

Introduction à l'apprentissage automatique - TD2

Caio Corro

1 Gradient computation (part 2)

In this exercise, we focus on multiclass classification with k classes. Losses for these problems are defined as functions of the form $\ell : E(k) \times \mathbb{R}^k \rightarrow \mathbb{R}_+$.

1. The negative log-likelihood for multiclass classification is defined as $\ell_{nll}(\mathbf{y}, \mathbf{w}) = -\mathbf{y}^\top \mathbf{w} + \log \sum_j \exp(w_j)$. In practice, we have $\mathbf{w} = \mathbf{A}\mathbf{x} + \mathbf{b}$. Compute gradient of the loss function wrt to \mathbf{A} and \mathbf{b} (note: as \mathbf{A} is a matrix, gradient is an abuse of terminology).
2. The hinge loss for multiclass classification is defined as:

$$\ell_{hinge}(\mathbf{y}, \mathbf{w}) = \max(0, -\mathbf{w}^\top \mathbf{y} + m + \max_{\mathbf{y}' \in E(k) \setminus \{\mathbf{y}\}} \mathbf{w}^\top \mathbf{y}')$$

where $m \in \mathbb{R}_+$ is the margin. We assume $m = 1$. Prove that the following formulation is equivalent:

$$\ell_{aug}(\mathbf{y}, \mathbf{w}) = -\mathbf{w}^\top \mathbf{y} + \max_{\mathbf{y}' \in E(k)} (\mathbf{w} + (\mathbf{1} - \mathbf{y}))^\top \mathbf{y}'$$

we call the maximization problem in this loss "loss-augmented inference".

3. The hinge-loss function is not differentiable. However, we saw in the course that we can "approximate" the gradient of a function defined as a maximum of differentiable functions. More precisely, let $f(\mathbf{u}) = \max(f_1(\mathbf{u}), f_2(\mathbf{u}), \dots, f_n(\mathbf{u}))$ be a convex function where each f_i is convex and differentiable. Let $I(\mathbf{u}) = \{i \in \{1..n\} | f_i(\mathbf{u}) = f(\mathbf{u})\}$ the set of function indices reaching the maximum value for input \mathbf{u} . Then, $\forall i \in I(\mathbf{u}) : \nabla f_i(\mathbf{u}) \in \partial f(\mathbf{u})$, i.e. the gradient of any function $f_i(\mathbf{u})$ s.t. $i \in I(\mathbf{u})$ is a subgradient of $f(\mathbf{u})$. Using this fact, compute a subgradient of the hinge loss for multiclass classification wrt to parameters \mathbf{A} and \mathbf{b} .

$$\ell_{aug}(\mathbf{y}, \mathbf{A}\mathbf{x} + \mathbf{b})$$

This subgradient will be used instead of the gradient for learning. We abuse notation and note these subgradient as gradients.

2 Linear regression

To study linear regression, it is often useful to represent data in a matrix-vector notation. Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ be a vector containing the input feature of datapoints, one datapoint per row, and $\mathbf{y} \in \mathbb{R}^n$ be the associated output values. We assume \mathbf{X} and \mathbf{y} are our training data. If we ignore the intercept term (or add it as a feature), learning an unregularized linear regression model reduces to solving the following optimization problem:

$$\mathbf{a}^* = \arg \min_{\mathbf{a} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{X}\mathbf{a} - \mathbf{y}\|_2^2$$

1. Use first-order optimality condition to compute a closed form solution for the training problem. Under what conditions is this approach feasible? Under what conditions is this approach computationally challenging?
2. A matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ is invertible if and only if its columns are linearly independent, i.e. there exist no $\mathbf{v} \in \mathbb{R}^d$, $\mathbf{v} \neq 0$, such that $\mathbf{X}\mathbf{v} = 0$. Show that if \mathbf{X} is not invertible, then $\mathbf{X}^\top \mathbf{X}$ is not invertible.
3. We denote \mathbf{X}_i the i -th row of matrix \mathbf{X} . We assume each row consists of the value 1 followed by a $d - 1$ one-hot vector, i.e. the one-hot encoding of a $d - 1$ categorical feature, i.e. the first column of \mathbf{X} correspond to an implicit bias term. Prove that the column of \mathbf{X} are not linearly independent. What can we deduce?

4. We now assume the following regularized training problem:

$$\mathbf{a}^* = \arg \min_{\mathbf{a} \in \mathbb{R}^d} \|\mathbf{X}\mathbf{a} - \mathbf{y}\|_2^2 + \beta \sum_{i=2}^d \mathbf{a}_i^2$$

Warning: note that we don't regularize the weight in \mathbf{a} associated with the implicit bias feature and we removed the $\frac{1}{2}$ term to simplify computation.

Prove that, with this additional regularization term, the problem now has a closed form solution. **Hint:** you need to rewrite the problem so it "looks like" an unregularized problem.

5. Let D be a training dataset and consider the following training problem:

$$\arg \min_{\mathbf{a} \in \mathbb{R}^d, b} \frac{1}{|D|} \sum_{(\mathbf{x}, y) \in D} \frac{1}{2} (y - (\langle \mathbf{a}, \mathbf{x} \rangle + b))^2$$

Derive the closed form solution for updates in the coordinate descent algorithm.

6. Derive the coordinate descent updates with additional L2 regularization in the objective.