

Week 5: Support Vector Machines and Tensor Backpropagation

<https://mlvu.github.io>

March 7, 2021

1 Support Vector Machines

1.1 Tensor loss

The basic optimization objective for Support Vector Machines is

$$\begin{aligned} &\text{minimize } \frac{1}{2}\|\mathbf{w}\|^2 + C \sum_i p_i \\ &\text{such that } y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - p_i \\ &\text{and } p_i \geq 0 \end{aligned}$$

question 1: What does the variable i index? How many terms does the sum have, and how many constraints are there?

What is the value of y_i in this expression? What is its function?

question 2: There are two common ways to rewrite this expression before implementing it. What are they (in general terms) and what are their benefits?

1.2 Lagrange multipliers

Lagrange multipliers are a useful trick to know. In the lectures we only had time to describe the trick itself, and the rules for applying it. That should be enough to do the exam questions and the questions below, but if you

want some more intuition for why Lagrange multipliers work the way they do please read [this article from the required reading](#).

We have the following optimization problem:

$$\begin{aligned} \text{minimize } f(\mathbf{a}, \mathbf{b}) &= \mathbf{a}^2 + 2\mathbf{b}^2 \\ \text{such that } \mathbf{a}^2 &= -\mathbf{b}^2 + 1 \end{aligned}$$

question 3: The first step is to rewrite the constraint so that the right side is equal to zero. Do so. What does the constraint say about the allowed inputs (what shape do the allowed inputs make in the (\mathbf{a}, \mathbf{b}) -plane)?

We now define a function $L(\mathbf{a}, \mathbf{b}, \alpha) = f(\mathbf{a}, \mathbf{b}) + \alpha G$, where G is the left hand side of the constraint equal to zero (how much any given \mathbf{a} and \mathbf{b} violate the constraint).¹

question 4: Write out $L(\mathbf{a}, \mathbf{b}, \alpha)$ for our problem.

We take the derivative of L with respect to each of its *three* parameters, and set these equal to zero.

question 5: Fill in the blanks

$$\mathbf{a}(\dots) = 0 \tag{1}$$

$$\mathbf{b}(\dots) = 0 \tag{2}$$

$$\mathbf{a}^2 + \mathbf{b}^2 = 1 \tag{3}$$

Note that the last line recovers the original constraint. We now have three equations with three unknowns, so we can solve for \mathbf{a} and \mathbf{b} . From the shape of the function (it's symmetric in both the \mathbf{a} and \mathbf{b} axes), we should expect at least two solutions.

We can get these from the above equations by noting that if \mathbf{a} and \mathbf{b} are both nonzero, we can derive a contradiction. Thus either \mathbf{a} or \mathbf{b} must be zero.

question 6: Give the solutions for both cases (remember that $x^2 = 1$ has *two* solutions).

Happily, [Wolfram Alpha](#) agrees with us (and provides some informative plots).

¹For plain Lagrange multipliers, where the constraints are all equalities, we can either add or subtract the term containing the constraint. For inequality constraints, it depends on whether we are maximizing or minimizing.

1.3 The kernel trick

The feature space of k is a projection of point \mathbf{a} to point \mathbf{a}' such that

$$k(\mathbf{a}, \mathbf{b}) = \mathbf{a}'^T \mathbf{b}' .$$

We have a dataset with two features Let $\mathbf{a} = \begin{pmatrix} a_1 \\ a_2 \end{pmatrix}$ and $\mathbf{b} = \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}$. We define the *kernel*

$$k_1(\mathbf{a}, \mathbf{b}) = (\mathbf{a}^T \mathbf{b})^2 .$$

question 7: Show that the feature space defined by k_1 for a vector $\begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$ is

$$\begin{pmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{pmatrix} .$$

Hint: start by writing out the definition as a scalar function. See if you can re-arrange this back into a dot product of two other vectors.

question 8: What is the feature space for the kernel

$$k_2(\mathbf{a}, \mathbf{b}) = (\mathbf{a}^T \mathbf{b} + 1)^2 \quad ?$$

2 Backpropagation Revisited

2.1 The multivariate chain rule

If we want to do backpropagation for a computation graph where an output variable depend in an input variable through multiple intermediate values (the graph contains a diamond), we require the *multivariate chain rule*. If you don't quite remember how it goes, re-read slides 24–29 fo the deep learning lecture before trying these questions.

We will use the backpropagation algorithm to find the derivative of the function

$$f(x) = \sin(x^2) \cos(x^3)$$

with respect to x .

question 9: First, we break the function up into modules. Fill in the blanks.

$$f = \mathbf{a} \mathbf{b}$$

$$\mathbf{a} = \dots$$

$$\mathbf{b} = \dots$$

$$\mathbf{c} = x^2$$

$$\mathbf{d} = x^3$$

question 10: Draw the computation graph with nodes x , \mathbf{c} , \mathbf{d} , \mathbf{a} , \mathbf{b} , f

question 11: Work out the *local* derivatives. Which is the correct expression for the gradient $\partial f / \partial x$?

This is a common exam question so make sure to practice if you're not sure. You can easily create new questions for yourself by coming up with any function $f(x)$ with a diamond shape in the computation (just make sure that x is used twice).

2.2 Tensor Backpropagation

In scalar backpropagation, we need to work out only the local derivatives. Once we have these, we can multiply them in any order to get the global derivative. If we want to apply backpropagation to tensors, things are not so easy.

question 12: Why not? What is it about the local derivatives in tensor backpropagation that makes it difficult to apply in this way?

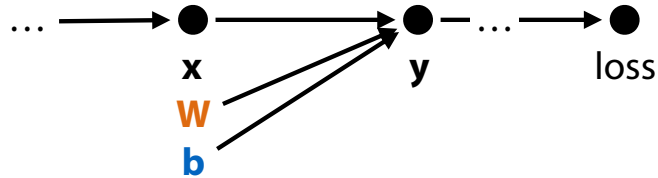
The solution is not to compute the local derivatives explicitly, but only to compute the gradients on the inputs over the loss, given the gradients of the outputs over the loss.

We will practice this for a simple feedforward layer (without activation). Consider a module f that computes:²

$$\mathbf{y} = f(\mathbf{x}, \mathbf{W}, \mathbf{b}) = \mathbf{W}\mathbf{x} + \mathbf{b}$$

²Normally, we would consider \mathbf{x} the input, and \mathbf{W} and \mathbf{b} the parameters, but for our AD engine, the distinction does not matter: everything going in to the computation is an input, and we may need the gradient over all three.

Assume that this module is part of a much larger network. Its output \mathbf{y} is fed as input to another module, which produces a new output that is fed to another module and so on. At the end, a single scalar loss L is produced.



Ultimately, we want to work out the derivative of this loss, with respect to our inputs:

$$\mathbf{x}^\nabla, \mathbf{W}^\nabla, \mathbf{b}^\nabla$$

Where the notation \mathbf{W}^∇ represents a matrix for the scalar derivatives of the loss with respect to elements of \mathbf{W} . For instance, element (2, 3) of matrix \mathbf{W}^∇ is the scalar derivative $\partial l / \partial W_{23}$.

We'll start with the bias term \mathbf{b} . Since matrix/vector calculus can get complicated, and doesn't translate to tensors of higher rank, we'll develop everything in terms of scalar calculus. Once we've taken the derivatives we'll transform everything to matrix operations to make the implementation efficient.

We're interested in $\frac{\partial L}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{b}_j}$, but we need matrix/vector calculus to do that. Instead we will consider $f(\mathbf{b})$ as a *scalar* computation, which has multiple intermediate values $y_1, y_2 \dots, y_k$. By the multivariate chain rule, we can just take the derivative over each path (through each value y_i), and sum them:

$$\frac{\partial l}{\partial \mathbf{b}_j} = \sum_i \frac{\partial l}{\partial y_i} \frac{\partial y_i}{\partial \mathbf{b}_j} = \sum_i y_i^\nabla \frac{\partial y_i}{\partial \mathbf{b}_j}$$

We will assume that \mathbf{y}^∇ is given to us by the automatic differentiation (AD) engine as a vector. All we need to work out is a simple matrix operation that computes $\partial l / \partial \mathbf{b}_j$ for all j and lays the results out in the same shape as \mathbf{b} . We'll start by working out the scalar derivative.

question 13: Fill in the gaps

$$\begin{aligned}\frac{\partial l}{\partial \mathbf{b}_j} &= \sum_i \frac{\partial l}{\partial y_i} \frac{\partial y_i}{\partial \mathbf{b}_j} = \sum_i y_i^\nabla \dots \\ &= \sum_i y_i^\nabla \frac{\partial [\dots]_i}{\partial \mathbf{b}_j} = \sum_i y_i^\nabla \frac{\partial [W\mathbf{x}]_i + \mathbf{b}_i}{\partial \mathbf{b}_j} \\ &= \sum_i y_i^\nabla \frac{\partial \mathbf{b}_i}{\partial \mathbf{b}_j} = \dots\end{aligned}$$

question 14: The backward() function for f if given \mathbf{y}^∇ . What should it return as the gradient for \mathbf{b} ?

We'll do the same for the gradient over \mathbf{x} .

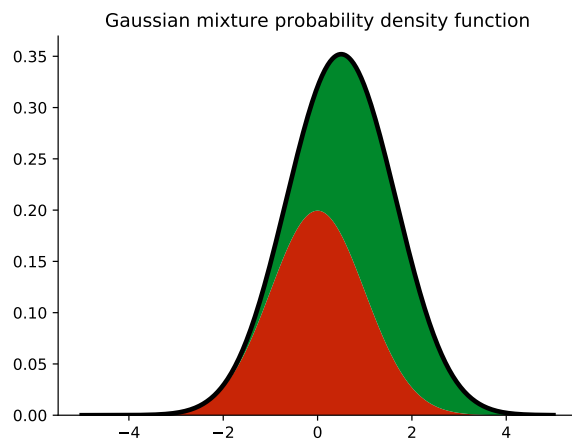
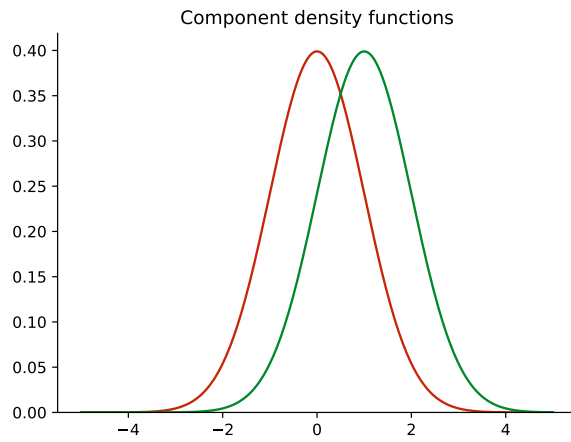
question 15: Work out the scalar derivative $\partial l / \partial x_j$ using the multivariate chain rule (again taking $y_i^\nabla = \partial l / \partial y_i$ as given).

question 16: What should the backward() function for f return as the gradient for \mathbf{x} ?

3 Bonus: Expectation Maximization

This is not an exam question, but it's helpful to do if you want to understand the EM algorithm.

Assume we have a Gaussian Mixture Model in one dimension with two components: $\mathcal{N}(0, 1)$ and $\mathcal{N}(1, 1)$. The weights w_1 and w_2 of the components are equal.



question 17: Compute the probability density of the point 0, under the Gaussian Mixture.

question 18: Under the EM algorithm, what responsibility is assigned to each component for the point 0?