

## HOMWORK SET 2

CSCI 5525 Advanced Machine Learning (Spring 2021)

**Due** 11:59 pm, Mar 19 2021

**Instruction** Typesetting your homework in  $\text{\LaTeX}$  is optional but encouraged, and you need to submit it as a single PDF file in Canvas. For programming, include all your codes and running results in a single Jupyter notebook file and submit it alongside the main PDF (since Jupyter notebook also allows text editing, feel free to put your textual answers inside the Jupyter notebook sometimes). No late submission will be accepted.

For each problem, you should acknowledge your collaborators if any. For problems containing multiple subproblems, there are often close logic connections between the subproblems. So always remember to build on previous ones, rather than work from scratch.

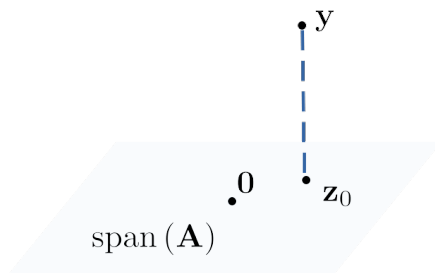
**Reminder about notations** We will use small letters (e.g.,  $u$ ) for scalars, small boldface letters (e.g.,  $\mathbf{a}$ ) for vectors, and capital boldface letters (e.g.,  $\mathbf{A}$ ) for matrices. For a matrix  $\mathbf{A}$ ,  $\mathbf{a}^i$  (superscripting) means its  $i$ -th row as a *row vector*, and  $\mathbf{a}_j$  (subscripting) means the  $j$ -th column as a column vector, and  $a_{ij}$  means its  $(i, j)$ -th element.  $\mathbb{R}$  is the set of real numbers.  $\mathbb{R}^n$  is the space of  $n$ -dimensional real vectors, and similarly  $\mathbb{R}^{m \times n}$  is the space of  $m \times n$  real matrices. The dotted equal sign  $\doteq$  means defining.

**Problem 1 (Least squares problem; 7/12)** Consider the least-squares problem

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) = \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2, \quad \mathbf{A} \in \mathbb{R}^{m \times n}. \quad (1)$$

In class, we showed that any local (also global, as  $f$  is convex) minimizer  $\mathbf{x}_0$  of  $f$  must obey the first-order necessary condition  $\nabla f(\mathbf{x}_0) = 2\mathbf{A}^\top(\mathbf{A}\mathbf{x}_0 - \mathbf{y}) = \mathbf{0}$ , or  $\mathbf{A}^\top\mathbf{A}\mathbf{x}_0 = \mathbf{A}^\top\mathbf{y}$ .

- (a) Assume  $m \geq n$ . Show that  $\mathbf{A}$  has column full rank, i.e., with linearly independent columns, if and only if  $\mathbf{A}^\top\mathbf{A}$  is invertible. (Hint: consider using compact SVD, or the fact that if a square matrix  $\mathbf{M}$  is invertible if and only if the only solution of  $\mathbf{M}\mathbf{z} = \mathbf{0}$  is  $\mathbf{z} = \mathbf{0}$ , i.e., if and only if  $\mathbf{M}$  has a trivial null space. (1/12)
- (b) Assume  $m \geq n$  and  $\mathbf{A}$  has column full rank. Then columns of  $\mathbf{A}$  span an  $n$ -dimensional subspace of  $\mathbb{R}^m$ .



Since any point on the subspace can be written in the form of  $\mathbf{A}\mathbf{x}$  for a certain  $\mathbf{x}$ , geometrically the optimization problem in Eq. (1) is to compute the distance of a given point  $\mathbf{y} \in \mathbb{R}^m$  to the subspace  $\text{span}(\mathbf{A})$ , i.e.,

$$\min_{\mathbf{z} \in \text{span}(\mathbf{A})} g(\mathbf{z}) = \|\mathbf{y} - \mathbf{z}\|_2^2. \quad (2)$$

Let  $\mathbf{z}_0$  be a minimizer for Eq. (2). Can we provide an analytic form for  $\mathbf{z}_0$ ? Is it unique and why or why not? In  $\mathbb{R}^m$ , two vectors  $\mathbf{u}, \mathbf{v}$  are orthogonal to each other if  $\langle \mathbf{u}, \mathbf{v} \rangle = 0$ . Prove that  $\langle \mathbf{z}_0 - \mathbf{y}, \mathbf{w} \rangle$  for any  $\mathbf{w} \in \text{span}(\mathbf{A})$ , i.e., the vector  $\mathbf{z}_0 - \mathbf{y}$  is orthogonal to the space  $\text{span}(\mathbf{A})$ —here,  $\mathbf{z}_0$  is called the orthogonal projection of  $\mathbf{y}$  onto the subspace. (1/12)

- (c) Fix a random seed, and set  $m = 100, n = 50$ . Generate  $\mathbf{A}$  and  $\mathbf{y}$  as iid standard normal. Implement gradient descent, with backtracking line search as the step size rule. Please submit your code, and plot the objective value vs. iterate. (1/12)
- (d) When  $m < n$ , there is no unique global minimizer for  $f$ . Assume  $\mathbf{x}_0$  is a global minimizer. Provide write down the set of all global minimizers and also prove the set indeed contains all global minimizers. (Hint: when constructing the set, start from the optimality condition and derive what condition the difference of two global minimizers must satisfy. ) (1/12)
- (e) Assume  $m < n$  and we run gradient descent to optimize  $f(\mathbf{x})$  starting with  $\mathbf{x}^{(0)} = \mathbf{0}$  and with a fixed step size  $t$ .
- (i) Write down the analytic forms of  $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}$ , and  $\mathbf{x}^{(3)}$  and prove that  $\mathbf{x}^{(k)} \in \text{row}(\mathbf{A})$  for all  $k \geq 0$ , i.e., all iterates will stay in  $\text{row}(\mathbf{A})$ . (1/12)
  - (ii) When  $t$  is sufficiently small, or backtracking line search is implemented, the sequence  $\mathbf{x}^{(0)}, \mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots$  will converge to a local (and also global) minimizer of  $f(\mathbf{w})$ . Obviously, the limit  $\mathbf{x}_*$  also lies in  $\text{row}(\mathbf{A})$ . Prove that  $\mathbf{x}_*$  is global minimizer of  $f(\mathbf{x})$  with the smallest  $\ell_2$  norm among all global minimizers. (Hint: any  $\mathbf{x} \in \mathbb{R}^n$  can be decomposed as two orthogonal components:  $\mathbf{x} = \mathbf{x}_r + \mathbf{x}_n$ , where  $\mathbf{x}_r \in \text{row}(\mathbf{A})$  and  $\mathbf{x}_n \in \text{null}(\mathbf{A})$ . ) (1/12)

Why is this interesting? Although in our setting  $f(\mathbf{x})$  does not have a unique minimizer, a specific algorithm—gradient descent, with specific initialization—zero (in fact, the above results hold for any initialization lying the  $\text{row}(\mathbf{A})$ ) and appropriate step sizes, returns a solution with a special property—least  $\ell_2$  norm, which is simple in some sense. This phenomenon of algorithmic biases toward certain “simple” solutions on problems with many solutions is called *implicit regularization*. This is the basis for one of many emerging theories explaining why deep neural networks generalize well despite that the network capacities well exceed the intrinsic data complexities.

- (f) Assume  $m < n$  again and so global minimizers for  $f$  abound. In statistics, this kind of problems where the number of data points is (far) less than the input (feature) dimension is called *high-dimensional problems* (other fields may use the term in different ways). Reasonable regularization is always needed for high-dimensional problems to bias the estimation procedure toward solutions with certain desired properties. For our least-squares problem Eq. (1), one possibility is seeking (approximately) sparse solutions by adding in  $\ell_1$  norm regularization:

$$\min_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{m} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1. \quad (3)$$

This is the famous *Lasso* problem. Suppose features (i.e., columns of  $\mathbf{A}$ ) are gene expressions of patients and the output is the probability of getting liver cancer. If the solution  $\mathbf{x}_*$  we find is indeed sparse, i.e., containing very few large entries and the rest entries negligibly small. Then only the genes/columns corresponding to the large entries of  $\mathbf{x}_*$  are important for predicting the cancer probability. So Lasso and its variants are often used for feature selection.

Fix a random seed, and set  $m = 30, n = 100$ . Generate  $\mathbf{A}$  as iid standard normal, a groundtruth  $\mathbf{x}_0$  as iid Bernoulli with rate 0.2 (i.e., assuming 1 with probability 0.2, and otherwise assuming 0), and so  $\mathbf{y} = \mathbf{A}\mathbf{x}_0 + \boldsymbol{\varepsilon}$ , where  $\boldsymbol{\varepsilon}$  is iid normal with mean 0 and variance 0.5, i.e.,  $\mathcal{N}(0, 0.5)$ . Implement Lasso with the CVXPY package, with reference to this example [https://www.cvxpy.org/examples/machine\\_learning/lasso\\_regression.html](https://www.cvxpy.org/examples/machine_learning/lasso_regression.html). Try 3 different values—ideally

with different orders of magnitude—for  $\lambda$ , and plot the groundtruth  $x_*$ , and the  $x$ 's you obtain *in the same stem plot*; check out this link [https://matplotlib.org/stable/api/\\_as\\_gen/matplotlib.pyplot.stem.html](https://matplotlib.org/stable/api/_as_gen/matplotlib.pyplot.stem.html) for making stem plots in Python. What do you observe with increasing  $\lambda$ ? (1/12)

## Problem 2 (Logistic regression; 5/12)

- (a) Recall the maximum likelihood estimation (MLE) principle: given observation data, maximize the probability density (or mass, depending on if the distribution is continuous or discrete) function  $p$  (observed data)—it is the *likelihood function* of unknown parameters  $\theta$ ; we often make the dependency on  $\theta$  explicit by writing  $p$  (observed data;  $\theta$ ).

For a supervised learning problem, let's assume the inputs  $x_1, \dots, x_N$  are sampled iid from a fixed but unknown distribution  $\mathcal{D}_x$ . Then, each  $y_i$  for  $i \in [N]$  is sampled from the conditional distribution  $\mathcal{D}_y(x_i, w)$ , which depends on both  $x_i$  and a set of parameters  $w$ . Show that

$$\max_w p\left(\{(x_i, y_i)\}_{i=1}^N; w\right) \quad (4)$$

is equivalent to

$$\max_w \prod_{i=1}^N p(y_i | x_i; w). \quad (5)$$

This is sometimes called *maximum conditional likelihood* (MCL) principle. (Hint: note that  $p(x, y) = p(x)p(y | x)$ . 1/12)

- (b) Now let's apply MLE to derive logistic regression. Since the logistic function  $\phi(z) = \frac{1}{1+e^{-z}}$  takes value in  $[0, 1]$ , we can use it to naturally model probability. Assume our  $\{x_i\}_{i=1}^N$  are sampled from  $\mathcal{D}_x$ , and  $y_i$  takes value in  $\{1, -1\}$  for all  $i \in [N]$  and obeys

$$p(y = 1 | x) = \phi \circ \langle w, x \rangle = \frac{1}{1 + \exp(-\langle w, x \rangle)}. \quad (6)$$

Obviously,  $p(y = -1 | x) = 1 - p(y = 1 | x)$ . Show that applying MCL above yields the logistic regression problem we introduced in class:

$$\min_w \sum_{i=1}^N \log(1 + \exp(-y_i \langle w, x_i \rangle)). \quad (7)$$

(Hint: think of how to derive a unified form of  $p(y | x)$  for convenient analysis; 1/12)

- (c) Assume we solve problem (7) and get  $w_0$ . We shall use this for classification on test data. For any data point  $x$ , a natural classification scheme is:

$$y = \begin{cases} 1 & p(y = 1 | x) > 0.5 \\ -1 & \text{otherwise} \end{cases}. \quad (8)$$

Show that this classifier is linear, i.e., taking the form  $x \mapsto \text{sign}(\langle w_0, x \rangle)$ . (1/12)

- (d) Derive the gradient of problem (7), and implement gradient with backtracking linear search step size—this will be used in the next subproblem. (1/12)

- (e) Load the breast cancer classification dataset in scikit-learn: [https://scikit-learn.org/stable/modules/generated/sklearn.datasets.load\\_breast\\_cancer.html#sklearn.datasets.load\\_breast\\_cancer](https://scikit-learn.org/stable/modules/generated/sklearn.datasets.load_breast_cancer.html#sklearn.datasets.load_breast_cancer). Test and report the classification accuracy of your logistic regression implementation, and compare it with that by the built-in logistic regression module [https://scikit-learn.org/stable/modules/generated/sklearn.linear\\_model.LogisticRegression.html](https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LogisticRegression.html). Don't panic if you see slightly different results, as scikit-learn implements a regularized version; see [https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression). (1/12)