

CS 446 / ECE 449 — Homework 2

ivette2

Version 1.0

Instructions.

- Homework is due **Tuesday, February 22, at noon CST**; no late homework accepted.
- Everyone must submit individually on Gradescope under **hw2** and **hw2code**.
- The “written” submission at **hw2** **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 **hw2**, Gradescope will ask you to select pages for each problem; 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 **hw2code** if we detect funny business (e.g., your solution lacks any algorithm and hard-codes answers you obtained from someone else, or simply via trial-and-error with the autograder).
- Coding problems come with suggested “library routines”; we include these to reduce your time fishing around APIs, but you are free to use other APIs.
- When submitting to **hw2code**, only upload the two python files **hw2.py** and **hw2_utils.py**. Don’t upload a zip file or additional files.

Version history.

1.0. Initial version.

1. SVM with Biases.

This problem is about SVMs over \mathbb{R}^d with linearly separable data (i.e., the hard margin SVM).

Our formulation of SVM required separators to pass through the origin, which does not provide a geometrically pleasing notion of maximum margin direction.

A first fix is provided by lecture 4: by appending a 1 to the inputs, we obtain the convex program

$$\begin{aligned} \min_{\mathbf{u}} \quad & \frac{1}{2} \|\mathbf{u}\|^2 \\ \text{subject to} \quad & \mathbf{u} \in \mathbb{R}^{d+1} \\ & y_i \begin{bmatrix} \mathbf{x}_i \\ 1 \end{bmatrix}^\top \mathbf{u} \geq 1 \quad \forall i, \end{aligned}$$

and let $\bar{\mathbf{u}}$ denote the optimal solution to this program.

A second standard fix is to incorporate the bias directly into the optimization problem:

$$\begin{aligned} \min_{\mathbf{v}, b} \quad & \frac{1}{2} \|\mathbf{v}\|^2 \\ \text{subject to} \quad & \mathbf{v} \in \mathbb{R}^d, b \in \mathbb{R} \\ & y_i (\mathbf{v}^\top \mathbf{x}_i + b) \geq 1 \quad \forall i, \end{aligned}$$

and let $(\bar{\mathbf{v}}, \bar{b}) \in \mathbb{R}^d \times \mathbb{R}$ denote an optimal solution to this program. This second version is standard, but we do not use it in lecture for various reasons.

- (a) In lecture, we stated that the first formulation is a *convex program* (formally defined in lecture 5). Show that the second formulation is also a convex program.

To be considered a convex program, we must show that the second formulation is minimizing a convex function over a convex set. $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if for any $x, x' \in \mathbb{R}^d$ and $\alpha \in [0, 1]$, $f((1-\alpha)x + \alpha x') \leq (1-\alpha)f(x) + \alpha f(x')$. For the second formulation, we are minimizing $\frac{1}{2} \|\mathbf{v}\|^2$, which gives us the inequality $\frac{1}{2} \left\| ((1-\alpha)\mathbf{v} + \alpha \mathbf{v}') \right\|^2 \leq (1-\alpha) \frac{1}{2} \|\mathbf{v}\|^2 + \alpha \frac{1}{2} \|\mathbf{v}'\|^2$. Given the triangle equality for real norms, $\|x + y\| \leq \|x\| + \|y\|$, and the fact that $\alpha \in [0, 1]$, we can say that $\left\| ((1-\alpha)\mathbf{v} + \alpha \mathbf{v}') \right\| \leq (1-\alpha) \|\mathbf{v}\| + \alpha \|\mathbf{v}'\|$ and $\left\| ((1-\alpha)\mathbf{v} + \alpha \mathbf{v}') \right\|^2 \leq ((1-\alpha) \|\mathbf{v}\| + \alpha \|\mathbf{v}'\|)^2$. This becomes $\left\| ((1-\alpha)\mathbf{v} + \alpha \mathbf{v}') \right\|^2 \leq ((1-\alpha)^2 \|\mathbf{v}\|^2 + 2(\alpha - \alpha^2) \|\mathbf{v}\| \|\mathbf{v}'\| + \alpha^2 \|\mathbf{v}'\|^2)$, which proves the initial inequality and the function to be convex.

To prove the set convex, we can use the definition: For every pair of points $\{x, x'\} \in S$, $[x, x'] = \{\alpha x + (1-\alpha)x' : \alpha \in [0, 1]\} \subseteq S$. We can use the subjected condition and prove that

$$\begin{aligned} y_i \left((\alpha \mathbf{v} + (1-\alpha)\mathbf{v}')^\top \mathbf{x}_i + (\alpha b + (1-\alpha)b') \right) \geq 1 \text{ has } (\mathbf{v}, b) \subseteq S. \text{ This expands out to} \\ y_i \left(\alpha \mathbf{v}^\top \mathbf{x}_i + (1-\alpha) (\mathbf{v}')^\top \mathbf{x}_i + \alpha b + (1-\alpha)b' \right), \text{ where we can look at each term and say } \alpha \mathbf{v}^\top \mathbf{x}_i \geq 0, \\ (1-\alpha) (\mathbf{v}')^\top \mathbf{x}_i \geq (\mathbf{v}')^\top \mathbf{x}_i, \alpha b \geq 0 \text{ and } (1-\alpha)b' \geq b', \text{ all given that } \alpha \in [0, 1]. \text{ This proves that} \\ y_i \left(\alpha \mathbf{v}^\top \mathbf{x}_i + (1-\alpha) (\mathbf{v}')^\top \mathbf{x}_i + \alpha b + (1-\alpha)b' \right) \geq y_i ((\mathbf{v}')^\top \mathbf{x}_i + b') \geq 1, \text{ and we have a convex set.} \end{aligned}$$

- (b) Suppose there is only one datapoint: $\mathbf{x}_1 = \mathbf{e}_1$, the first standard basis vector, with label $y_1 = +1$. The first formulation will have a unique solution $\bar{\mathbf{u}}$, as discussed in lecture. Show that the second formulation does not have a unique solution.

If we have only one data point with $x_1 = e_1, y_1 = +1$, we are subject to condition $y_i (v^T x_i + b) \rightarrow v_1 + b \geq 1$. With this, $\min_{v,b} \frac{1}{2} \|v\|^2$ will be minimized for any $v = 0$ and $b \geq 1$, resulting in infinite solutions.

- (c) Let's add another datapoint: $x_2 = -ae_1$ for some $a \geq 3$, with label $y_2 = -1$. Now that we have two data points, both of the convex programs now have two constraints. Write out the explicit constraints to the first convex program.

Adding another data point, $x_2 = -ae_1, a \geq 3, y_1 = -1$. The constraints for the first program are $u \in \mathbb{R}^{d+1}, y_i \begin{bmatrix} x_i \\ 1 \end{bmatrix}^T u \geq 1$. Writing these out for the two data points, we have $1 \begin{bmatrix} e_1 \\ 1 \end{bmatrix}^T u \geq 1$ and $-1 \begin{bmatrix} -ae_1 \\ 1 \end{bmatrix}^T u \geq 1$. Knowing that e_1 is a standard basis vector, we can explicitly write these constraints as $u_1 + u_{d+1} \geq 1$ and $au_1 - u_{d+1} \geq 1$, respectively.

- (d) Using these two constraints, show that the first coordinate \bar{u}_1 of the optimal solution $\bar{\mathbf{u}}$ satisfies $\bar{u}_1 \geq \frac{2}{a+1}$.

If we simply combine the two constraint inequalities on an optimal solution \bar{u}_1 , we have $u_1 + au_1 \geq 2$. This is the same as saying $u_1(1+a) \geq 2$, and thus $u_1 \geq \frac{2}{a+1}$.

- (e) Using parts (c) and (d), find optimal solutions $\bar{\mathbf{u}}$ and $(\bar{\mathbf{v}}, \bar{b})$, and prove they are in fact optimal.
Hint: If you are stuck, first try the case $d = 1$. Then study what happens for $d = 2, d = 3, \dots$
Hint: $(\bar{\mathbf{v}}, \bar{b})$ will be unique.

The optimal solutions $\bar{\mathbf{u}}$ and $(\bar{\mathbf{v}}, \bar{b})$ can be determined using the constraints and the minimization problems themselves, $\min_{\mathbf{u}} \frac{1}{2} \|\mathbf{u}\|^2$ and $\min_{\mathbf{v}, b} \frac{1}{2} \|\mathbf{v}\|^2$. For $\bar{\mathbf{u}}$, the norm is minimized with the smallest possible components \bar{u}_i , so $\bar{\mathbf{u}} = \left(\frac{2}{a+1}, 0, \dots, \frac{a-1}{a+1} \right)$, as these are the smallest components that still satisfy the constraints.

For $(\bar{\mathbf{v}}, \bar{b})$, we can take the same approach. The components that minimize the norm of \mathbf{v} go as $(\bar{\mathbf{v}}) = \left(\frac{1}{1+a}, 0, \dots \right)$, and $(\bar{b}) = \left(1 - \frac{a}{1+a} \right)$.

- (f) Now we will consider the behavior of $\bar{\mathbf{u}}$ and $\bar{\mathbf{v}}$ as a increases; to this end, write $\bar{\mathbf{u}}_a$ and $\bar{\mathbf{v}}_a$, and consider $a \rightarrow \infty$. Determine and formally prove the limiting behavior of $\lim_{a \rightarrow \infty} \frac{1}{2} \|\bar{\mathbf{u}}_a\|^2$ and $\lim_{a \rightarrow \infty} \frac{1}{2} \|\bar{\mathbf{v}}_a\|^2$.
Hint: The two limits will not be equal.

$$\lim_{a \rightarrow \infty} \frac{1}{2} \|\bar{\mathbf{u}}_a\|^2 = (0, 0, \dots, 1), \text{ and } \lim_{a \rightarrow \infty} \frac{1}{2} \|\bar{\mathbf{v}}_a\|^2 = (0, 0, \dots).$$

- (g) Between the two versions of SVM with bias, which do you prefer? Any answer which contains at least one complete sentence will receive full credit.

Remark: Initially it may have seemed that both optimization problems have the same solutions; the purpose of this problem was to highlight that small differences in machine learning methods can lead

to observably different performance.

I prefer the second convex program, obtained by incorporating the bias into the optimization problem. I thought the notion of infinite solutions for insufficient data (like the one datapoint proposed in 1b) felt more intuitive than the first formulation with appended ones.

Solution.

2. SVM Implementation.

Recall that the dual problem of an SVM is

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

where the domain $\mathcal{C} = [0, \infty)^n = \{\alpha : \alpha_i \geq 0\}$ for a hard-margin SVM, and $\mathcal{C} = [0, C]^n = \{\alpha : 0 \leq \alpha_i \leq C\}$ for a soft-margin SVM. Equivalently, we can frame this as the minimization problem

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

This can be solved by projected gradient descent, which starts from some $\alpha_0 \in \mathcal{C}$ (e.g., $\mathbf{0}$) and updates via

$$\alpha_{t+1} = \Pi_{\mathcal{C}} [\alpha_t - \eta \nabla f(\alpha_t)],$$

where $\Pi_{\mathcal{C}}[\alpha]$ is the *projection* of α onto \mathcal{C} , defined as the closest point to α in \mathcal{C} :

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

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

(a) Prove that

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

and

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

Hint: Show that the i th component of any other $\alpha' \in \mathcal{C}$ is further from the i th component of α than the i th component of the projection is. Specifically, show that $|\alpha'_i - \alpha_i| \geq |\max\{0, \alpha_i\} - \alpha_i|$ for $\alpha' \in [0, \infty)^n$ and that $|\alpha'_i - \alpha_i| \geq |\min\{\max\{0, \alpha_i\}, C\} - \alpha_i|$ for $\alpha' \in [0, C]^n$.

Looking at $|\alpha'_i - \alpha_i| \geq |\max\{0, \alpha_i\} - \alpha_i|$ for $\alpha' \in [0, \infty)^n$, we know that $\alpha'_i \geq 0$ for any $\alpha_i \geq 0$. We can say that $\left(\Pi_{[0, \infty)^n}[\alpha] \right)_i = \begin{cases} \alpha_i & : \alpha_i \geq 0 \\ 0 & : \alpha_i < 0 \end{cases} = \max\{0, \alpha_i\}$

We can say that $|\alpha'_i - \alpha_i| \geq |\max\{0, \alpha_i\} - \alpha_i|$ for $\alpha' \in [0, \infty)^n$, because $\alpha'_i \geq 0$.

The same argument goes for $|\alpha'_i - \alpha_i| \geq |\min\{\max\{0, \alpha_i\}, C\} - \alpha_i|$ for $\alpha' \in [0, C]^n$.

(b) Implement an `svm_solver()`, using projected gradient descent formulated as above. Initialize your α to zeros. See the docstrings in `hw2.py` for details.

Remark: Consider using the `.backward()` function in pytorch. However, then you may have to use in-place operations like `clamp_()`, otherwise the gradient information is destroyed.

Library routines: `torch.outer`, `torch.clamp`, `torch.autograd.backward`, `torch.tensor(..., requires_grad=True)`, with `torch.no_grad():`, `torch.tensor.grad.zero_`, `torch.tensor.detach`.

(c) Implement an `svm_predictor()`, using an optimal dual solution, the training set, and an input. See the docstrings in `hw2.py` for details.

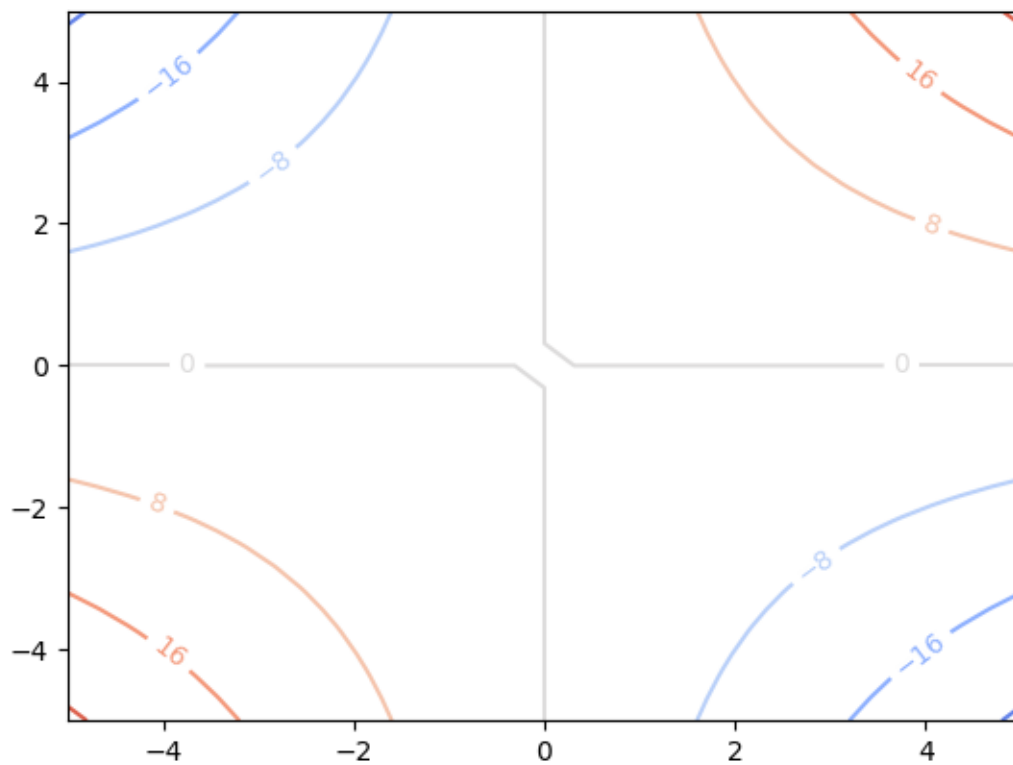
Library routines: `torch.empty`.

(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 `hw2_utils.py`. Learning rate 0.1 and 10000 steps should be enough. To draw the contour lines, you can use `hw2_utils.svm_contour()`.

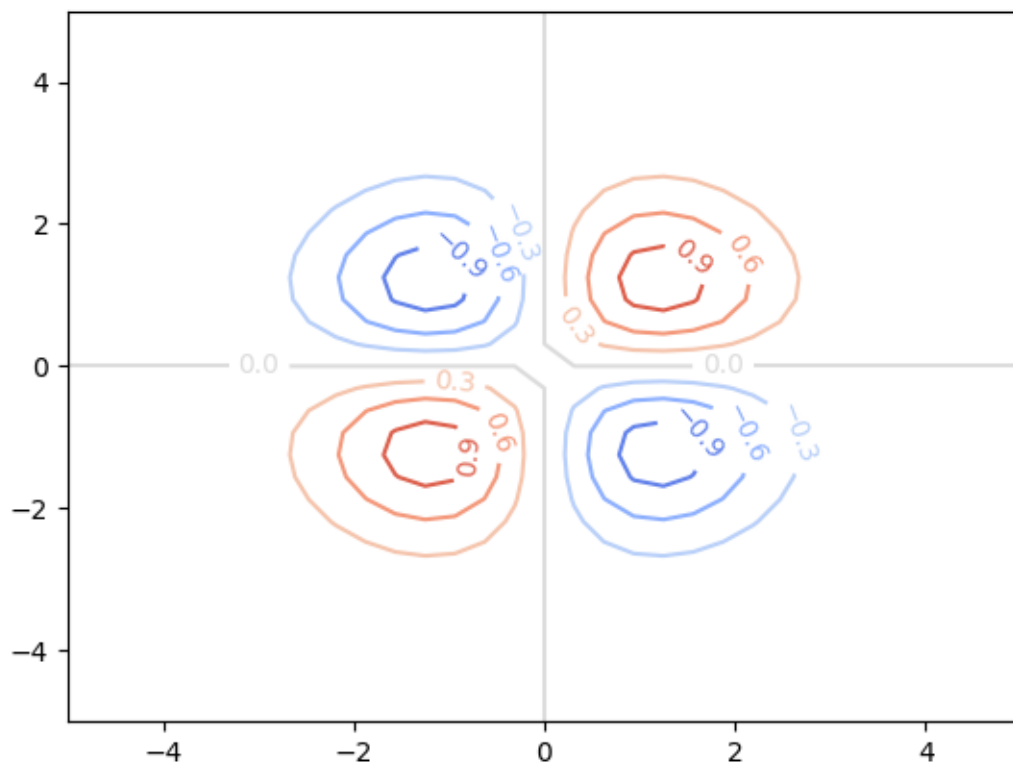
- The polynomial kernel with degree 2.
- The RBF kernel with $\sigma = 1$.
- The RBF kernel with $\sigma = 2$.
- The RBF kernel with $\sigma = 4$.

Include these four plots in your written submission.

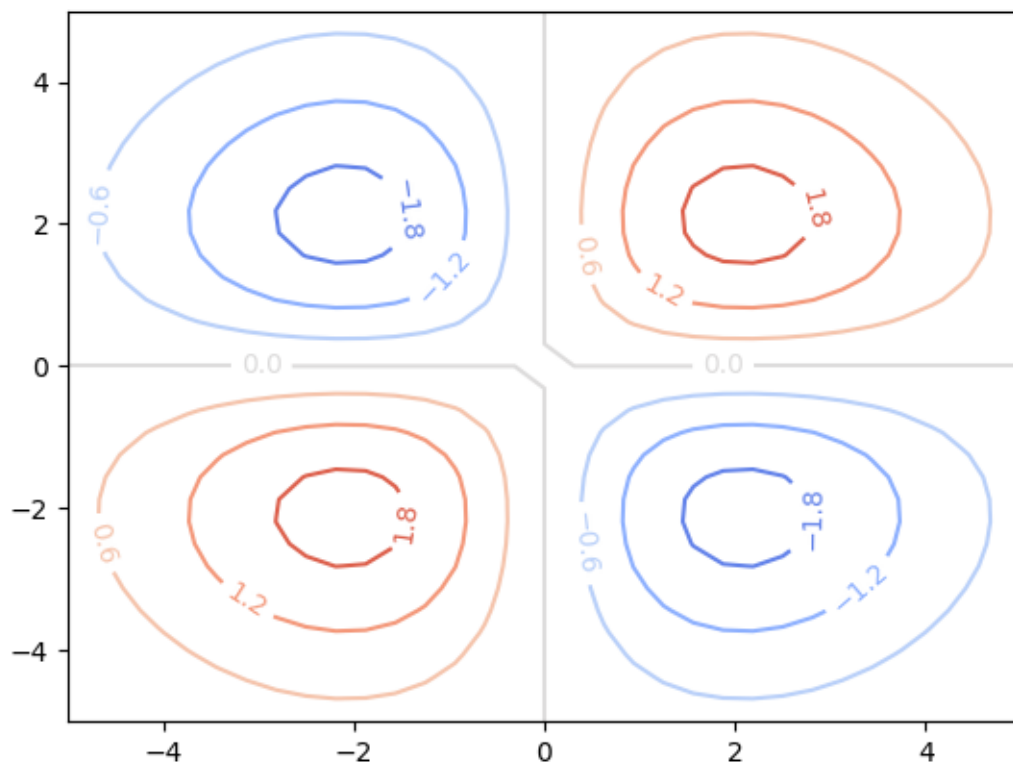
Polynomial kernel with degree 2



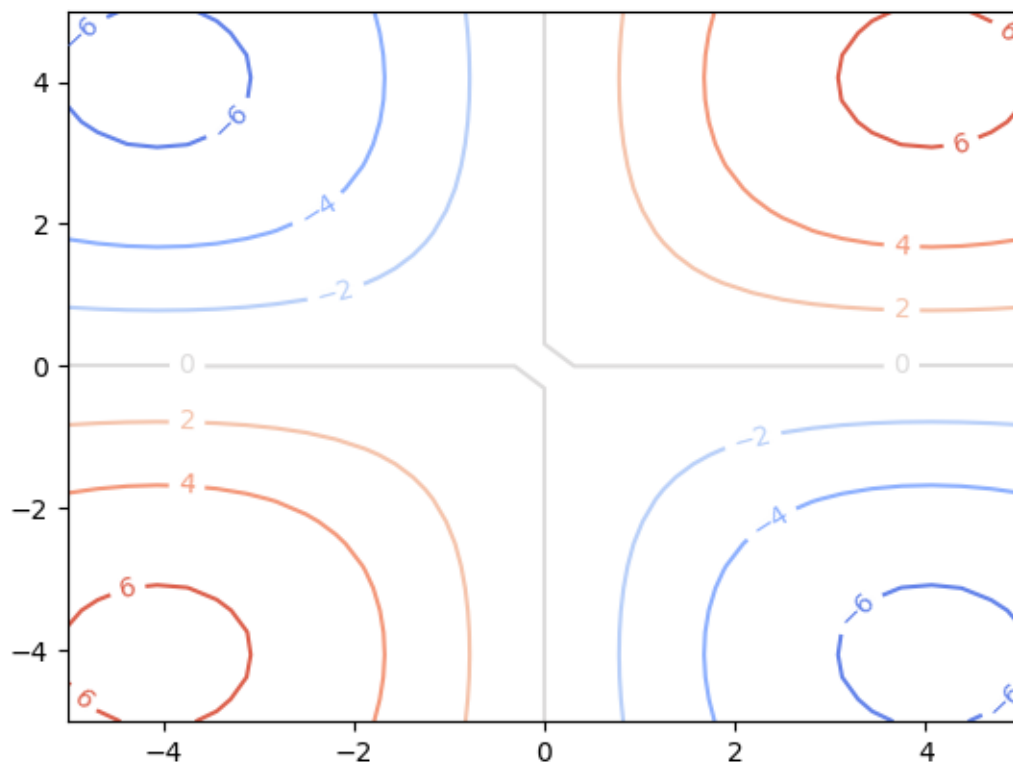
RBF kernel with $\sigma = 1$



RBF kernel with $\sigma = 2$



RBF kernel with $\sigma = 4$



Solution.

3. Neural Networks for Emotion Classification

In this problem you will build a single-layer neural network that classifies pictures into one of six categories: anger, disgusted, happy, maudlin, fear, and surprise. The CAFE ¹ dataset included in this homework's zip file provides a set of grayscale facial images expressing the described emotions. This will also serve as an introduction to writing your own neural networks in PyTorch! Consider the single layer neural network below

$$\mathbf{x} \mapsto \sigma(\mathbf{W}\mathbf{x} + \mathbf{b}),$$

where σ is the softmax activation and we use cross entropy loss to train the network.

- (a) Implement your network in the class CAFENet. You will need to modify the `__init__` and `forward` methods. Due to numerical issues, do not include an explicit softmax layer in your network. Instead, your implementation should output the raw logits (meaning $\mathbf{W}\mathbf{x} + \mathbf{b}$); in part (b), the network will be fit to data with `torch.nn.CrossEntropyLoss`, which implicitly applies a softmax as discussed in lecture. Refer to `IMAGE_DIMS`, `load_cafe`, and `get_cafe_data` in `hw2_utils.py` for how the images and labels will be passed to the network as tensors.

Library routines: `torch.nn.Linear`, `torch.nn.Module.forward`.

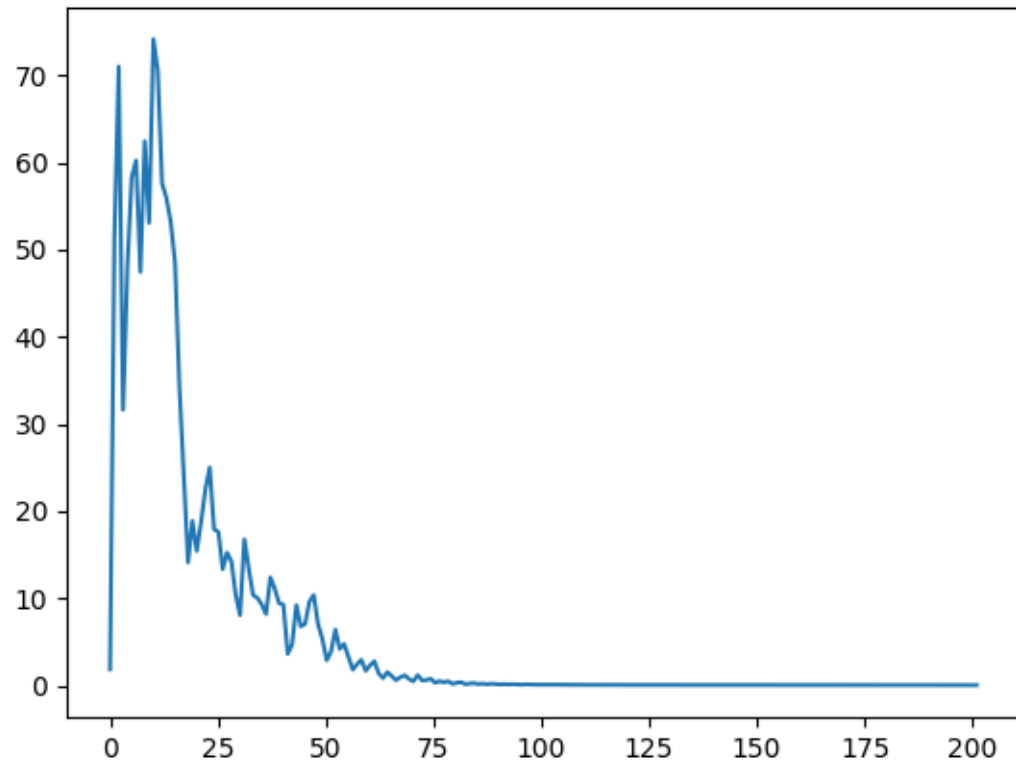
- (b) Implement `fit` to train the input network for `n_epochs` epochs. Use cross entropy loss and an Adam optimizer.

Library routines: `torch.nn.Module.forward`, `torch.nn.Loss.backward`, `torch.optim.Adam`, `torch.optim.Optimizer.step`, `torch.optim.Optimizer.zero_grad`, `torch.nn.CrossEntropyLoss`.

- (c) Implement and run the `plot_cafe_loss` function. Specifically, use `hw2_utils.get_cafe_data()` to load the training set, then train a CAFENet via your `fit` function for 201 epochs. Plot the empirical risk (in terms of cross entropy loss) across these first 201 epochs, and include the resulting plot in your written handin. Lastly, use `torch.save` to save your model in order to use it in the next two problem parts.

Library routines: `plt.plot`, `torch.save`.

¹Inspiration for this problem from Garrison Cottrell's neural networks course. See ? for more info on the CAFE dataset.



- (d) Let's see how well our model predicts labels. We will use a confusion matrix to visualize how well it does for each category. Implement `print_confusion_matrix` to print out two confusion matrices for your model, one on the training set and one on the test set. Use `hw2.utils.get_cafe_data("test")` to load the test set. Include both matrices in your writeup, along with 1-3 sentences discussing differences in the matrices and what might cause them.

Library routines: `torch.load`, `torch.argmax`, `sklearn.metrics.confusion_matrix`.

```
[[6 0 0 0 0 0]
 [0 6 0 0 0 0]
 [0 0 6 0 0 0]
 [0 0 0 6 0 0]
 [0 0 0 0 6 0]
 [0 0 0 0 0 6]]
[[2 0 0 1 0 0]
 [1 0 1 0 1 0]
 [0 0 2 0 1 0]
 [2 0 0 0 1 0]
 [0 0 0 0 2 1]
 [0 0 0 0 1 2]]
```

We can see that the first confusion matrix on the training set has all positives, as each diagonal entry $C_{i,i}$ is fully confident in the 6 classes. This makes sense given that the model was trained on the training data, so it performs well on predictions it has already seen. The second matrix for the test set represents significantly less confidence, as expected with data it hasn't seen before.

- (e) Now let's visualize the model's weights by implementing the `visualize_weights` method. For each of the 91,200-dimensional weights of your CAFENet's six output nodes, linearly map them to the grayscale range `[0, 255]` by performing the following transformations:
- Compute the minimum and maximum weights across all six output nodes, denoted `min_weight`, `max_weight` respectively.
 - Transform the weights `w` by `w = (w - min_weight) * 255 / (max_weight - min_weight)` to linearly map `w` into the range `[0, 255]`.
 - Cast the weights to integers.

Then, reshape the weights to the image dimensions `380 x 240` and plot them in grayscale. Include all six plots in your writeup. What do you see? Why might the weights appear this way?

Library routines: `torch.load`, `torch.nn.Module.parameters`, `torch.nn.tensor.min`, `torch.nn.tensor.max`, `torch.nn.tensor.int`, `torch.nn.tensor.reshape`, `torch.tensor.detach`, `plt.imshow(..., cmap='gray')`.

Note: In practice, for simple neural networks like this we would use `torch.nn.Sequential`.

Solution.

4. Shallow Network Random 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 $\partial \widehat{\mathcal{R}} / \partial \mathbf{w}_j$. Note that the first is a derivative with respect to a scalar (so the answer should be a scalar), and the second is a derivative with respect to a vector (so the answer should be a vector).

$$\begin{aligned} \widehat{\mathcal{R}}(\mathbf{W}, \mathbf{v}) &= \ell(f(\mathbf{x}; \mathbf{W}, \mathbf{v}), y) = \ell\left(\sum_{j=1}^m v_j \sigma(\langle \mathbf{w}_j, \mathbf{x} \rangle), y\right), \text{ and we can differentiate using chain rule} \\ \text{and linearity. } \partial \widehat{\mathcal{R}} / \partial v_j &= (\partial \ell / \partial f) (\partial / \partial v_j) = (\partial \ell / \partial f) \left(\sum_{j=1}^m \sigma(\langle \mathbf{w}_j, \mathbf{x} \rangle)\right) \\ \text{Similarly, and assuming } \sigma \text{ is an arbitrary activation function,} \\ \partial \widehat{\mathcal{R}} / \partial \mathbf{w}_j &= (\partial \ell / \partial f) (\partial f / \partial \sigma) (\partial \sigma / \partial \langle \mathbf{w}_j, \mathbf{x} \rangle) (\partial \langle \mathbf{w}_j, \mathbf{x} \rangle / \partial \mathbf{w}_j) = (\partial \ell / \partial f) (\partial \sigma / \partial \langle \mathbf{w}_j, \mathbf{x} \rangle) \left(\sum_{j=1}^m v_j\right) (\mathbf{x}^T) \end{aligned}$$

- (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 exist two hidden units $p, q \in \{1, 2, \dots, m\}$ and t such that $\mathbf{w}_p^{(t)} = \mathbf{w}_q^{(t)}$ and $v_p^{(t)} = v_q^{(t)}$. Show that $\mathbf{w}_p^{(t+1)} = \mathbf{w}_q^{(t+1)}$ and $v_p^{(t+1)} = v_q^{(t+1)}$.

$$\begin{aligned} \mathbf{w}_p^{(t+1)} &= \mathbf{w}_p^{(t)} - \eta \frac{\partial \widehat{\mathcal{R}}}{\partial \mathbf{w}_p^{(t)}} = \mathbf{w}_p^{(t)} - \eta \left((\partial \ell / \partial f) (\partial \sigma / \partial \langle \mathbf{w}_p, \mathbf{x} \rangle) \left(\sum_{j=1}^m v_j\right) (\mathbf{x}^T) \right) \\ \mathbf{w}_q^{(t+1)} &= \mathbf{w}_q^{(t)} - \eta \frac{\partial \widehat{\mathcal{R}}}{\partial \mathbf{w}_q^{(t)}} = \mathbf{w}_q^{(t)} - \eta \left((\partial \ell / \partial f) (\partial \sigma / \partial \langle \mathbf{w}_q, \mathbf{x} \rangle) \left(\sum_{j=1}^m v_j\right) (\mathbf{x}^T) \right) \end{aligned}$$

The updates $\mathbf{w}_p^{(t+1)}$ and $\mathbf{w}_q^{(t+1)}$ will be equal so long that $\mathbf{w}_p^{(t)} = \mathbf{w}_q^{(t)}$

$$v_p^{(t+1)} = v_p^{(t)} - \eta \frac{\partial \widehat{\mathcal{R}}}{\partial v_p^{(t)}} = v_p^{(t)} - \eta \left((\partial \ell / \partial f) \left(\sum_{j=1}^m \sigma(\langle \mathbf{w}_j, \mathbf{x} \rangle)\right) \right) = v_q^{(t+1)}$$

Similarly, v_p and v_q will follow the same update as long as $v_p^{(t)} = v_q^{(t)}$.

- (c) Suppose there exist two hidden units $p, q \in \{1, 2, \dots, m\}$ such that $\mathbf{w}_p^{(0)} = \mathbf{w}_q^{(0)}$ and $v_p^{(0)} = v_q^{(0)}$. Using induction, conclude 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.

For the base case of $t = 0$, we're given that $\mathbf{w}_p^{(0)} = \mathbf{w}_q^{(0)}$. For any step $t \geq 0$, we retain the update rule derived in 4b such that $\mathbf{w}_p^{(t+1)} = \mathbf{w}_q^{(t+1)}$ and $v_p^{(t+1)} = v_q^{(t+1)}$, so we can say that $\mathbf{w}_p^{(t)} = \mathbf{w}_q^{(t)}$ and $v_p^{(t)} = v_q^{(t)}$, $\forall t \geq 0$.

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

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 0 and variance $1/m$). We will show that

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

For convenience, define $\mathbf{W} := \mathbf{W}^{(0)}$.

- (d) Let \mathbf{w}^\top be an arbitrary row of \mathbf{W} . Prove that

$$\mathbb{E} \left[(\mathbf{w}^\top \mathbf{x})^2 \right] = \mathbb{E} \left[\sum_{i=1}^d w_i^2 x_i^2 + \sum_{\substack{i,j=1 \\ i \neq j}}^d w_i w_j x_i x_j \right].$$

We can prove this knowing that $(\mathbf{w}^\top \mathbf{x})$ is simply the inner product, where $(\mathbf{w}^\top \mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle = \sum_{i=1}^d w_i x_i$. With this, the result becomes trivial as $(\mathbf{w}^\top \mathbf{x})^2 = \left(\sum_{i=1}^d w_i x_i \right)^2 = \sum_{i=1}^d w_i^2 x_i^2 + \sum_{\substack{i,j=1 \\ i \neq j}}^d w_i w_j x_i x_j$.

- (e) Using linearity of expectation, prove that

$$\mathbb{E} \left[\sum_{i=1}^d w_i^2 x_i^2 + \sum_{\substack{i,j=1 \\ i \neq j}}^d w_i w_j x_i x_j \right] = \frac{1}{m} \|\mathbf{x}\|^2.$$

Hint: It may be helpful to recall that for independent random variables X, Y , we have $\mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y]$ and $\text{Var}(X) = \mathbb{E}[X^2] - (\mathbb{E}[X])^2$.

Knowing that $(\mathbf{w}^\top \mathbf{x})^2 = \sum_{i=1}^d w_i^2 x_i^2 + \sum_{\substack{i,j=1 \\ i \neq j}}^d w_i w_j x_i x_j$, we can rewrite as $\mathbb{E} \left[\sum_{i=1}^d w_i^2 x_i^2 + \sum_{\substack{i,j=1 \\ i \neq j}}^d w_i w_j x_i x_j \right] = \mathbb{E} \left[(\mathbf{w}^\top \mathbf{x})^2 \right]$. We can also say that $\mathbb{E}[w_i] = 0$ and thus $\mathbb{E}[w_i^2] = 0$ and the first summation goes to zero. By the remaining summation and given means, we have $\mathbb{E} \left[\sum_{\substack{i,j=1 \\ i \neq j}}^d w_i w_j x_i x_j \right] = \frac{1}{m} \|\mathbf{x}\|^2$.

(f) Using parts (d) and (e), prove that

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

Remark: A similar property holds with the ReLU activation.

With the same reasoning as in 4e, $\mathbb{E} [\|\mathbf{W}\mathbf{x}\|^2] = \|\mathbf{x}\|^2$ is simply equivalent to $\mathbb{E} [\|\mathbf{W}\mathbf{x}\|^2] = \sum_i^d x_i^2$, as all \mathbf{W} components will cancel and the squared norm of \mathbf{x} remains.

Solution.