

CS 446 MJT — Homework 3

your NetID here

Version 1

Instructions.

- Homework is due **Tuesday, March 12, at 11:59pm**; no late homework accepted.
- Everyone must submit individually at gradescope under **hw3** and **hw3code**.
- The “written” submission at **hw3 must be typed**, and submitted in any format gradescope accepts (to be safe, submit a PDF). You may use L^AT_EX, markdown, google docs, MS word, whatever you like; but it must be typed!
- When submitting at **hw3**, gradescope will ask you to mark out boxes around each of your answers; please do this precisely!
- Please make sure your NetID is clear and large on the first page of the homework.
- Your solution **must** be written in your own words. Please see the course webpage for full academic integrity information. Briefly, you may have high-level discussions with at most 3 classmates, whose NetIDs you should place on the first page of your solutions, and you should cite any external reference you use; despite all this, your solution must be written in your own words.
- We reserve the right to reduce the auto-graded score for **hw3code** if we detect funny business (e.g., rather than implementing an algorithm, you keep re-submitting the assignment to the auto-grader, eventually completing a binary search for the answers).
- There are **no regrade requests** on **hw3code**, which is the code auto-grader; however, you can re-submit and re-grade as many times as you like before the deadline! Start early and report any issues on piazza!
- Methods and functions in the template and utility code include docstrings describing the inputs and outputs. The autograder relies on correct implementations of these methods. Follow the docstrings to avoid failing tests.

1. **The ln-sum-exp and cross entropy.**

- (a) Given $\mathbf{z} \in \mathbb{R}^k$, prove that $g(\mathbf{z}) = \ln \sum_{j=1}^k \exp(z_j)$ is convex.

Hint: prove that the Hessian matrix is positive semi-definite, meaning $\nabla^2 g(\mathbf{z}) \succeq 0$.

- (b) Recall that given a data example (\mathbf{x}, y) where $\mathbf{x} \in \mathbb{R}^d$ and $y \in \{1, 2, \dots, k\}$, and a classifier $f : \mathbb{R}^d \rightarrow \mathbb{R}^k$, the cross entropy loss is defined as

$$\ell_{\text{ce}}(f(\mathbf{x}), y) = -\ln \left(\frac{\exp(f(\mathbf{x})_y)}{\sum_{j=1}^k \exp(f(\mathbf{x})_j)} \right) = -f(\mathbf{x})_y + \ln \sum_{j=1}^k \exp(f(\mathbf{x})_j).$$

Let data examples $((\mathbf{x}_i, y_i))_{i=1}^n$ be given, where $\mathbf{x}_i \in \mathbb{R}^d$ and $y_i \in \{1, 2, \dots, k\}$. Consider the linear predictor $\mathbf{x} \mapsto \mathbf{W}\mathbf{x}$, where $\mathbf{W} \in \mathbb{R}^{k \times d}$. Prove that the empirical risk

$$\widehat{\mathcal{R}}(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \ell_{\text{ce}}(\mathbf{W}\mathbf{x}_i, y_i)$$

is convex.

Hint: use part (a) and the fact that convexity is preserved under affine composition and nonnegative combination. (It doesn't matter that this part uses matrix variables; this affine composition property holds for convex functions over matrices, and then when applying the previous part, its input is a vector after the affine transformation.)

- (c) For $\mathbf{z} \in \mathbb{R}^k$ and $r > 0$, let $g_r(\mathbf{z}) = \frac{1}{r} \ln \sum_{j=1}^k \exp(rz_j)$. Prove that

$$\lim_{r \rightarrow \infty} g_r(\mathbf{z}) = \max_{1 \leq j \leq k} z_j.$$

- (d) As a corollary, for $z \in \mathbb{R}$ and $r > 0$, the logistic loss $\ell(z) = \ln(1 + \exp(-z))$, and $\ell_r(z) = \frac{1}{r} \ln(1 + \exp(-rz))$, prove that

$$\lim_{r \rightarrow \infty} \ell_r(z) = \max\{0, -z\} = \text{ReLU}(-z).$$

Solution. (Your solution here.)

2. On initialization.

Consider a 2-layer network

$$f(\mathbf{x}; \mathbf{W}, \mathbf{v}) = \sum_{j=1}^m v_j \sigma(\langle \mathbf{w}_j, \mathbf{x} \rangle),$$

where $\mathbf{x} \in \mathbb{R}^d$, $\mathbf{W} \in \mathbb{R}^{m \times d}$ with rows \mathbf{w}_j^\top , and $\mathbf{v} \in \mathbb{R}^m$. For simplicity, the network has a single output, and bias terms are omitted.

Given a data example (\mathbf{x}, y) and a loss function ℓ , consider the empirical risk

$$\widehat{\mathcal{R}}(\mathbf{W}, \mathbf{v}) = \ell(f(\mathbf{x}; \mathbf{W}, \mathbf{v}), y).$$

Only a single data example will be considered in this problem; the same analysis extends to multiple examples by taking averages.

- (a) For each $1 \leq j \leq m$, derive $\partial \widehat{\mathcal{R}} / \partial v_j$ and $\partial \widehat{\mathcal{R}} / \partial \mathbf{w}_j$.
- (b) Consider gradient descent which starts from some $\mathbf{W}^{(0)}$ and $\mathbf{v}^{(0)}$, and at step $t \geq 0$, updates the weights for each $1 \leq j \leq m$ as follows:

$$\mathbf{w}_j^{(t+1)} = \mathbf{w}_j^{(t)} - \eta \frac{\partial \widehat{\mathcal{R}}}{\partial \mathbf{w}_j^{(t)}}, \quad \text{and} \quad v_j^{(t+1)} = v_j^{(t)} - \eta \frac{\partial \widehat{\mathcal{R}}}{\partial v_j^{(t)}}.$$

Suppose there exists two hidden units $p, q \in \{1, 2, \dots, m\}$ such that $\mathbf{w}_p^{(0)} = \mathbf{w}_q^{(0)}$ and $v_p^{(0)} = v_q^{(0)}$. Prove by induction that for any step $t \geq 0$, it holds that $\mathbf{w}_p^{(t)} = \mathbf{w}_q^{(t)}$ and $v_p^{(t)} = v_q^{(t)}$.

Remark: as a result, if the neural network is initialized symmetrically, then such a symmetry may persist during gradient descent, and thus the representation power of the network will be limited.

- (c) Random initialization is a good way to break symmetry. Moreover, proper random initialization also preserves the squared norm of the input, as formalized below.

First consider the identity activation $\sigma(z) = z$. For each $1 \leq j \leq m$ and $1 \leq k \leq d$, initialize $w_{j,k}^{(0)} \sim \mathcal{N}(0, 1/m)$ (i.e., normal distribution with mean $\mu = 0$ and variance $\sigma^2 = 1/m$). Prove that

$$\mathbb{E} \left[\left\| \mathbf{W}^{(0)} \mathbf{x} \right\|_2^2 \right] = \|\mathbf{x}\|_2^2.$$

Next consider the ReLU activation $\sigma_r(z) = \max\{0, z\}$. For each $1 \leq j \leq m$ and $1 \leq k \leq d$, initialize $w_{j,k}^{(0)} \sim \mathcal{N}(0, 2/m)$. Prove that

$$\mathbb{E} \left[\left\| \sigma_r(\mathbf{W}^{(0)} \mathbf{x}) \right\|_2^2 \right] = \|\mathbf{x}\|_2^2.$$

Hint: linear combinations of Gaussians are again Gaussian! For the second part (with ReLU), consider the symmetry of a Gaussian around 0.

Solution. (Your solution here.)

3. ResNet.

In this problem, you will implement a simplified ResNet. You do not need to change arguments which are not mentioned here (but you of course could try and see what happens).

- (a) Implement a class `Block`, which is a building block of ResNet. It is described in (He et al., 2016) Figure 2.

The input to `Block` is of shape (N, C, H, W) , where N denotes the batch size, C denotes the number of channels, and H and W are the height and width of each channel. For each data example \mathbf{x} with shape (C, H, W) , the output of `block` is

$$\text{Block}(\mathbf{x}) = \sigma_r(\mathbf{x} + f(\mathbf{x})),$$

where σ_r denotes the ReLU activation, and $f(\mathbf{x})$ also has shape (C, H, W) and thus can be added to \mathbf{x} . In detail, f contains the following layers.

- i. A `Conv2d` with C input channels, C output channels, kernel size 3, stride 1, padding 1, and no bias term.
- ii. A `BatchNorm2d` with C features.
- iii. A ReLU layer.
- iv. Another `Conv2d` with the same arguments as i above.
- v. Another `BatchNorm2d` with C features.

Because 3×3 kernels and padding 1 are used, the convolutional layers do not change the shape of each channel. Moreover, the number of channels are also kept unchanged. Therefore $f(\mathbf{x})$ does have the same shape as \mathbf{x} .

Additional instructions are given in docstrings in `hw3.py`.

- (b) Explain why a `Conv2d` layer does not need a bias term if it is followed by a `BatchNorm2d` layer.
- (c) Implement a (shallow) `ResNet` consists of the following parts:
- i. A `Conv2d` with 1 input channel, C output channels, kernel size 3, stride 2, padding 1, and no bias term.
 - ii. A `BatchNorm2d` with C features.
 - iii. A ReLU layer.
 - iv. A `MaxPool2d` with kernel size 2.
 - v. A `Block` with C channels.
 - vi. An `AdaptiveAvgPool2d` which for each channel takes the average of all elements.
 - vii. A `Linear` with C inputs and 10 outputs.

Additional instructions are given in docstrings in `hw3.py`.

- (d) Train a `ResNet` with 16 channels on the data given by `hw3_utils.torch_digits()`, using the cross entropy loss and SGD with learning rate 0.005 and batch size 16, for 30 epochs. You can just use your `fit_and_validate()` in `hw2`. Plot the epochs vs training and validation cross entropy losses. Since there is some inconsistency due to random initialization, try 3 runs and have 3 plots.

Solution. (*Your solution here.*)

4. Kernel properties.

- (a) Let $\mathbf{A} \in \mathbb{R}^{n \times n}$ denote a symmetric positive semi-definite matrix of rank r . Prove that there exists n vectors $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_n \in \mathbb{R}^r$, such that $A_{i,j} = \langle \mathbf{z}_i, \mathbf{z}_j \rangle$.

Hint: use the eigendecomposition $\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top$.

- (b) On the other hand, given n vectors $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_n \in \mathbb{R}^m$, prove that the matrix \mathbf{A} defined by $A_{i,j} = \langle \mathbf{z}_i, \mathbf{z}_j \rangle$ is positive semi-definite.

Remark: note that the rank of \mathbf{A} is at most m , and it could be strictly less than m . In particular, m could be larger than n .

- (c) Using (a) and (b), prove that if $K_1(\mathbf{x}, \mathbf{y})$ and $K_2(\mathbf{x}, \mathbf{y})$ are kernels, then $K(\mathbf{x}, \mathbf{y}) = K_1(\mathbf{x}, \mathbf{y})K_2(\mathbf{x}, \mathbf{y})$ is also a kernel.
- (d) Using (c), prove that $K(\mathbf{x}, \mathbf{y}) = (1 + \langle \mathbf{x}, \mathbf{y} \rangle)^r$ is a kernel for any positive integer r . (It is the polynomial kernel with degree r .)
- (e) Assume $K(\mathbf{x}, \mathbf{y}) = \exp(-\|\mathbf{x} - \mathbf{y}\|_2^2 / 2\sigma^2)$ is a kernel in the 1-dimensional case (i.e., $x, y \in \mathbb{R}$). Using (c), prove that $K(\mathbf{x}, \mathbf{y})$ is indeed a kernel for any dimension. (It is the Gaussian / RBF kernel.)

Solution. (*Your solution here.*)

5. RBF kernel and nearest neighbors.

- (a) Recall that given data examples $((\mathbf{x}_i, y_i))_{i=1}^n$ and an optimal dual solution $(\hat{\alpha}_i)_{i=1}^n$, the RBF kernel SVM makes a prediction as follows:

$$f_\sigma(\mathbf{x}) = \sum_{i=1}^n \hat{\alpha}_i y_i \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}_i\|_2^2}{2\sigma^2}\right) = \sum_{i \in S} \hat{\alpha}_i y_i \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}_i\|_2^2}{2\sigma^2}\right),$$

where $S \subset \{1, 2, \dots, n\}$ is the set of indices of support vectors.

Given an input \mathbf{x} , let $T := \arg \min_{i \in S} \|\mathbf{x} - \mathbf{x}_i\|_2$ denote the set of closest support vectors to \mathbf{x} , and let $\rho := \min_{i \in S} \|\mathbf{x} - \mathbf{x}_i\|_2$ denote this smallest distance. (In other words, $T := \{i \in S : \|\mathbf{x} - \mathbf{x}_i\| = \rho\}$.) Prove that

$$\lim_{\sigma \rightarrow 0} \frac{f_\sigma(\mathbf{x})}{\exp(-\rho^2/2\sigma^2)} = \sum_{i \in T} \hat{\alpha}_i y_i.$$

Remark: in other words, when the bandwidth σ becomes small enough, RBF kernel SVM is almost the 1-nearest neighbor predictor with the set of support vectors as the training set.

- (b) Consider the XOR dataset:

$$\begin{aligned} \mathbf{x}_1 &= (+1, +1), & y_1 &= +1, \\ \mathbf{x}_2 &= (-1, +1), & y_2 &= -1, \\ \mathbf{x}_3 &= (-1, -1), & y_3 &= +1, \\ \mathbf{x}_4 &= (+1, -1), & y_4 &= -1. \end{aligned}$$

Verify that $\hat{\alpha} = (1/\alpha, 1/\alpha, 1/\alpha, 1/\alpha)$ is an optimal dual solution to the RBF kernel SVM, where

$$\alpha = \left(1 - \exp\left(-\frac{\|\mathbf{x}_1 - \mathbf{x}_2\|_2^2}{2\sigma^2}\right)\right)^2 = \left(1 - \exp\left(-\frac{2}{\sigma^2}\right)\right)^2 > 0.$$

Hint: prove that the gradient of the dual function is $\mathbf{0}$ at $\hat{\alpha}$. Since the dual function is concave, and $\hat{\alpha} > \mathbf{0}$, it follows that $\hat{\alpha}$ is an optimal dual solution.

Remark: in other words, all four data examples are mapped to support vectors in the reproducing kernel Hilbert space. In light of (a), when σ is small enough, $f_\sigma(\mathbf{x})$ is almost the 1-nearest neighbor predictor on the XOR dataset. In fact, it is also true for large σ , due to the symmetry of the XOR data.

Solution. (Your solution here.)

6. SVM implementation.

Recall that the dual problem of SVM is

$$\max_{\boldsymbol{\alpha} \in \mathcal{C}} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j).$$

where the domain $\mathcal{C} = [0, \infty)^n = \{\boldsymbol{\alpha} : \alpha_i \geq 0\}$ for hard-margin SVM, and $\mathcal{C} = [0, C]^n = \{\boldsymbol{\alpha} : 0 \leq \alpha_i \leq C\}$ for soft-margin SVM.

Equivalently, it can be formulated as a minimization problem

$$\min_{\boldsymbol{\alpha} \in \mathcal{C}} f(\boldsymbol{\alpha}) := \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) - \sum_{i=1}^n \alpha_i.$$

It can be solved by projected gradient descent, which starts from some $\boldsymbol{\alpha}_0 \in \mathcal{C}$ (e.g., $\mathbf{0}$) and updates as follows

$$\boldsymbol{\alpha}_{t+1} = \Pi_{\mathcal{C}} [\boldsymbol{\alpha}_t - \eta \nabla f(\boldsymbol{\alpha}_t)].$$

Here $\Pi_{\mathcal{C}}[\boldsymbol{\alpha}]$ is the *projection* of $\boldsymbol{\alpha}$ onto \mathcal{C} , defined as the closet point to $\boldsymbol{\alpha}$ in \mathcal{C} :

$$\Pi_{\mathcal{C}}[\boldsymbol{\alpha}] := \arg \min_{\boldsymbol{\alpha}' \in \mathcal{C}} \|\boldsymbol{\alpha}' - \boldsymbol{\alpha}\|_2.$$

If \mathcal{C} is convex, the projection is uniquely defined.

(a) Prove that

$$\left(\Pi_{[0, \infty)^n}[\boldsymbol{\alpha}] \right)_i = \max\{\alpha_i, 0\},$$

and

$$\left(\Pi_{[0, C]^n}[\boldsymbol{\alpha}] \right)_i = \min\{\max\{0, \alpha_i\}, C\}.$$

- (b) Implement an `svm_solver()`, using projected gradient descent formulated as above. See the docstrings in `hw3.py` for details.
- (c) Implement an `svm_predictor()`, using an optimal dual solution, the training set, and an input. See the docstrings in `hw3.py` for details.
- (d) On the area $[-5, 5] \times [-5, 5]$, plot the contour lines of the following kernel SVMs, trained on the XOR data. Different kernels and the XOR data are provided in `hw3_utils.py`. Learning rate 0.1 and 10000 steps should be enough. To draw the contour lines, you can use `hw3.svm_contour()`.
- The polynomial kernel with degree 2.
 - The RBF kernel with $\sigma = 1$.
 - The RBF kernel with $\sigma = 2$.
 - The RBF kernel with $\sigma = 4$.

Solution. (*Your solution here.*)

References

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.