

1 Week 2

You'd like to use polynomial regression to predict a student's final exam score from their midterm exam score. Concretely, suppose you want to fit a model of the form $h(\theta)$, where x is the midterm score and x^2 is (midterm score)². Further, you plan to use both feature scaling (dividing by the "max-min", or range, of a feature) and mean normalization. What is the normalized feature $x(1)$? (Hint: midterm = 72, final = 74 is training example 2.)

Please round your answer to two decimal places and enter in the text box below.

Suppose you have training examples with features (excluding the additional all-ones feature for the intercept term, which you should add). The normal equation is $(X^T X + \lambda I)^{-1} X^T y$. For the given values of X and y , what are the dimensions of θ , λ , and I in this equation?

1.1 Question 4

Suppose you have a dataset with $m = 1000000$ examples and $n = 200000$ features for each example. You want to use multivariate linear regression to fit the parameters θ to our data. Should you prefer gradient descent or the normal equation?

1.2 Question 1.

Suppose $m=4$ students have taken some class, and the class had a midterm exam and a final exam. You have collected a dataset of their scores on the two exams, which is as follows:

	midterm exam	(midterm exam) ²	final exam
follows:	89	7921	96
	72	5184	74
	94	8836	87
	69	4761	78

You'd like to use polynomial regression to predict a student's final exam score from their midterm exam score. Concretely, suppose you want to fit a model of the form $h(\theta) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$, where x_1 is the midterm score and x_2 is (midterm score)². Further, you plan to use both feature scaling (dividing by the "max-min", or range, of a feature) and mean normalization.

What is the normalized feature $x(1)$? (Hint: midterm = 89, final = 96 is training example 1.) Please round off your answer to two decimal places and enter in the text box below.

Enter answer here

1.3 Question 2.

You run gradient descent for 15 iterations

with $\alpha = 0.3$ and compute

$J(\theta)$ after each iteration. You find that the

value of $J(\theta)$ **decreases** quickly then levels

off. Based on this, which of the following conclusions seems most plausible?

$\alpha = 0.3$ is an effective choice of learning rate. [YES]

Rather than use the current value of α , it'd be more promising to try a smaller value of α (say $\alpha=0.1$).

Rather than use the current value of α , it'd be more promising to try a larger value of α (say $\alpha=1.0$).

1.4 Question 3.

Suppose you have $m=23$ training examples with $n=5$ features (excluding the additional all-ones feature for the intercept term, which you should add).

The normal equation is $\theta=(X^T X)^{-1} X^T y$. For the given values of m and n , what are the dimensions of θ , X , and y in this equation?

- X is 23×6 , y is 23×1 , θ is 6×1 [YES]
- X is 23×5 , y is 23×1 , θ is 5×5
- X is 23×5 , y is 23×1 , θ is 5×1
- X is 23×6 , y is 23×6 , θ is 6×6

1.5 Question 4.

Suppose you have a dataset with $m=1000000$ examples and $n=200000$ features for each example. You want to use multivariate linear regression to fit the parameters θ to our data. Should you prefer gradient descent or the normal equation?

- The normal equation, since it provides an efficient way to directly find the solution. [NO]
- Gradient descent, since it will always converge to the optimal θ .
- Gradient descent, since $(X^T X)^{-1}$ will be very slow to compute in the normal equation.
- The normal equation, since gradient descent might be unable to find the optimal θ .

1.6 Question 5.

Which of the following are reasons for using feature scaling?

- It speeds up solving for θ using the normal equation.
- It prevents the matrix $X^T X$ (used in the normal equation) from being non-invertible (singular/degenerate).
- It speeds up gradient descent by making it require fewer iterations to get to a good solution.
- It is necessary to prevent gradient descent from getting stuck in local optima.

4
0.5
3

NO No matter how θ_0 and θ_1 are initialized, so long as α is sufficiently small, you can safely expect gradient descent to converge to the same solution.

Correct 0.25

This is not true, because depending on the initial condition, gradient descent may end up at a local minimum.

YES If the learning rate is too small, then gradient descent may take a very long time to converge.

Correct 0.25

If the learning rate is small, gradient descent ends up taking an extremely small step on each iteration.

YES If θ_0 and θ_1 are initialized at the global minimum, the one iteration will not change their values.

Correct 0.25

At the global minimum, the derivative (gradient) is zero, so gradient descent will not change the parameters.

NO Setting the learning rate α to be very small is not harmful, and can only speed up convergence.

Correct 0.25

If the learning rate is small, gradient descent ends up taking an extremely small step on each iteration.

If θ_0 and θ_1 are initialized at a local minimum, the one iteration will not change their values. Incorrect 0.00

At a local minimum, the derivative (gradient) is zero, so gradient descent will not change the parameters.

YES If the first few iterations of gradient descent cause $f(\theta_0, \theta_1)$ to increase, then the algorithm has converged to a local minimum.

Incorrect 0.00 If α were small enough, then gradient descent should always successfully find the global minimum.

YES If θ_0 and θ_1 are initialized at the global minimum, the one iteration will not change their values.

Correct 0.25 At the global minimum, the derivative (gradient) is zero, so gradient descent will not change the parameters.

YES No matter how θ_0 and θ_1 are initialized, so long as α is sufficiently small, you can safely expect gradient descent to converge to the same solution.

Correct 0.25 This is not true, because depending on the initial condition, gradient descent may end up at a local minimum.

NO Even if the learning rate α is very large, every iteration of gradient descent will move the parameters in the direction of the negative gradient.

Incorrect 0.00 If the learning rate α is too large, one step of gradient descent can overshoot the minimum.

YES If θ_0 and θ_1 are initialized so that $\theta_0 = \theta_1$, then by symmetry (because the cost function is symmetric), the updates to θ_0 and θ_1 will be the same.

Incorrect 0.00 The updates to θ_0 and θ_1 are different (even though we're doing simultaneous updates).

Suppose that for some linear regression problem (say, predicting housing prices as in the lecture), you have managed to find some θ_0, θ_1 such that $J(\theta_0, \theta_1) = 0$. Which of the following statements is true?

Your Answer Score Explanation

NO We can perfectly predict the value of y even for new examples that we have not yet seen.
Incorrect 0.00 Even though we can fit our training set perfectly, this does not mean that

NO This is not possible: By the definition of $J(\theta_0, \theta_1)$, it is not possible for $J(\theta_0, \theta_1) = 0$.
Correct 0.25 If all of our training examples lie perfectly on a line, then $J(\theta_0, \theta_1) = 0$.

YES Our training set can be fit perfectly by a straight line, i.e., all of our training examples lie on a line.
Incorrect 0.00

If $J(\theta_0, \theta_1) = 0$, that means the line defined by the equation " $y = \theta_0 + \theta_1 x$ " perfectly fits

NO Gradient descent is likely to get stuck at a local minimum and fail to find the global minimum.
Incorrect 0.00 The cost function $J(\theta_0, \theta_1)$ for linear regression has no local optima.

NO For this to be true, we must have $y(i) = 0$ for every value of $i = 1, 2, \dots, m$.
Correct 0.25 So long as all of our training examples lie on a straight line, we will be able to find a line that fits them perfectly.
NO For this to be true, we must have $\theta_0 = 0$ and $\theta_1 = 0$ so that $h_{\theta}(x) = 0$ for all x .
Correct 0.25

If $J(\theta_0, \theta_1) = 0$, that means the line defined by the equation " $y = \theta_0 + \theta_1 x$ " perfectly fits all of our data.
There's no particular reason to expect that the values of θ_0 and θ_1 that achieve this are both zero.

NO This is not possible: By the definition of $J(\theta_0, \theta_1)$, it is not possible for $J(\theta_0, \theta_1) = 0$.
Correct 0.25 If all of our training examples lie perfectly on a line, then $J(\theta_0, \theta_1) = 0$.

YES Gradient descent is likely to get stuck at a local minimum and fail to find the global minimum.
Incorrect 0.00 The cost function $J(\theta_0, \theta_1)$ for linear regression has no local optima, so gradient descent will not get stuck at a bad local minimum.

NO Our training set can be fit perfectly by a straight line, i.e., all of our training examples lie on a line.
Incorrect 0.00
If $J(\theta_0, \theta_1) = 0$, that means the line defined by the equation " $y = \theta_0 + \theta_1 x$ " perfectly fits all of our data.

YES For these values of θ_0 and θ_1 that satisfy $J(\theta_0, \theta_1) = 0$, we have that $h_{\theta}(x) = y$ for all x .
Incorrect 0.00 $J(\theta_0, \theta_1) = 0$, that means the line defined by the equation " $y = \theta_0 + \theta_1 x$ " perfectly fits all of our data.

2 Week 3 - Logistic Regression

2.1 Logistic Function

- A logistic function (or logistic curve) is a common sigmoid function, given its name (in reference to its S-shape) in 1844 or 1845 by Pierre François Verhulst who studied it in relation to population growth.
- A generalized logistic curve can model the "S-shaped" behaviour (abbreviated S-curve) of growth of some population P.
- The initial stage of growth is approximately exponential; then, as saturation begins, the growth slows, and at maturity, growth stops. The logistic function is the sigmoid curve with equation:

$$f(x) = \frac{1}{1 + e^{-x}}$$

2.2 Question 1.

Suppose that you have trained a logistic regression classifier, and it outputs on a new example x a prediction $h\theta(x) = 0.4$. This means (check all that apply):

2.3 Question 2.

Suppose you have the following training set, and fit a logistic regression classifier $h\theta(x)=g(\theta_0+\theta_1x_1+\theta_2x_2)$. Which of the following are true? Check all that apply.

2.4 Question 3.

For logistic regression, the gradient is given by EQUATION

Which of these is a correct gradient descent update for logistic regression with a learning rate of α ? Check all that apply.

2.5 Question 4.

Which of the following statements are true? Check all that apply.

Incorrect

2.6 Question 5.

Suppose you train a logistic classifier $h\theta(x)=g(\theta_0+\theta_1x_1+\theta_2x_2)$. Suppose $\theta_0=-6, \theta_1=0, \theta_2=1$.

Which of the following figures represents the decision boundary found by your classifier?

2.7 Question 1

Suppose that you have trained a logistic regression classifier, and it outputs on a new example a prediction $h\theta(x) = 0.7$. This means (check all that apply):

- Our estimate for $\Pr(y = 1|x; \theta)$ is 0.7.
- Our estimate for $\Pr(y = 0|x; \theta)$ is 0.3.
- Our estimate for $\Pr(y = 1|x; \theta)$ is 0.3.
- Our estimate for $\Pr(y = 0|x; \theta)$ is 0.7.

Solution

Our estimate for $\Pr(y = 1|x; \theta)$ is 0.7. T $h\theta(x)$ is precisely $\Pr(y = 1|x; \theta)$, so each is 0.7. Our estimate for $\Pr(y = 0|x; \theta)$ is 0.3. T Since we must have $\Pr(y = 0|x; \theta) = 1 - \Pr(y = 1|x; \theta)$, the former is $1 - 0.7 = 0.3$. Our estimate for $\Pr(y = 1|x; \theta)$ is 0.3. F $h\theta(x)$ gives $\Pr(y = 1|x; \theta)$, not $1 - \Pr(y = 1|x; \theta)$. Our estimate for $\Pr(y = 0|x; \theta)$ is 0.7. F $h\theta(x)$ is $\Pr(y = 1|x; \theta)$, not $\Pr(y = 0|x; \theta)$.

2.8 Question 3

Suppose you have the following training set, and fit a logistic regression classifier $h\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$.

Which of the following are true? Check all that apply.

Adding polynomial features (e.g., instead using $h(x) = g(0 + 1x_1 + 2x_2 + 3x_1^2 + 4x_1x_2 + 5x_2^2)$) could increase how well we can fit the training data. The positive and negative examples cannot be separated using a straight line. So, gradient descent will fail to converge.

At the optimal value of θ (e.g., found by `fminunc`), we will have $J(\theta) \geq 0$. Because the positive and negative examples cannot be separated using a straight line, linear regression will perform as well as logistic regression on this data. solution

1) Adding polynomial features (e.g., instead using $h(x) = g(0 + 1x_1 + 2x_2 + 3x_1^2 + 4x_1x_2 + 5x_2^2)$) could increase how well we can fit the training data. TRUE Adding new features can only improve the fit on the training set: since setting $\theta_3 = \theta_4 = \theta_5 = 0$ makes the hypothesis the same as the original one, gradient descent will use those features (by making the corresponding θ_j non-zero) only if doing so improves the training set fit.

2) The positive and negative examples cannot be separated using a straight line. So, gradient descent will fail to converge. FALSE While it is true they cannot be separated, gradient descent will still converge to the optimal fit. Some examples will remain misclassified at the optimum.

3) At the optimal value of θ (e.g., found by `fminunc`), we will have $J(\theta) \geq 0$. TRUE The cost function $J(\theta)$ is always non-negative for logistic regression.

4) Because the positive and negative examples cannot be separated using a straight line, linear regression will perform as well as logistic regression on this data. FALSE While it is true they cannot be separated, logistic regression will outperform linear regression since its cost function focuses on classification, not prediction.

2.9 Question 5

Which of the following statements are true? Check all that apply.

- For logistic regression, sometimes gradient descent will converge to a local minimum (and fail to find the global minimum). This is the reason we prefer more advanced optimization algorithms such as fminunc (conjugate gradient/BFGS/L-BFGS/etc).
- Since we train one classifier when there are two classes, we train two classifiers when there are three classes (and we do one-vs-all classification).
- The sigmoid function $g(z)$ is never greater than one ($g(z) \leq 1$). $g(z) = \frac{1}{1+e^{-z}}$
- The one-vs-all technique allows you to use logistic regression for problems in which each $y(i)$ comes from a fixed, discrete set of values.

Solutions 1) For logistic regression, sometimes gradient descent will converge to a local minimum (and fail to find the global minimum). This is the reason we prefer more advanced optimization algorithms such as fminunc (conjugate gradient/BFGS/L-BFGS/etc). 0.00

The cost function for logistic regression is convex, so gradient descent will always converge to the global minimum. We still might use a more advanced optimization algorithm since they can be faster and don't require you to select a learning rate.

2) Since we train one classifier when there are two classes, we train two classifiers when there are three classes (and we do one-vs-all classification). We need to train three classifiers if there are three classes; each one treats one of the three classes as the $y=1$ examples and the rest as the $y=0$ examples. 3) The sigmoid function $g(z)$ is never greater than one ($g(z) \leq 1$). $g(z) = \frac{1}{1+e^{-z}}$

The denominator ranges from 1 to e^z as z grows, so the result is always in $(0,1)$. 4) The one-vs-all technique allows you to use logistic regression for problems in which each $y(i)$ comes from a fixed, discrete set of values. If each $y(i)$ is one of k different values, we can give a label to each $y(i) \in \{1, 2, \dots, k\}$ and use one-vs-all as described in the lecture.

The mean of x_2 is 6675.5 and the range is $8836-4761=4075$ So $x(1)/4075 = 0.47$.

Question 2

You run gradient descent for 15 iterations with $\alpha=0.3$ and compute $J(\theta)$ after each iteration. You find that the value of $J(\theta)$ decreases quickly then levels off. Based on this, which of the following conclusions seems most plausible? Your Answer Score Explanation

- Rather than use the current value of α , it'd be more promising to try a larger value of α (say $\alpha=1.0$).
Incorrect 0.00 A larger value for α will make it more likely that $J(\theta)$ diverges.

- Rather than use the current value of α , it'd be more promising to try a smaller value of α (say $\alpha=0.1$).
- $\alpha=0.3$ is an effective choice of learning rate.

Total 0.00 / 1.00

Rather than use the current value of α , it'd be more promising to try a smaller value of α (say $\alpha=0.1$). Correct 1.00 Since the cost function is increasing, we know that gradient descent is diverging, so we need a lower learning rate.

Question 3

Suppose you have $m=28$ training examples with $n=4$ features (excluding the additional all-ones feature for the intercept term, which you should add). The normal equation is $\theta=(X^T X)^{-1} X^T y$. For the given values of m and n , what are the dimensions of θ , X , and y in this equation?

- X is 28×4 , y is 28×1 , θ is 4×4
- X is 28×4 , y is 28×1 , θ is 4×1
- X is 28×5 , y is 28×1 , θ is 5×1 **Correct 1.0**
- X is 28×5 , y is 28×5 , θ is 5×5

Total 1.00 / 1.00 Question Explanation

X has m rows and $n+1$ columns (+1 because of the $x_0=1$ term). y is an m -vector. θ is an $(n+1)$ -vector.

Question 4

Suppose you have a dataset with $m=50$ examples and $n=200000$ features for each example. You want to use multivariate linear regression to fit the parameters θ to our data. Should you prefer gradient descent or the normal equation?

- Gradient descent, since $(X^T X)^{-1}$ will be very slow to compute in the normal equation.
Correct
- The normal equation, since it provides an efficient way to directly find the solution.
- The normal equation, since gradient descent might be unable to find the optimal θ .

Incorrect 0.00

For an appropriate choice of α , gradient descent can always find the optimal θ .

- Gradient descent, since it will always converge to the optimal θ .

Question 5

Which of the following are reasons for using feature scaling?

- FALSE It prevents the matrix XTX (used in the normal equation) from being non-invertable (singular/degenerate).
Correct 0.25 XTX can be singular when features are redundant or there are too few examples. Feature scaling does not solve these problems.
- FALSE It speeds up gradient descent by making it require fewer iterations to get to a good solution.
Correct 0.25 Feature scaling speeds up gradient descent by avoiding many extra iterations that are required when one or more features take on much larger values than the rest.
- TRUE It speeds up solving for θ using the normal equation.
Incorrect 0.00 The magnitude of the feature values are insignificant in terms of computational cost.
- FALSE It is necessary to prevent gradient descent from getting stuck in local optima.
Correct 0.25 The cost function $J(\theta)$ for linear regression has no local optima.