

week 3: ROC curves, Normalization and PCA

<http://mlvu.github.io>

February 18, 2020

1 Exam questions

At this point in the course, it's a good idea to have a look at the practice exam. The homework prepares you for the exam, but the format of the homework questions is different from that of the exam questions. Here are some questions you should be able to answer based on the material covered so far. You can find two complete practice exams on Canvas.

1. What is the advantage of gradient descent over random search?
 - A In gradient descent, parallel searches are allowed to communicate.
 - B Gradient descent is less likely to get stuck in local minima.
 - C Gradient descent computes the direction of steepest descent, random search approximates it.
 - D Gradient descent is easier to parallelise.
2. Why is accuracy a bad loss function to use in gradient descent?
 - A It is expensive to compute.
 - B It makes the gradient zero almost everywhere.
 - C It is unreliable in situations with high class imbalance.
 - D The confidence interval is high on small test sets.
3. I'm performing spam classification. I represent each email by three numbers: how often the word *pill* occurs, how often the word *hello* occurs and how often the word *congratulations* occurs. What are these three attributes called?
 - A The instances
 - B The classes
 - C The features
 - D The principal components

Here we see the derivation of the gradient of the squared-error loss for linear regression.

$$\frac{\partial \frac{1}{2} \sum_i (f(x_i) - y_i)^2}{\partial b} = \frac{1}{2} \frac{\partial \sum_i (x_i w + b - y_i)^2}{\partial b} \quad (1)$$

$$= \frac{1}{2} \sum_i \frac{\partial (x_i w + b - y_i)^2}{\partial b} \quad (2)$$

$$= \frac{1}{2} \sum_i \frac{\partial (x_i w + b - y_i)^2}{\partial (x_i w + b - y_i)} \frac{\partial (x_i w + b - y_i)}{\partial b} \quad (3)$$

$$= \sum_i (x_i w + b - y_i) \frac{\partial (x_i w + b - y_i)}{\partial b} \quad (4)$$

$$= \dots \quad (5)$$

4. To get from line (1) to line (2), we use the
A Chain rule **B** Product rule **C** Exponent rule **D** Sum rule
5. To get from line (2) to line (3), we use the
A Chain rule **B** Product rule **C** Exponent rule **D** Sum rule
6. In line 5, the correct result is
A $\sum_i (x_i^2 w + b - y_i)$
B $\sum_i x_i (x_i w + b - y_i)$
C $\sum_i (x_i w + b - y_i)$
D $\sum_i (x_i w + b - y_i)^2$

2 ROC Curves and confusion matrices

In this section, we'll look at evaluation metrics. This section covers the lecture "Methodology 1" and section 2.2 of Peter Flach's *Machine Learning* (see the PDF on Canvas).

As a running example, we will use the following classification dataset:

	x_1	x_2	label
a	1	0	Neg
b	2	0	Pos
c	3	0	Pos
d	1/2	1	Pos
e	2	2	Pos
f	0	3	Neg
g	1	3	Neg
h	2	3	Neg

If you have trouble with the exercises below, we recommend drawing the data in its feature space, and drawing the decision boundary of each classifier.

2.1 A linear classifier

question 1: As a first classifier, c_{lin} , we will use a diagonal line crossing all the points where $x_1 = x_2$. Points above this line will be classified as **Neg**, points below or on the line as **Pos**. How do we describe this classifier mathematically? Fill in the blanks:

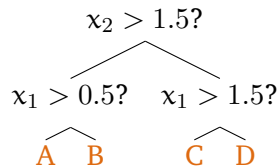
$$c_{\text{lin}}(x_1, x_2) = \begin{cases} \text{Pos} & \text{if } \dots \\ \text{Neg} & \text{otherwise} \end{cases} .$$

c_{lin} is a simple classifier. Let's try a more complex one before we move on.

question 2: Let c_2 be the classifier whose decision boundary crosses the points $(0, 1)$ and $(1, 1.5)$, with the class **Neg** above the line. Draw this line first. How do we define this classifier?

$$c_2(x_1, x_2) = \begin{cases} \text{Pos} & \text{if } \dots \\ \text{Neg} & \text{otherwise} \end{cases}$$

question 3: We will compare the linear classifier c_{lin} with a *decision tree* classifier, c_{tree} . Here is the decision tree:



If the inequality on the node is true, move right, otherwise move left.

This tree divides the feature space into four **segments** (A, B, C and D). To each segment we assign the majority class in that segment (using Pos for a draw). Label the leaves A, B, C and D with classes.

question 4: In this exercise we don't do any training (we just evaluate the given classifiers), so we'll evaluate all metrics on a single given dataset.

If you do train, you would *split* your dataset. You would compute your evaluation metrics on the ... during hyperparameter selection and on the ... during final testing.

Evaluating on the training data is the worst sin you can commit in Machine Learning. Can you think of a situation when you would still want to *compute* the training loss?

question 5: Give the *confusion matrices* for classifiers c_{lin} and c_{tree} .

question 6: From the confusion matrices, we can compute several metrics. Compute the accuracy, precision, recall, true positive rate and false positive rate

To plot a ROC and precision/recall curves, and to compute AUC metrics, we need to turn our simple binary classifier into a *ranking* classifier. That is, we don't just want classifications, we want it to rank the data from most Neg to most Pos. How we do this depends on the classifier.

question 7: How do we turn a linear classifier into a ranking classifier? Give the ranking that c_{lin} assigns to the training data.

question 8: How do we turn a decision tree classifier into a ranking classifier? Give the ranking that c_{tree} assigns to the training data.

We can't tell whether our ranking is perfect, because the training data doesn't give us a correct *ranking*, only a correct *labeling*. But for some pairs of instances we *do* know how they should be ranked: all pairs of positive and negative examples.

We can visualize this in a matrix with **negative** instances on the horizontal axis, and **positive** instances on the vertical, both sorted with the most positive (according to the classifier's ranking) in the bottom left corner. The

matrix is usually colored **green** for pairs that are ranked correctly, **red** for pairs ranked incorrectly, and **yellow** for pairs that are ranked equal.

question 9:

1. Draw the coverage matrices of c_{lin} and c_{tree} .
2. How many ranking errors does each classifier make on the training data?
3. If we have class imbalance (one class has more instances than the other), how can we tell by the coverage matrix?

question 10: How do we convert a coverage matrix to an **ROC plot**? How does the green area in a coverage matrix relate to the ROC plot? How do they differ?

3 Whitening

Basis systems The *standard basis system* (in 2D) is just the vectors $b_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$ and $b_2 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$. When a point has coordinates (x, y) , that just means that we can get to the point by adding x times the b_1 vector to y times the b_2 vector.

For the standard basis, we just end up at the point $\begin{pmatrix} x \\ y \end{pmatrix}$. But we can also use a nonstandard basis. If a point p has coordinates (x, y) in the coordinate system with basis vectors b_1, b_2 , it means that p has coordinates $xb_1 + yb_2$ in the standard basis system.

question 11: Convert the following points to a representation in the standard basis system.

1. The point $\begin{pmatrix} 1 \\ 1 \end{pmatrix}$ expressed in the coordinate system with basis vectors $b_1 = \begin{pmatrix} 2 \\ 0 \end{pmatrix}, b_2 = \begin{pmatrix} 0 \\ 2 \end{pmatrix}$.

2. The point $\begin{pmatrix} 5 \\ 3 \end{pmatrix}$ expressed in the coordinate system with basis vectors $\mathbf{b}_1 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$, $\mathbf{b}_2 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$.
3. The point $\begin{pmatrix} 1 \\ 2 \end{pmatrix}$ expressed in the coordinate system with basis vectors $\mathbf{b}_1 = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$, $\mathbf{b}_2 = \begin{pmatrix} 3 \\ 4 \end{pmatrix}$.

question 12: A basis system is **orthonormal** if two conditions hold:

1. ...
2. ...

Consider the following three bases:

$$\mathbf{A} = \begin{pmatrix} 0 & 1 \\ 1 & 1 \end{pmatrix} \quad \mathbf{B} = \begin{pmatrix} 1 & -1 \\ 1 & 1 \end{pmatrix} \quad \mathbf{C} = \sqrt{1/2} \begin{pmatrix} 1 & -1 \\ 1 & 1 \end{pmatrix}$$

The columns of each matrix are the basis vectors (it may help to draw them in standard coordinates). Which of these have orthogonal basis vectors? Which one is orthonormal?

question 13: Why is it useful to have an orthonormal basis?

Covariance and normal distributions PCA is closely related to the normal distribution. We can think of PCA as a basis transformation that makes it look like the data was sampled from a standard multivariate normal distribution: it makes the variance 1 in every dimension, and the covariance 0 for every pair of dimensions.

question 14: Here is some data.

x_1	x_2
0	2
2	0
0	2
2	0

Compute the sample covariance.

Let $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{t}$ be an *affine* transformation (a linear transformation \mathbf{A} followed by a translation \mathbf{t}). We draw \mathbf{x} from a standard normal distribution

$$\mathbf{N}(\mathbf{0}, \mathbf{I}) = \mathbf{N}\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}\right)$$

and compute \mathbf{y} . This is equivalent to sampling from a normal distribution with mean $\boldsymbol{\mu} = \mathbf{t}$ and covariance $\boldsymbol{\Sigma} = \mathbf{A}\mathbf{A}^T$.

question 15: Let

$$\mathbf{A} = \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix}, \mathbf{t} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

Sampling \mathbf{x} from $\mathbf{N}(\mathbf{0}, \mathbf{I})$ is easy: we just sample x_1 and x_2 independently from a *univariate* standard normal distributions. Here are some **standard normally distributed numbers**:

$$0.5, 1.1, 0.1, 0.9, -0.2, 1.3$$

Turn these into three samples from $\mathbf{N}(\mathbf{0}, \mathbf{I})$ and transform them to samples from $\mathbf{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$.

Bonus question: Compute the sample covariance of the resulting data, and the actual covariance $\boldsymbol{\Sigma}$.