

HOMWORK SET 1

CSCI 5525 Advanced Machine Learning (Spring 2021)

Due 11:59 pm, Feb 10 2021

Instruction Typesetting your homework in \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 (supscripting) 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 (Matrix norms, inner products, traces; 6/12) Recall that for any vector $\mathbf{v} \in \mathbb{R}^n$, the ℓ_p norm of \mathbf{v} is defined as $\|\mathbf{v}\|_p \doteq (\sum_i |v_i|^p)^{1/p}$. The cases when $p = 1, 2, \infty$ are often used. When $p = 2$, it is also called the Euclidean norm. Similar norms can be defined for matrices. Particularly, the direct generalization of the vector Euclidean norm is the *Frobenius norm* defined as

$$\|\mathbf{M}\|_F \doteq \sqrt{\sum_{ij} m_{ij}^2}$$

for a matrix \mathbf{M} . On the other hand, the inner product of matrices is defined similarly to that of vectors. For \mathbf{A}, \mathbf{B} of the same size, $\langle \mathbf{A}, \mathbf{B} \rangle \doteq \sum_{ij} a_{ij} b_{ij}$. Obviously, $\langle \mathbf{A}, \mathbf{B} \rangle = \langle \mathbf{B}, \mathbf{A} \rangle$ and $\|\mathbf{M}\|_F = \sqrt{\langle \mathbf{M}, \mathbf{M} \rangle}$. A third notion of interest is the matrix trace, $\text{tr}(\mathbf{M}) = \sum_i m_{ii}$, i.e., sum of the diagonal entries, which is only defined for square matrices.

- Show that $\langle \mathbf{A}, \mathbf{B} \rangle = \text{tr}(\mathbf{A}^\top \mathbf{B})$ and so $\|\mathbf{M}\|_F = \sqrt{\text{tr}(\mathbf{M}^\top \mathbf{M})}$. (1/12)
- Show that $\text{tr}(\mathbf{A}^\top \mathbf{B}) = \text{tr}(\mathbf{B}^\top \mathbf{A})$. (1/12)
- Assume \mathbf{A} and \mathbf{B} have the same size. In general, $\mathbf{A}\mathbf{B}^\top$ and $\mathbf{B}^\top \mathbf{A}$ have different sizes, but $\text{tr}(\mathbf{A}\mathbf{B}^\top) = \text{tr}(\mathbf{B}^\top \mathbf{A})$. Show it! (1/12)
- Show that $\text{tr}(\mathbf{M}_1 \mathbf{M}_2 \mathbf{M}_3) = \text{tr}(\mathbf{M}_3 \mathbf{M}_1 \mathbf{M}_2) = \text{tr}(\mathbf{M}_2 \mathbf{M}_3 \mathbf{M}_1)$, assuming that the sizes of \mathbf{M}_1 , \mathbf{M}_2 and \mathbf{M}_3 are compatible with all the matrix multiplications. This is known as the *cyclic property* of matrix traces. (Hint: think of (c)) (1/12)
- For any matrices $\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}$ of compatible sizes, we always have $\langle \mathbf{A}\mathbf{C}\mathbf{B}, \mathbf{D} \rangle = \langle \mathbf{C}\mathbf{B}, \mathbf{A}^\top \mathbf{D} \rangle = \langle \mathbf{A}\mathbf{C}, \mathbf{D}\mathbf{B}^\top \rangle$, i.e., we can always move the **leading** matrix of one side of the inner product to the other side as **leading matrix once transposed** (if these matrices are complex-valued, should be conjugate transposed), and similarly the **trailing** matrix to the other side as **trailing matrix once transposed**. Why? (Hint: think of the above results and also try to remember this important property that will be useful for calculation later) (1/12)

- (f) For M , let's perform a *compact SVD* (if not sure, check up Wikipedia! https://en.wikipedia.org/wiki/Singular_value_decomposition#Compact_SVD) to obtain $M = U\Sigma V^\top$, so that U and V are orthonormal (not necessarily square) matrices, i.e., $U^\top U = I$ and $V^\top V = I$. Use the cyclic property of trace and that $\|M\|_F = \sqrt{\text{tr}(M^\top M)}$ to show that

$$\|M\|_F = \sqrt{\sum_{i=1}^r \sigma_i^2},$$

assuming the rank of M is r . Here σ_i 's are the singular values of M . (1/12)

Problem 2 (Computation of Jacobian, gradient, and Hessian; 4/12) For all the subproblems, you are free to deploy any techniques or their mixtures that we describe in the class. However, the perturbation-expansion (Taylor-expansion) technique could be more efficient for some cases.

Background on convexity: A twice-differentiable function $f(x)$ is convex if $\nabla^2 f(x)$ is positive semidefinite for all x . If $\nabla^2 f(x)$ is positive definite for all x , then f is said to be strongly convex and f has a unique minimizer.

- (a) Let A be a square matrix. Deriving the gradient and Hessian of the quadratic function $f(x) = x^\top A x + b^\top x$. Please include your calculation details. (Hint: note that Hessian must be a symmetric matrix.) (1/12)
- (b) Let $p(x; \beta) = \frac{e^{\beta^\top x}}{1 + e^{\beta^\top x}}$. The log-likelihood for logistic regression with two classes is (assuming N samples of the form (x_i, y_i))

$$\begin{aligned} f(\beta) &= \sum_{i=1}^N [y_i \log p(x_i; \beta) + (1 - y_i) \log (1 - p(x_i; \beta))] \\ &= \sum_{i=1}^N [y_i \beta^\top x_i - \log (1 + e^{\beta^\top x_i})]. \end{aligned}$$

Derive the gradient and Hessian of $f(\beta)$. Please include your calculation details. (1/12) For logistic regression, we are going to maximize $f(\beta)$, which is equivalent to minimize $-f(\beta)$. Does the minimization problem has a unique minimizer or not? (0.5/12)

- (c) Let $A \in \mathbb{R}^{m \times n}$ with $m \leq n$, then given a $y \in \mathbb{R}^m$, the least-squares problem $\min_x \|y - Ax\|_2^2$ has infinitely many solutions. Now let's say we want a solution with a small ℓ_2 norm, then it is reasonable to put a penalty on the ℓ_2 norm:

$$\min_x \|y - Ax\|_2^2 + \lambda \|x\|_2^2$$

with a chosen $\lambda > 0$. This is *ridge regression*. Now we know that for an unconstrained first-order differentiable function $g(x)$, any of its local minimizer x_* must satisfy the *first-order optimality condition*: $\nabla g(x_*) = 0$. Use this to derive x_* (1/12). Is the x_* unique? Why? (0.5/12)

Problem 3 (Robust linear regression; 2/12) Given data points $\{(x_i, y_i)\}_{i=1}^N$. Linear regression tries to find a linear function of x , i.e., $w^\top x$, so that the collective approximation error $\sum_{i=1}^N \ell(w^\top x_i, y_i)$ is minimized. Geometrically, this fits a linear subspace to the point cloud $\{(x_i, y_i)\}_{i=1}^N$ so that the cumulative error in the y direction can be minimized, as illustrated in Fig. 1.

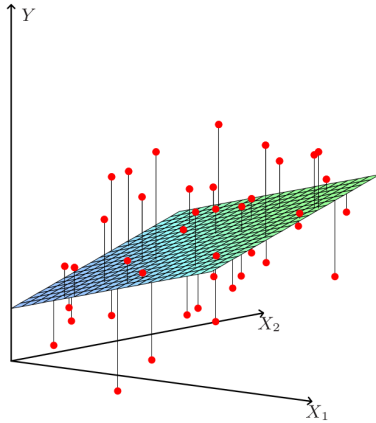


Figure 1: Illustration of linear regression. Figure reproduced from Figure 3.1 of [HTF09].

If the data points $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$ do follow a linear relationship with deviations of comparable magnitude, e.g., for an underlying \mathbf{w}_*

$$y_i = \mathbf{w}_*^\top \mathbf{x}_i + \varepsilon_i \quad \forall i = 1, \dots, N$$

with ε_i iid Gaussian with a small variance, the typical least-squares objective, e.g.,

$$\min_{\mathbf{w}} \sum_{i=1}^N (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 \quad (1)$$

is reasonable. But if ε_i 's are very different in magnitude across the i 's, e.g., coming from heavy-tailed distributions, or due to irregular measurement corruption, the $(\cdot)^2$ as loss function may not be appropriate. This is because one extremely large ε_i may rule the total loss and ruin the estimation of \mathbf{w}_* .

To cure this, an alternative is to use the absolute value, instead of the squares:

$$\min_{\mathbf{w}} \sum_{i=1}^N |\mathbf{w}^\top \mathbf{x}_i - y_i|. \quad (2)$$

Compared to Eq. (1), the influence of terms with potential large errors is suppressed and with small errors amplified, so the estimation procedure tends to be more stable. This is often called *least absolute deviations* (LAD) estimation. For more information on this, check out here https://en.wikipedia.org/wiki/Least_absolute_deviations.

Now let's explore the benefit of LAD. *For the sake of reproducibility, please fix a random seed before you start to generate any data.*

- (a) Let's first generate an iid normal (i.e., $\mathcal{N}(0, 1)$) vector $\mathbf{w}_* \in \mathbb{R}^{20}$ and 100 iid normal vectors $\mathbf{x}_i \in \mathbb{R}^{20}$. Now produce

$$y_i = \mathbf{w}_*^\top \mathbf{x}_i + \varepsilon_i, \quad (3)$$

where ε_i 's are iid Laplace $(0, 0.5)$. (0.5/12)

- (b) CVXPY is an excellent Python-based modeling framework for solving small-scale convex optimization problems. Follow the instruction here <https://www.cvxpy.org/install/index.html> and install it in your workspace. (0.5/12)
- (c) Read this CVXPY example on solving least squares https://www.cvxpy.org/examples/basic/least_squares.html and adapt it for solving our least-squares with the data of (a). To solve LAD, you only need to change the line

```
cost = cp.sum_squares(A @ x - b)
```

to

```
cost = cp.norm1(A @ x - b).
```

For the $\hat{\mathbf{w}}$ estimated from both least-squares and LAD, compute $\|\hat{\mathbf{w}} - \mathbf{w}_*\|_2$, i.e., estimate error for \mathbf{w}_* . Which method leads to smaller estimation error? (1/12)

References

- [HTF09] Trevor Hastie, Robert Tibshirani, and Jerome Friedman, *The elements of statistical learning: Data mining, inference, and prediction, second edition*, SPRINGER NATURE, 2009.