have discussed in lecture: decision trees, k -NN, logistic regression, perceptrons. If you are required to train a model that predicts the probability that the patient has cancer, which of these would you prefer, and why?

Logistic Regression.

Because only this model has probability interpretation.
The range of the result is from 0 to 1.

(b) (2 pts) A model that attains 100% accuracy on the training set and 70% accuracy on the test set is better than a model that attains 80% accuracy on the training set and 75% accuracy on the test set.

True

False

Memorization allows a model to attain 100% accuracy on training set, yet it doesn't guarantee accuracy on testing set (there might be problems such as overfitting, etc.). Not useful either.

Therefore, accuracy on test set indicates a model's performance much better than training set. The latter has 75% accuracy on training, while the former one only has 70%.
The latter model is better.

True/False

3. (2 pts) You are given a training dataset with attributes A_1, \dots, A_m and instances $x^{(1)}, \dots, x^{(n)}$ and you use the ID3 algorithm to build a decision tree D_1 . You then take one of the instances, add a copy of it to the training set (so your new training set will have $n + 1$ instances), and rerun the decision tree learning algorithm (with the same random seed) to create D_2 . D_1 and D_2 are necessarily identical decision trees.

True

False

Adding a copy to the training set might affect the entropy of some/multiple expansions, thus influencing the information gain. This could make D_2 different from D_1 .

4. (2 pts) Stochastic Gradient Descent is faster per iteration than Batch Gradient Descent.

True

False

Time complexity Stochastic Gradient Descent: $O(D)$
Batch Gradient Descent: $O(ND)$

5. (2 pts) You run the PerceptronTrain algorithm with $maxIter = 100$. The algorithm terminates at the end of 100 iterations with a classifier that attains a training error of 1%. This means that the training data is not linearly separable.

True

False

it might take more than 100 iterations for the algorithm to converge. It is possible that after some more iterations, training error lowers down to 0%.

6. (2 pts) We want to learn a non-linear regression function to predict y from \mathbf{x} where $y \in \mathbb{R}$, $\mathbf{x} \in \mathbb{R}^D$ given training data $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$. To do so, we transform \mathbf{x} by a function $\phi(\mathbf{x})$ and minimize the residual sum of squares objective function on the transformed features: $\sum_{i=1}^n (y_i - \theta^T \phi(\mathbf{x}_i))^2$. This optimization problem is convex.

True

False

$$(y_i - \theta^T \phi(\mathbf{x}_i))^2$$

the outer layer of the composed function is $(y_i - \theta^T f_i)^2$,
 which is convex, the whole function is convex. *where $f_i = \phi(\mathbf{x}_i)$*

7. (2 pts) We want to use 1-Nearest Neighbors (1-NN) to classify houses into one of two classes (cheap vs expensive) given a single feature that measures the area of the house. The predictions made by the 1-NN classifier data can change if the area of the house is measured in square metres instead of square feet. (You can neglect the effect of ties i.e., two training instances that are both nearest neighbors to a test instance.)

True

False

If all houses are measured in m^2 ,
 the ratio remains the same.

The nearest neighbors to the new house remains the same.

Thus predictions remain the SAME.

8. (2 pts) You run gradient descent to minimize the function $f(x) = (2x-3)^2$. Assume the step size has been chosen appropriately and you run gradient descent till convergence. Then gradient descent will return the global minimum of f .

True

False

$$f(x) = 4x^2 - 12x + 9$$

$$f'(x) = 8x - 12$$

$$f''(x) = 8 > 0$$

$f(x)$ is convex, so the local minima is also a
 global minima

Multiple choice

9. (2 pts) In k -nearest neighbor classification, which of the following statements are true? (circle all that are correct)

- (a) The decision boundary is smoother with smaller values of k .
- (b) k -NN does not require any parameters to be learned in the training step (for a fixed value of k and a fixed distance function).
- (c) If we set k equal to the number of instances in the training data, k -NN will predict the same class for any input.
- (d) For larger values of k , it is more likely that the classifier will overfit than underfit.

10. (2 pts) Assume we are given a set of one-dimensional inputs and their corresponding output (that is, a set of $\{(x_i, y_i)\}, x_i \in \mathbb{R}, y_i \in \mathbb{R}$). We would like to compare the following two models on our input dataset where $\theta \in \mathbb{R}$:

$$A : y = \theta^2 x$$

$$B : y = \theta x$$

For each model, we split into training and testing set to evaluate the learned model. Which of the following is correct? Choose the answer that best describes the outcome, and provide justification.

- (a) There are datasets for which A would be more *accurate* than B.
- (b) There are datasets for which B would be more *accurate* than A.
- (c) Both (a) and (b) are correct.
- (d) They would perform equally well on all datasets.

for A, θ^2 is always non-negative ($\theta^2 \geq 0$)

for B, θ could be any real number.

this makes A unable to perfectly fit some data set



$\theta = -1$ for B fits well,
yet no possible value of θ allows
A to fit well

11. (3 pts) If your model is overfitting, increasing the training set size (by drawing more instances from the underlying distribution) will tend to result in which of the following? (circle the best answer for each)

- (a) training error will ... increase / decrease / unknown
- (b) test error will ... increase / decrease / unknown
- (c) overfitting will ... increase / decrease / unknown

For these problems, you must show your work to receive credit!

Maximum likelihood

12. We observe the following data consisting of four independent random variables $X_n, n \in \{1, \dots, 4\}$ drawn from the same Bernoulli distribution with parameter θ (i.e., $P(X_n = 1) = \theta$): $(X_1, X_2, X_3, X_4) = (1, 1, 0, 1)$.

- (a) Give an expression for the log likelihood $l(\theta)$ as a function of θ given this specific dataset. [2 pts]

$$L(\theta) = P(X_1, X_2, X_3, X_4; \theta) \\ = \prod_{i=1}^4 P(X_i; \theta)$$

$$l(\theta) = \log(L(\theta)) = \log \left(\prod_{i=1}^4 P(X_i; \theta) \right) \\ = \sum_{i=1}^4 \log(P(X_i; \theta))$$

$$= \log(1 \cdot \theta) + \log(1 \cdot \theta) + \log(1 \cdot (1 - \theta)) + \log(1 \cdot \theta) \\ = 3 \log \theta + \log(1 - \theta)$$

- (b) Give an expression for the derivative of the log likelihood for this specific dataset. [2 pts]

$$\cancel{l(\theta) = \frac{3}{\theta}}$$

$$l'(\theta) = \frac{d}{d\theta} (3 \log \theta + \log(1 - \theta))$$

$$= \frac{3}{\theta} - \frac{1}{1 - \theta}$$

$$= \frac{3}{\theta} + \frac{1}{\theta - 1}$$

(c) What is the maximum likelihood estimate $\hat{\theta}$ of θ ? [1 pts]

$$\frac{3}{\theta} + \frac{1}{\theta-1} = 0$$

$$-\frac{3}{\theta} = \frac{1}{\theta-1}$$

$$3\theta - 3 = -\theta$$

$$4\theta = 3$$

$$\hat{\theta} = \frac{3}{4}$$

Decision Trees

13. We would like to learn a decision tree given the following pairs of training instances with attributes (a_1, a_2) and target variable Y .

Instance number	a_1	a_2	Y
1	T	T	T
2	T	T	T
3	T	F	F
4	F	F	T
5	F	T	F
6	F	T	F

For reference, for a random variable X that takes on two values with probability p and $1 - p$, here are some values of the entropy function (we use **log to the base 2** in this question):

$$p = \frac{1}{2} : H(X) = 1$$

$$p \in \{\frac{1}{3}, \frac{2}{3}\} : H(X) \approx .92$$

- (a) What is the entropy of Y ? [1 pts]

Y : TTTFFF

$$H(Y) = \sum_{i=1}^2 -P(Y=T) \log_2(P(Y=T))$$

$$= -\frac{1}{2} \cdot (-1) - \frac{1}{2} \cdot (-1)$$

$$= 1$$

$$/ P = \frac{1}{2} \cdot H(Y) = 1$$

(b) What is the information gain of each of the attributes a_1 and a_2 relative to Y ? [4 pts]

expand on a_1 : $\begin{matrix} TTF \\ FTF \end{matrix}$ $F: \frac{1}{3}$ $p \in \{\frac{1}{3}, \frac{2}{3}\} = H(a_1) = 0.92$

information gain $G = 1 - 0.92 = 0.08$

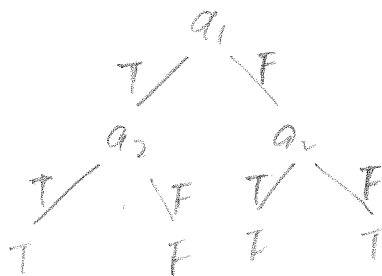
expand on a_2 : $\begin{matrix} TTF \\ TF \end{matrix}$ $F: \frac{1}{2}$ $p \in \{\frac{1}{2}\} = H(a_2) = 1$

information gain $G = 1 - 1 = 0$

(c) Using information gain, which attribute will the ID3 decision tree learning algorithm choose as the root? [1 pts]

will choose a_1

(d) Construct a decision tree with zero training error on this training data. [2 pts]



(e) Change exactly one of the instances (by changing either the attributes or labels but not both) so that **no decision tree can attain zero training error** on this dataset (indicate the instance number and the change). [2 pts]

Change the result (T) of the first instance (1) to F

Weighted linear regression

14. In the problem set, we considered weighted linear regression where the input features are 1-dimensional. We now extend this to D -dimensional features. Thus, we want to find θ that minimizes the cost function

$$J(\theta) = \sum_{n=1}^N w_n (y_n - \theta^T x_n)^2$$

Here $w_n > 0$, $x_n \in \mathbb{R}^{D+1}$, $\theta \in \mathbb{R}^{D+1}$. $X = \begin{pmatrix} x_1^T \\ \vdots \\ x_N^T \end{pmatrix} \in \mathbb{R}^{N \times (D+1)}$, $y \in \mathbb{R}^N$. For this problem, assume that the intercept term is included in the θ and that the linear regression solution exists in this setting.

Questions:

(a) Show that $J(\theta)$ can also be written as:

$$J(\theta) = (y - X\theta)^T W (y - X\theta)$$

Here W is a diagonal matrix where the entry on the diagonal on row n , column n is w_n . [3 pts]

$$X\theta = \begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{pmatrix} \begin{pmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_D \end{pmatrix} = \begin{pmatrix} x_1^T \theta \\ x_2^T \theta \\ \vdots \\ x_N^T \theta \end{pmatrix}$$

$$\therefore y - X\theta = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} - \begin{pmatrix} x_1^T \theta \\ x_2^T \theta \\ \vdots \\ x_N^T \theta \end{pmatrix} = \begin{pmatrix} y_1 - x_1^T \theta \\ y_2 - x_2^T \theta \\ \vdots \\ y_N - x_N^T \theta \end{pmatrix}$$

$$W(y - X\theta) = \begin{bmatrix} w_1 & 0 & \dots \\ 0 & w_2 & \dots \\ \vdots & \vdots & \ddots \\ 0 & 0 & \dots \\ \vdots & \vdots & \dots \\ 0 & 0 & \dots \end{bmatrix} \cdot \begin{pmatrix} y_1 - x_1^T \theta \\ y_2 - x_2^T \theta \\ \vdots \\ y_N - x_N^T \theta \end{pmatrix} = \begin{pmatrix} w_1 (y_1 - x_1^T \theta) \\ w_2 (y_2 - x_2^T \theta) \\ \vdots \\ w_N (y_N - x_N^T \theta) \end{pmatrix}$$

$$(y - X\theta)^T W (y - X\theta) = [y_1 - x_1^T \theta, y_2 - x_2^T \theta, \dots, y_N - x_N^T \theta] \cdot \begin{pmatrix} w_1 (y_1 - x_1^T \theta) \\ w_2 (y_2 - x_2^T \theta) \\ \vdots \\ w_N (y_N - x_N^T \theta) \end{pmatrix}$$

$$= w_1 (y_1 - x_1^T \theta)^2 + w_2 (y_2 - x_2^T \theta)^2 + \dots + w_N (y_N - x_N^T \theta)^2$$

$$= \sum_{n=1}^N w_n (y_n - x_n^T \theta)^2$$

$$= \sum_{n=1}^N w_n (y_n - \theta^T x_n)^2$$

same as $J(\theta)$

- (b) Show that the optimal value for $\hat{\theta} = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{y}$. For reference, here are some useful gradient identities (where \mathbf{x}, \mathbf{b} are vectors and \mathbf{A} is a symmetric matrix).

$$\begin{aligned} f(\mathbf{x}) &= \mathbf{b}^T \mathbf{x} & \nabla f(\mathbf{x}) &= \mathbf{b} \\ f(\mathbf{x}) &= \mathbf{x}^T \mathbf{A} \mathbf{x} & \nabla f(\mathbf{x}) &= 2\mathbf{A} \mathbf{x} \end{aligned}$$

[5 pts]

$$\nabla J(\theta) = \frac{\partial (y - \mathbf{x}\theta)^T \mathbf{W} (y - \mathbf{x}\theta)}{\partial \theta}$$

$$\begin{aligned} & \mathbf{x}^T \cdot 2\mathbf{W} (y - \mathbf{x}\theta) \\ & 2\mathbf{x}^T \mathbf{W} (y - \mathbf{x}\theta) \\ & 2\mathbf{x}^T \mathbf{W} y - 2\mathbf{x}^T \mathbf{W} \mathbf{x} \theta \end{aligned}$$

$$\text{let } f(\theta) = y - \mathbf{x}\theta \quad f'(\theta) = -\mathbf{x}^T$$

$$\nabla J(\theta) = \frac{\partial f(\theta)^T \mathbf{W} f(\theta)}{\partial \theta}$$

$$= f'(\theta) \cdot 2 \cdot \mathbf{W} f(\theta)$$

$$= -\mathbf{x}^T \cdot 2 \cdot \mathbf{W} \cdot (y - \mathbf{x}\theta)$$

$$= -2\mathbf{x}^T \mathbf{W} y + 2\mathbf{x}^T \mathbf{W} \mathbf{x} \theta$$

$$-2\mathbf{x}^T \mathbf{W} y + 2\mathbf{x}^T \mathbf{W} \mathbf{x} \hat{\theta} = 0$$

$$2\mathbf{x}^T \mathbf{W} \mathbf{x} \hat{\theta} = 2\mathbf{x}^T \mathbf{W} y$$

$$\hat{\theta} = \frac{\mathbf{x}^T \mathbf{W} y}{\mathbf{x}^T \mathbf{W} \mathbf{x}}$$

$$= (\mathbf{x}^T \mathbf{W} \mathbf{x})^{-1} \cdot \mathbf{x}^T \mathbf{W} y$$

- (c) In class, we provided a probabilistic interpretation of ordinary least squares (OLS). We now try to provide a probabilistic interpretation of weighted linear regression. Consider a model where each of the N samples is independently drawn according to a normal distribution

$$P(y_n | \mathbf{x}_n, \theta) = \frac{1}{\sqrt{2\pi\sigma_n^2}} \exp\left(-\frac{(y_n - \theta^T \mathbf{x}_n)^2}{2\sigma_n^2}\right)$$

In this model, each y_n is drawn from a normal distribution with mean $\theta^T \mathbf{x}_n$ and variance σ_n^2 . The σ_n^2 are known. Write the log likelihood of this model as a function of θ . [3 points]

$$\begin{aligned} \ell(\theta) &= \sum_{n=1}^N \log(P(y_n | \mathbf{x}_n, \theta)) \\ &= \sum_{n=1}^N \left[\log\left(\frac{1}{\sqrt{2\pi\sigma_n^2}}\right) + \log e^{-\frac{(y_n - \theta^T \mathbf{x}_n)^2}{2\sigma_n^2}} \right] \\ &= \sum_{n=1}^N \left[\log\left(\frac{1}{\sqrt{2\pi\sigma_n^2}}\right) - \frac{(y_n - \theta^T \mathbf{x}_n)^2}{2\sigma_n^2} \right] \\ &= -\frac{1}{2} \sum_{n=1}^N \log(2\pi\sigma_n^2) - \sum_{n=1}^N \frac{(y_n - \theta^T \mathbf{x}_n)^2}{2\sigma_n^2} \end{aligned}$$

- (d) Show that finding the maximum likelihood estimate of θ leads to the same answer as solving a weighted linear regression. How do σ_n^2 relate to w_n ? [5 points]

$$\begin{aligned} \ell'(\theta) &= -\frac{1}{\sigma_n^2} \sum_{n=1}^N \frac{(y_n - \mathbf{x}_n^T \theta)^2}{2\sigma_n^2} \\ &= -\sum_{n=1}^N -(\mathbf{x}_n^T) \cdot \frac{1}{\sigma_n^2} \cdot 2(y_n - \mathbf{x}_n^T \theta) \\ &= \sum_{n=1}^N \mathbf{x}_n^T \cdot \frac{1}{\sigma_n^2} \cdot (y_n - \mathbf{x}_n^T \theta) \\ \sum_{n=1}^N \mathbf{x}_n^T \cdot \frac{1}{\sigma_n^2} (y_n - \mathbf{x}_n^T \theta) &= 0 \\ \sum_{n=1}^N \mathbf{x}_n^T \cdot \frac{1}{\sigma_n^2} \cdot y_n \hat{\theta} &= \sum_{n=1}^N \mathbf{x}_n^T \cdot \frac{1}{\sigma_n^2} \cdot y_n \\ \hat{\theta} &= \sum_{n=1}^N (\mathbf{x}_n^T \cdot \frac{1}{\sigma_n^2} \cdot \mathbf{x}_n)^{-1} \cdot (\mathbf{x}_n^T \cdot \frac{1}{\sigma_n^2} \cdot y_n) \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \cdot \mathbf{X}^T \mathbf{W} \mathbf{y} \end{aligned}$$

$\frac{1}{\sigma_n^2} = w_n$

(Blank page provided for your work)