

# HOMWORK SET 0

CSCI 5980/8980 Think Deep Learning (Fall 2020)

**Due** 11:59 pm, Sep 30 2020

**Instruction** Please typeset your homework in  $\text{\LaTeX}$  and submit it as a single PDF file in Canvas. 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.

**Notation** 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.  $\mathbb{R}$  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 (Chain rules, gradient and Hessian)** Recall from calculus that for a multivariate function  $f(\mathbf{x})$  mapping from  $\mathbb{R}^n$  to  $\mathbb{R}$ , i.e.,  $f: \mathbb{R}^n \mapsto \mathbb{R}$ , the  $i$ -th partial derivative of  $f$  is defined as  $\frac{\partial f}{\partial x_i}$ , i.e., the univariate derivative with respect to the  $i$ -th variable while holding the other variables constant. This generalizes naturally to the matrix case, where we consider  $f(\mathbf{X})$  with  $\mathbf{X} \in \mathbb{R}^{m \times n}$ . The  $(i, j)$ -th partial derivative of  $f(\mathbf{X})$  is then  $\frac{\partial f}{\partial X_{ij}}$ . For either case, the *gradient* is the ordered collection of the partial derivatives. A common convention is to order the partial derivatives into the same shape as the input variable. So for  $f(\mathbf{x})$ , the gradient  $\nabla f(\mathbf{x}) \in \mathbb{R}^n$  is

$$\nabla f(\mathbf{x}) \doteq \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_i} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix}.$$

Similarly for  $f(\mathbf{X})$ , the gradient  $\nabla f(\mathbf{X}) \in \mathbb{R}^{m \times n}$  is

$$\nabla f(\mathbf{X}) \doteq \begin{bmatrix} \frac{\partial f}{\partial X_{11}} & \cdots & \frac{\partial f}{\partial X_{1j}} & \cdots & \frac{\partial f}{\partial X_{1n}} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \frac{\partial f}{\partial X_{i1}} & \cdots & \frac{\partial f}{\partial X_{ij}} & \cdots & \frac{\partial f}{\partial X_{in}} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \frac{\partial f}{\partial X_{m1}} & \cdots & \frac{\partial f}{\partial X_{mj}} & \cdots & \frac{\partial f}{\partial X_{mn}} \end{bmatrix}.$$

Each of the  $\frac{\partial f}{\partial x_i}$  (respectively  $\frac{\partial f}{\partial X_{ij}}$ ) is a real-valued function of  $\mathbf{x}$  (respectively  $\mathbf{X}$ ) again. So we can take the partial derivatives with respect to entries in  $\mathbf{x}$  (respectively  $\mathbf{X}$ ) again to obtain the 2nd order partial derivatives like

$$\frac{\partial \left( \frac{\partial f}{\partial x_i} \right)}{\partial x_\ell},$$

which we often write neatly as

$$\frac{\partial^2 f}{\partial x_\ell \partial x_i}.$$

For  $f(\mathbf{x})$ , there are two indices  $(\ell, j)$  to index all the 2nd order partial derivatives, and so it is convenient to form an  $n \times n$  matrix whose  $(\ell, j)$ -th entry is  $\frac{\partial^2 f}{\partial x_\ell \partial x_j}$  to collect them. We shall write this matrix compactly as  $\left[ \frac{\partial^2 f}{\partial x_\ell \partial x_j} \right]_{\ell, j} \in \mathbb{R}^{n \times n}$ , and this is the *Hessian* for  $f(\mathbf{x})$ :

$$\nabla^2 f(\mathbf{x}) \doteq \left[ \frac{\partial^2 f}{\partial x_\ell \partial x_j} \right]_{\ell, j}.$$

Under very mild conditions, Hessian for  $f(\mathbf{x})$  is symmetric, that is,  $\frac{\partial^2 f}{\partial x_\ell \partial x_j} = \frac{\partial^2 f}{\partial x_j \partial x_\ell}$  for all  $i$  and  $\ell$ . The Hessian for  $f(\mathbf{X})$  is more complicated. There are four indices to index all the 2nd order partial derivatives, e.g.,  $\frac{\partial^2 f}{\partial X_{\ell, k} \partial X_{i, j}}$ . So it is natural to use a 4-dimensional array (matrices are 2-dimensional arrays) to collect them here, i.e., we need 4-D tensors. We will cover more about tensors when we come to deep learning on images. In practice, one would never explicitly form the Hessian of  $f(\mathbf{X})$  for computation. We will talk more about this in future lectures and homework sets on numerical optimization.

Let  $y = f(u)$  be differentiable with respect to  $u$  and  $u = g(x)$  be differentiable with respect to  $x$ . The univariate chain rule says

$$\frac{dy}{dx} = \frac{dy}{du} \frac{du}{dx}.$$

This can be extended to the multivariate case. Let  $y = f(u_1, \dots, u_j, \dots, u_m)$  be differentiable and  $u_1, \dots, u_j, \dots, u_m$  be differentiable functions of  $(x_1, \dots, x_i, \dots, x_n)$ . Then, the multivariate chain rule says that

$$\frac{dy}{dx_i} = \sum_j \frac{\partial y}{\partial u_j} \frac{\partial u_j}{\partial x_i} \quad \forall i.$$

We can also write the multivariate version in matrix notations. Write  $y = f(\mathbf{u})$  where  $\mathbf{u} \in \mathbb{R}^m$ , and  $\mathbf{u}$  as a differentiable function of  $\mathbf{x} \in \mathbb{R}^n$ . Then,

$$\nabla_{\mathbf{x}} f = [\mathbf{J}_{\mathbf{u}}(\mathbf{x})]^T \nabla_{\mathbf{u}} f,$$

where  $\mathbf{J}_{\mathbf{u}}(\mathbf{x}) \in \mathbb{R}^{m \times n}$  is the Jacobian matrix of  $\mathbf{u}$  with respect to  $\mathbf{x}$ :

$$\mathbf{J}_{\mathbf{u}}(\mathbf{x}) = \left[ \frac{\partial u_j}{\partial x_i} \right]_{j, i}.$$

- (a) Let  $p(\mathbf{x}; \boldsymbol{\beta}) = \frac{e^{\boldsymbol{\beta}^T \mathbf{x}}}{1 + e^{\boldsymbol{\beta}^T \mathbf{x}}}$ . The log-likelihood for logistic regression with two classes is (assuming  $N$  samples of the form  $(\mathbf{x}_i, y_i)$ )

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

Derive the gradient and Hessian of  $f(\boldsymbol{\beta})$ . Please include your calculation details. (1/12)

- (b) Let  $\sigma = \tanh(x)$ , i.e., the *hyperbolic tangent function*. Derive the gradient of the matrix-variable function  $g(\mathbf{W}) = \|\mathbf{y} - \sigma(\mathbf{W}\mathbf{x})\|_2^2$ , where  $\sigma$  is applied to the vector  $\mathbf{W}\mathbf{x}$  elementwise. This is regression based on a one-layer network. Please include your calculation details. (1/12)

**Problem 2 (Matrix norms, inner products, traces)** Recall that for any vector  $\mathbf{v} \in \mathbb{R}^n$ , the  $\ell_p$  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

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

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

- Show that  $\langle A, B \rangle = \text{tr}(A^T B)$  and so  $\|M\|_F = \sqrt{\text{tr}(M^T M)}$ . (1/12)
- Show that  $\text{tr}(A^T B) = \text{tr}(B^T A)$ . (1/12)
- Assume  $A$  and  $B$  have the same size. In general,  $AB^T$  and  $B^T A$  have different sizes, but  $\text{tr}(AB^T) = \text{tr}(B^T A)$ . Show it! (1/12)
- Show that  $\text{tr}(M_1 M_2 M_3) = \text{tr}(M_3 M_1 M_2) = \text{tr}(M_2 M_3 M_1)$ , assuming that the sizes of  $M_1, M_2$  and  $M_3$  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  $A, B, C, D$  of compatible sizes, we always have  $\langle ACB, D \rangle = \langle CB, A^T D \rangle = \langle AC, DB^T \rangle$ , 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)
- For  $M$ , let's perform a *compact SVD* (if not sure, check up Wikipedia! [https://en.wikipedia.org/wiki/Singular\\_value\\_decomposition#Compact\\_SVD](https://en.wikipedia.org/wiki/Singular_value_decomposition#Compact_SVD)) to obtain  $M = U \Sigma V^T$ , so that  $U$  and  $V$  are orthonormal (not necessarily square) matrices, i.e.,  $U^T U = I$  and  $V^T V = I$ . Use the cyclic property of trace and that  $\|M\|_F = \sqrt{\text{tr}(M^T M)}$  to show that

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

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

**Problem 3 (Taylor expansion of multivariate functions)** We learned Taylor's theorem for single variable functions in our first calculus course. For a function  $f$  that is twice-differentiable around  $x$ ,

$$f(x + \delta) = f(x) + f'(x) \delta + \frac{1}{2} f''(x) \delta^2 + o(|\delta|^2) \quad \text{as } \delta \rightarrow 0,$$

i.e., the approximation error of the 2nd order Taylor expansion  $f(x) + f'(x) \delta + \frac{1}{2} f''(x) \delta^2$  is lower order than  $|\delta|^2$  when  $\delta$  is sufficiently small. The result generalizes naturally to the multivariate case. For  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  that is twice-differentiable around  $\mathbf{x}$ ,

$$f(\mathbf{x} + \boldsymbol{\delta}) = f(\mathbf{x}) + [\nabla f(\mathbf{x})]^T \boldsymbol{\delta} + \frac{1}{2} \boldsymbol{\delta}^T \nabla^2 f(\mathbf{x}) \boldsymbol{\delta} + o(\|\boldsymbol{\delta}\|_2^2).$$

Similarly, for  $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$  that is twice-differentiable around  $\mathbf{X}$ ,

$$f(\mathbf{X} + \Delta) = f(\mathbf{X}) + \text{tr}([\nabla f(\mathbf{X})]^\top \Delta) + \frac{1}{2} \text{tr}(\Delta^\top \nabla^2 f(\mathbf{X}) \Delta) + o(\|\Delta\|_F^2).$$

Consider the one-layer network regression again:  $g(\mathbf{W}) = \|\mathbf{y} - \sigma(\mathbf{W}\mathbf{x})\|_2^2$  with  $\sigma(x) = \tanh(x)$ . Let's try to work out its 1st order Taylor expansion by direct expansion as follows.

- Show that  $\sigma((\mathbf{W} + \Delta)\mathbf{x}) = \sigma(\mathbf{W}\mathbf{x}) + \sigma'(\mathbf{W}\mathbf{x}) \odot (\Delta\mathbf{x}) + o(\|\Delta\|_F)$  when  $\Delta \rightarrow \mathbf{0}$ . Here, both  $\sigma$  and  $\sigma'$  are applied elementwise, and  $\odot$  denotes the elementwise (Hadamard) product. (1/12)
- So  $\mathbf{y} - \sigma((\mathbf{W} + \Delta)\mathbf{x}) = (\mathbf{y} - \sigma(\mathbf{W}\mathbf{x})) - \sigma'(\mathbf{W}\mathbf{x}) \odot (\Delta\mathbf{x}) - o(\|\Delta\|_F)$  when  $\Delta \rightarrow \mathbf{0}$ . Substitute this back into the square and use the identity  $\|\mathbf{a} + \mathbf{b} + \mathbf{c}\|_2^2 = \|\mathbf{a}\|_2^2 + \|\mathbf{b}\|_2^2 + \|\mathbf{c}\|_2^2 + 2\mathbf{a}^\top \mathbf{b} + 2\mathbf{a}^\top \mathbf{c} + 2\mathbf{b}^\top \mathbf{c}$  to obtain the first-order approximation to  $g(\mathbf{W} + \Delta)$ . Remember that any terms lower order than  $\|\Delta\|_F$  are not interesting and we can always assume  $\Delta$  as small as needed. (1/12)
- Substitute the result from **Problem 1(b)** into the 1st order Taylor expansion formula above and compare it to the result obtained here. Are they equal or not? (0.5/12)

**Problem 4 (Conditional probability and Bayes' Rule)** Let  $P(A | B)$  be the probability of event  $A$  given event  $B$ . Bayes' theorem states that

$$P(A | B) = \frac{P(B | A) P(A)}{P(B)}.$$

In general, for a partition  $\{A_i\}$  of the sample space,

$$P(A_i | B) = \frac{P(B | A_i) P(A_i)}{\sum_i P(B | A_i) P(A_i)}.$$

- Suppose that women who drink alcohol often are 10 times more likely to develop a stomach cancer than women who do not, while the corresponding risk for men is 20 (drinker v.s. non-drinker). Statistical data show that 40% of women drink alcohol often. Now, a woman has been diagnosed with the stomach cancer. But due to some privacy issues, it is impossible for you to get other information about this woman. Based on the limited information, is it possible to calculate the probability that she is a drinker? If yes, what is the probability? If no, explain why and what extra information you will need.<sup>1</sup> (1/12)
- Doctors apply a standardized test for a certain disease. If the patient has the disease, the test shows a positive result with a 99% chance. However, with 2% probability a healthy patient can have a positive test. Statistical data show that 1 out of 1000 in the population have the disease. What is the probability for a patient with a positive test to be affected by the disease? (0.5/12)

---

<sup>1</sup>All data and information in this question are made up. No medical studies/institutions have proven it.