

## HANDOUT 1

We review probability theory and statistics in this handout. We will also have a basic introduction to Stochastic Gradient Descent, an extremely important algorithm for machine learning.

### 1 Basic Probability

#### PROBLEM 1: THE MOST USEFUL BOUND IN PROBABILITY THEORY

Recall the definition of a probability measure. In this problem, we will rigorously prove the most useful bound in probability theory, the Union Bound, stating that given any  $n$  events  $E_1, E_2, \dots, E_n$ , we have

$$P(\cup_{i=1}^n E_i) \leq \sum_{i=1}^n P(E_i).$$

Prove the following statements.

(a) For any two events  $A$  and  $B$  such that  $A \subseteq B$ , prove that  $P(A) \leq P(B)$ .

(b) Prove the union bound  $P(\cup_{i=1}^n E_i) \leq \sum_{i=1}^n P(E_i)$ .

(Hint: For any two events  $E_1$  and  $E_2$ , prove that  $P(E_1 \cup E_2) \leq P(E_1) + P(E_2)$ . For this, you can use a useful decomposition rule of sets:  $E_1 \cup E_2 = (E_1 \cap E_2^c) \cup E_2$ . Applying this bound recursively will finish the proof.

To use the recursive arguments, you need to write the union of  $n$  events as a union of 2 events.)

### 2 “Modeling” in the Absence of Models; Neural Networks as Universal Approximators

Recall that the least-square estimator naturally arises from the maximum likelihood estimation of

A1. a true signal  $\mathbf{x}^\dagger$ ,

A2. with a linear model,

A3. of added Gaussian noise.

That is, if we assume that the data  $(\mathbf{a}_i, b_i) \in (\mathbb{R}^p \times \mathbb{R})$  is generated by the following relation:

$$b_i = \langle \mathbf{a}_i, \mathbf{x}^\dagger \rangle + \mathbf{W}_i, \quad \mathbf{W}_i \sim \mathcal{N}(0, \sigma^2),$$

then

$$\mathbf{x}_{ML}^* = \arg \min_{\mathbf{x} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2^2$$

where  $\mathbf{A}^\top = [\mathbf{a}_1 \ \mathbf{a}_2 \ \dots \ \mathbf{a}_n]$ .

Although the importance of LS estimators can hardly be overemphasized, it is still widely criticized by its restrictive setting; i.e., assuming A1-3. A natural question is: What if we drop these assumptions? For instance, what if the model is nonlinear, or there is no such thing as a true signal, or the noise is far away from Gaussian?

Neural networks [2] present a powerful framework to address these concerns with the following philosophy: First, denote the (possibly very complicated) relation between input and output by  $b = h(\mathbf{a})$ , where  $h$  is a function from  $\mathbb{R}^p \rightarrow \mathbb{R}$  and  $(\mathbf{a}, b)$  follows some **unknown** distribution  $\mathbb{P}$ . Then the function  $h$  should satisfy

$$h = \arg \min_g \frac{1}{2} \mathbb{E}_{\mathbb{P}} (g(\mathbf{a}) - b)^2, \quad g \text{ any function.} \quad (1)$$

**(Remark:** We focus on the regression example here, while the very same idea works for classification just as well.)

Now, there are two reasons why solving (1) is not possible:

1. We do not know the distribution  $\mathbb{P}$ ; instead, we only have samples  $(\mathbf{a}_i, b_i)$  from  $\mathbb{P}$ .
2. We do not know how to optimize over the set of all functions.

For the first issue, we already know that we can replace the true average by the empirical average, leading to an empirical risk minimization problem. For the second, the key idea of deep learning relies on the following theorem [3]: Informally speaking,

*Any function  $f$  can be approximated arbitrarily well by a neural network, as long as the network size is big enough.*

Combining these, we are lead to the empirical risk minimization of neural networks:

$$\mathbf{x}^* = \arg \min_{\mathbf{x}} \left\{ \frac{1}{n} \sum_{i=1}^n (b_i - h_{\mathbf{x}}(\mathbf{a}_i))^2 \right\}, \quad h_{\mathbf{x}}(\mathbf{a}) = \sigma(\mathbf{W}_k \sigma(\cdots \sigma(\mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{a} + \boldsymbol{\mu}_1) + \boldsymbol{\mu}_2) \cdots) + \boldsymbol{\mu}_k); \quad (2)$$

where  $\sigma$  is a "proper" activation function that requires some condition on continuity,  $\mathbf{x} = (\mathbf{W}_1, \boldsymbol{\mu}_1, \mathbf{W}_2, \boldsymbol{\mu}_2, \dots, \mathbf{W}_k, \boldsymbol{\mu}_k)$ ,  $\mathbf{W}_i \in \mathbb{R}^{d_i \times d_{i-1}}$ ,  $\boldsymbol{\mu}_i \in \mathbb{R}^{d_i}$  and  $\mathbf{a} \in \mathbb{R}^p$  (more on the notation in the lectures). The hope is that, as long as we have enough parameters and data, the learned neural network is a good approximator for the function  $h$ :

$$h_{\mathbf{x}^*} \simeq h,$$

which is true by the theorem of [3]. There are however trade-offs involved in scaling up over-parametrized networks: width helps robustness; depth helps robustness under a certain initialization but hurts it under another [5].

In conclusion,

*Neural networks can be viewed as a "universal modeling" scheme where no assumption about the data distribution is made.*

Finally, all the above reasoning relies on the big "if" that we can optimize (2). It is an important fact that one can efficiently compute the (stochastic) **gradients** of (2) via the so-called *backpropagation* [4], and therefore one can run first-order algorithms. The details will be covered in the coming lectures and exercises.

### 3 Randomness in Statistical Learning Problems, and Stochastic Gradient Descent

#### PROBLEM 3: RECOGNIZING DIFFERENT RANDOMNESS

There are many different randomness in modern data science or machine learning problems. The purpose of this exercise is to help you get a deeper understanding of them.

This course is all about inferring from data, and the data from real world is often random. Besides this intrinsic randomness, another common source of randomness in modern applications is the *randomized algorithms*. It is extremely important that you have a clear picture of what randomness is truly relevant for statistical inference, and what is only for computational purposes.

Consider the Gaussian linear model from Lecture 1: Let  $\mathbf{x}^{\mathbf{h}} \in \mathbb{R}^p$  and  $\mathbf{A} \in \mathbb{R}^{n \times p}$ . We have observations of the form

$$\mathbf{b} = \mathbf{A}\mathbf{x}^{\mathbf{h}} + \mathbf{w}, \quad (3)$$

where  $\mathbf{w} \sim \mathcal{N}(0, \sigma^2 I)$  is the Gaussian noise vector. We aim to solve the maximum likelihood estimator

$$\mathbf{x}_{ML}^* = \arg \min_{\mathbf{x} \in \mathbb{R}^p} \left\{ f(\mathbf{x}) := \frac{1}{2n} \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2^2 \right\}, \quad (4)$$

where we have normalized the the loss function by the number of measurements  $n$  (the number  $\frac{1}{2}$  is just for convenience later).

- (a) So far the only random component is the noise  $\mathbf{w}$ . Let  $\mathbb{E}_{\mathbf{w}}$  denote the expectation with respect to the randomness of  $\mathbf{w}$ . Compute  $\mathbb{E}_{\mathbf{w}} \|\mathbf{b} - \mathbf{A}\mathbf{x}^{\mathbf{h}}\|_2^2$ .
- (b) In practice, the measurement matrix  $\mathbf{A}$  is often random. Assume that the entries of  $\mathbf{A}$  are independent random variables with mean 0 and variance 1, and are independent of the noise  $\mathbf{w}$ . Let  $\mathbb{E}_{\mathbf{A}}$  denote the expectation with respect to the randomness of  $\mathbf{A}$ . Show that

$$\frac{1}{n} \mathbb{E}_{\mathbf{A}} \|\mathbf{A}\mathbf{x}\|_2^2 = \|\mathbf{x}\|_2^2 \quad (5)$$

for all  $\mathbf{x} \in \mathbb{R}^p$ .

(Hint: Let  $\mathbf{a}_i^{\top}$  be a row of  $\mathbf{A}$ . What is  $\mathbb{E}_{\mathbf{A}} |\langle \mathbf{a}_i, \mathbf{x} \rangle|^2$ ? Can you compute  $\mathbb{E}_{\mathbf{A}} \|\mathbf{A}\mathbf{x}\|_2^2$  through  $\mathbb{E}_{\mathbf{A}} |\langle \mathbf{a}_i, \mathbf{x} \rangle|^2$ ?)

- (c) Show that the following useful basic inequality holds:

$$\|\mathbf{A}(\mathbf{x}_{ML}^* - \mathbf{x}^{\mathbf{h}})\|_2^2 \leq 2 \langle \mathbf{w}, \mathbf{A}(\mathbf{x}_{ML}^* - \mathbf{x}^{\mathbf{h}}) \rangle. \quad (6)$$

(Hint: The maximum likelihood estimator minimizes the loss function, so you can compare the values of the loss function when substituting in  $\mathbf{x}_{ML}^*$  and any other  $\mathbf{x}$ .)

- (d) What we ultimately care about is the estimation error:  $\mathbb{E}_{\mathbf{A}, \mathbf{w}} \|\mathbf{x}_{ML}^* - \mathbf{x}^b\|_2^2$ . In view of (b) and (c), one might be tempted to conclude that

$$\mathbb{E}_{\mathbf{A}, \mathbf{w}} \|\mathbf{x}_{ML}^* - \mathbf{x}^b\|_2^2 \leq 2\mathbb{E}_{\mathbf{A}, \mathbf{w}} \left\langle \mathbf{w}, \frac{1}{n} \mathbf{A}(\mathbf{x}_{ML}^* - \mathbf{x}^b) \right\rangle. \quad (7)$$

Please argue carefully why this argument is NOT true.

(Hint: In part (b) we considered a fixed, deterministic  $\mathbf{x}$ . Is  $\mathbf{x}_{ML}^* - \mathbf{x}^b$  deterministic? If not, what randomness does it depend on?)

- (e) We introduce the important **Stochastic Gradient Descent (SGD)** in the exercise.

Recall Gradient Descent (GD) for minimizing (4):

- Choose  $\mathbf{x}_0$  arbitrarily.
- Do  $\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)$  for some predetermined step-sizes  $\alpha_k > 0$ .

Recall also that for (4), the gradient at a point  $\mathbf{x}$  is  $\nabla f(\mathbf{x}) = \frac{1}{n} \mathbf{A}^\top (\mathbf{A}\mathbf{x} - \mathbf{b})$ .

Consider a randomized algorithm as follows. At the current iterate  $\mathbf{x}_k$ , an index  $i \in \{1, 2, \dots, n\}$  is selected uniformly at random. We then replace the gradient at  $\mathbf{x}_k$  by the vector  $(\langle \mathbf{a}_i, \mathbf{x}_k \rangle - b_i) \mathbf{a}_i$ , where  $\mathbf{a}_i^\top$  is the  $i$ -th row of  $\mathbf{A}$  and  $b_i$  is the  $i$ -th element of  $\mathbf{b}$ . All other steps are the same as GD. Let  $\mathbb{E}_{SGD}$  denote the expectation with respect to the randomness of this algorithm. Show that

$$\mathbb{E}_{SGD} ((\langle \mathbf{a}_i, \mathbf{x}_k \rangle - b_i) \mathbf{a}_i) = \frac{1}{n} \mathbf{A}^\top (\mathbf{A}\mathbf{x}_k - \mathbf{b}). \quad (8)$$

That is, the algorithm is a randomized version of Gradient Descent, hence the name.

(Hint: Recall that  $\mathbf{a}_i^\top$  is a row of  $\mathbf{A}$ , and therefore it is a column of  $\mathbf{A}^\top$ .)

*Remark:* There are many reasons for using SGD instead of GD; we refer the interested students to [1] for a gentle introduction of SGD. Please do bear in mind that SGD is extremely important in practice. Ever heard of deep learning? AlphaGo? Your interest might be piqued if you know that SGD is the go-to algorithm for these state-of-the-art learning machines.

- (f) Under the assumptions in (b), show that

$$\mathbb{E}_{\mathbf{A}, \mathbf{w}} ((\langle \mathbf{a}_i, \mathbf{x} \rangle - b_i) \mathbf{a}_i) = \mathbf{x} - \mathbf{x}^b \quad (9)$$

for any  $\mathbf{x}$ .

- (g) Under the assumptions in (b), show that, for any  $\mathbf{x}$ ,

$$\frac{1}{n} \mathbb{E}_{\mathbf{A}, \mathbf{w}} \mathbf{A}^\top (\mathbf{A}\mathbf{x} - \mathbf{b}) = \mathbf{x} - \mathbf{x}^b. \quad (10)$$

(Hint: There are many ways of proving this. Combining (e) and (f) gives a very simple proof.)

## 4 Multinomial logistic regression and language model training

### PROBLEM 4: TOWARDS TRAINING LANGUAGE MODELS

Recall that in the lecture we talk about logistic regression for binary classification problems. Let  $\mathbf{x}^b \in \mathbb{R}^p$ . Let  $\mathbf{a}_1, \dots, \mathbf{a}_n \in \mathbb{R}^p$  be given. The sample is given by  $\mathbf{b} := (b_1, \dots, b_n) \in \{0, 1\}^n$ . The classifier  $h_x$  estimates the probability  $P(b = 1)$  by outputting a scalar  $h_x$  such that  $h_x := [1 + \exp(-\langle \mathbf{a}, \mathbf{x} \rangle)]^{-1} = P(b = 1)$ .

- (a) Show that maximizing the likelihood is equivalent to minimizing the cross-entropy between the real distribution and estimated distribution.
- (b) In logistic regression, we only deal with two classes. But in reality, there are many situations where multiple classes are involved. In this case, we need to generalize the logistic regression to multinomial logistic regression (sometimes called the softmax regression).

Suppose there are  $K$  classes ( $K \geq 2$ ). The classifier aims to output a vector  $\mathbf{h}_x \in \mathbb{R}^K$  to estimate the probability of each class such that its  $k$  element is  $P(b = k | \mathbf{a})$ . Note that in this case, the classifier is parameterized by a learnable matrix  $\mathbf{X} \in \mathbb{R}^{K \times p}$  to estimate the probability as follows:  $P(b = k | \mathbf{a}) = (\text{Softmax}(\mathbf{X}\mathbf{a}))^{[k]}$ , for  $k \in [K]$ , where Softmax is defined by:

$$\text{Softmax}(\mathbf{X}\mathbf{a}) = \begin{bmatrix} \frac{\exp((\mathbf{X}\mathbf{a})^{[1]})}{\sum_{i=1}^K \exp((\mathbf{X}\mathbf{a})^{[i]})} \\ \vdots \\ \frac{\exp((\mathbf{X}\mathbf{a})^{[K]})}{\sum_{i=1}^K \exp((\mathbf{X}\mathbf{a})^{[i]})} \end{bmatrix} \in \mathbb{R}^K.$$

Again, show that maximizing the likelihood is equivalent to minimizing the cross entropy between the real distribution and the estimated distribution.

- (c) To some extent, training a language model is essentially tackling a multinomial logistic regression problem where the model needs to predict next word given previous words. A neural network  $\mathbf{h}_x$  parameterized by  $\mathbf{x}$  is used as an ML estimator, as mentioned in the lecture. Specifically,  $\mathbf{h}_x$  aims to predict the  $t$  token given  $t - 1$  tokens.

$$\begin{aligned} \min_{\mathbf{x}} -\log p_x(\mathbf{b}_{1:T}) &= -\log \left( \prod_{t=1}^T p_x(\mathbf{b}_t | \mathbf{b}_{1:t-1}) \right) = \sum_{t=1}^T (-\log p_x(\mathbf{b}_t | \mathbf{b}_{1:t-1})) \\ &= \sum_{t=1}^T (-\log \mathbf{h}_x(\mathbf{b}_{1:t-1})^{["\mathbf{b}_t"]}) = \text{cross entropy loss.} \end{aligned}$$

Write down the SGD algorithm for training a language model given a corpus that consists of  $n$  sentences.

## References

- [1] LÉON BOTTOU *Stochastic Gradient Descent Tricks, Neural Networks: Tricks of the Trade*, page 421-436, Springer 2012.
- [2] SIMON HAYKIN *Neural networks: a comprehensive foundation*, Prentice Hall PTR 1994.
- [3] ANDREW BARRON *Universal approximation bounds for superpositions of a sigmoidal function*, *IEEE Transactions on Information theory*, 1993.
- [4] DAVID RUMELHART, GEOFFREY HINTON, AND RONALD WILLIAMS *Learning representations by back-propagating errors*, *Cognitive modeling*, 1988.
- [5] Zhu, Z., Liu, F., Chrysos, G. & Cevher, V. Robustness in deep learning: The good (width), the bad (depth), and the ugly (initialization). *ArXiv Preprint ArXiv:2209.07263*. (2022)