

# NYU Center for Data Science: DS-GA 1003

## Machine Learning and Computational Statistics (Spring 2019)

Brett Bernstein

January 30, 2019

**Instructions:** Following most lab and lecture sections, we will be providing concept checks for review. Each concept check will:

- List the lab/lecture learning objectives. You will be responsible for mastering these objectives, and demonstrating mastery through homework assignments, exams (midterm and final), and on the final course project.
- Include concept check questions. These questions are intended to reinforce the lab/lectures, and help you master the learning objectives.

You are strongly encourage to complete all concept check questions, and to discuss these (and related) problems on Piazza and at office hours. However, problems marked with a (★) are considered optional.

## Lecture 1: Introduction to Statistical Learning Theory

### Topic 1: Statistical Learning Theory

#### Learning Objectives

1. Identify the input, action, and outcome spaces for a given machine learning problem.
2. Provide an example for which the action space and outcome spaces are the same and one for which they are different.
3. Explain the relationships between the decision function, the loss function, the input space, the action space, and the outcome space.
4. Define the risk of a decision function and a Bayes decision function.
5. Provide example decision problems for which the Bayes risk is 0 and the Bayes risk is nonzero.
6. Know the Bayes decision functions for square loss and multiclass 0/1 loss.

7. Define the empirical risk for a decision function and the empirical risk minimizer.
8. Explain what a hypothesis space is, and how it can be used with constrained empirical risk minimization to control overfitting.

### Concept Check Questions

1. Suppose  $\mathcal{A} = \mathcal{Y} = \mathbb{R}$  and  $\mathcal{X}$  is some other set. Furthermore, assume  $P_{\mathcal{X} \times \mathcal{Y}}$  is a discrete joint distribution. Compute a Bayes decision function when the loss function  $\ell : \mathcal{A} \times \mathcal{Y} \rightarrow \mathbb{R}$  is given by

$$\ell(a, y) = \mathbf{1}(a \neq y),$$

the 0 – 1 loss.

*Solution.* The Bayes decision function  $f^*$  satisfies

$$f^* = \arg \min_f R(f) = \arg \min_f \mathbb{E}[\mathbf{1}(f(X) \neq Y)] = \arg \min_f P(f(X) \neq Y),$$

where  $(X, Y) \sim P_{\mathcal{X} \times \mathcal{Y}}$ . Let

$$f_1(x) = \arg \max_y P(Y = y \mid X = x),$$

the maximum a posteriori estimate of  $Y$ . If there is a tie, we choose any of the maximizers. If  $f_2$  is another decision function we have

$$\begin{aligned} P(f_1(X) \neq Y) &= \sum_x P(f_1(x) \neq Y \mid X = x)P(X = x) \\ &= \sum_x (1 - P(f_1(x) = Y \mid X = x))P(X = x) \\ &\leq \sum_x (1 - P(f_2(x) = Y \mid X = x))P(X = x) \quad (\text{Defn of } f_1) \\ &= \sum_x P(f_2(x) \neq Y \mid X = x)P(X = x) \\ &= P(f_2(X) \neq Y). \end{aligned}$$

Thus  $f^* = f_1$ .

2. (★) Suppose  $\mathcal{A} = \mathcal{Y} = \mathbb{R}$ ,  $\mathcal{X}$  is some other set, and  $\ell : \mathcal{A} \times \mathcal{Y} \rightarrow \mathbb{R}$  is given by  $\ell(a, y) = (a - y)^2$ , the square error loss. What is the Bayes risk and how does it compare with the variance of  $Y$ ?

*Solution.* From Homework 1 we know that the Bayes decision function is given by  $f^*(x) = \mathbb{E}[Y \mid X = x]$ . Thus the Bayes risk is given by

$$\mathbb{E}[(f^*(X) - Y)^2] = \mathbb{E}[(\mathbb{E}[Y \mid X] - Y)^2] = \mathbb{E}[\mathbb{E}[(\mathbb{E}[Y \mid X] - Y)^2 \mid X]] = \mathbb{E}[\text{Var}(Y \mid X)],$$

where we applied the law of iterated expectations. The law of total variance states that

$$\text{Var}(Y) = \mathbb{E}[\text{Var}(Y \mid X)] + \text{Var}[\mathbb{E}(Y \mid X)].$$

This proves the Bayes risk satisfies

$$\mathbb{E}[\text{Var}(Y|X)] = \text{Var}(Y) - \text{Var}[\mathbb{E}(Y|X)] \leq \text{Var}(Y).$$

Recall from Homework 1 that  $\text{Var}(Y)$  is the Bayes risk when we estimate  $Y$  without any input  $X$ . This shows that using  $X$  in our estimation reduces the Bayes risk, and that the improvement is measured by  $\text{Var}[\mathbb{E}(Y|X)]$ . As a sanity check, note that if  $X, Y$  are independent then  $\mathbb{E}(Y|X) = \mathbb{E}(Y)$  so  $\text{Var}[\mathbb{E}(Y|X)] = 0$ . If  $X = Y$  then  $\mathbb{E}(Y|X) = Y$  and  $\text{Var}[\mathbb{E}(Y|X)] = \text{Var}(Y)$ .

The prominent role of variance in our analysis above is due to the fact that we are using the square loss.

3. Let  $\mathcal{X} = \{1, \dots, 10\}$ , let  $\mathcal{Y} = \{1, \dots, 10\}$ , and let  $A = \mathcal{Y}$ . Suppose the data generating distribution,  $P$ , has marginal  $X \sim \text{Unif}\{1, \dots, 10\}$  and conditional distribution  $Y|X = x \sim \text{Unif}\{1, \dots, x\}$ . For each loss function below give a Bayes decision function.

- (a)  $\ell(a, y) = (a - y)^2$ ,
- (b)  $\ell(a, y) = |a - y|$ ,
- (c)  $\ell(a, y) = \mathbf{1}(a \neq y)$ .

*Solution.*

- (a) From Homework 1 we know that  $f^*(x) = \mathbb{E}[Y|X = x] = (x + 1)/2$ .
- (b) From Homework 1, we know that  $f^*(x)$  is the conditional median of  $Y$  given  $X = x$ . If  $x$  is odd, then  $f^*(x) = (x + 1)/2$ . If  $x$  is even, then we can choose any value in the interval

$$\left[ \left\lfloor \frac{x+1}{2} \right\rfloor, \left\lceil \frac{x+1}{2} \right\rceil \right].$$

- (c) From question 1 above, we know that  $f^*(x) = \arg \max_y P(Y = y|X = x)$ . Thus we can choose any integer between 1 and  $x$ , inclusive, for  $f^*(x)$ .

4. Show that the empirical risk is an unbiased and consistent estimator of the Bayes risk. You may assume the Bayes risk is finite.

*Solution.* We assume a given loss function  $\ell$  and an i.i.d. sample  $(x_1, y_1), \dots, (x_n, y_n)$ . To show it is unbiased, note that

$$\begin{aligned} \mathbb{E}[\hat{R}_n(f)] &= \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i) \right] \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\ell(f(x_i), y_i)] \quad (\text{Linearity of } \mathbb{E}) \\ &= \mathbb{E}[\ell(f(x_1), y_1)] \quad (\text{i.i.d.}) \\ &= R(f). \end{aligned}$$

For consistency, we must show that as  $n \rightarrow \infty$  we have  $\hat{R}_n(f) \rightarrow R(f)$  with probability 1. Letting  $z_i = \ell(f(x_i), y_i)$ , we see that the  $z_i$  are i.i.d. with finite mean. Thus consistency follows by applying the strong law of large numbers.

5. Let  $\mathcal{X} = [0, 1]$  and  $\mathcal{Y} = \mathcal{A} = \mathbb{R}$ . Suppose you receive the  $(x, y)$  data points  $(0, 5)$ ,  $(.2, 3)$ ,  $(.37, 4.2)$ ,  $(.9, 3)$ ,  $(1, 5)$ . Throughout assume we are using the 0 – 1 loss.
  - (a) Suppose we restrict our decision functions to the hypothesis space  $\mathcal{F}_1$  of constant functions. Give a decision function that minimizes the empirical risk over  $\mathcal{F}_1$  and the corresponding empirical risk. Is the empirical risk minimizing function unique?
  - (b) Suppose we restrict our decision functions to the hypothesis space  $\mathcal{F}_2$  of piecewise-constant functions with at most 1 discontinuity. Give a decision function that minimizes the empirical risk over  $\mathcal{F}_2$  and the corresponding empirical risk. Is the empirical risk minimizing function unique?

*Solution.*

- (a) We can let  $\hat{f}(x) = 5$  or  $\hat{f}(x) = 3$  and obtain the minimal empirical risk of 3/5. Thus the empirical risk minimizer is not unique.
  - (b) One solution is to let  $\hat{f}(x) = 5$  for  $x \in [0, .1]$  and  $\hat{f}(x) = 3$  for  $x \in (.1, 1]$  giving an empirical risk of 2/5. There are uncountably many empirical risk minimizers, so again we do not have uniqueness.
6. (★) Let  $\mathcal{X} = [-10, 10]$ ,  $\mathcal{Y} = \mathcal{A} = \mathbb{R}$  and suppose the data generating distribution has marginal distribution  $X \sim \text{Unif}[-10, 10]$  and conditional distribution  $Y|X = x \sim \mathcal{N}(a + bx, 1)$  for some fixed  $a, b \in \mathbb{R}$ . Suppose you are also given the following data points:  $(0, 1)$ ,  $(0, 2)$ ,  $(1, 3)$ ,  $(2.5, 3.1)$ ,  $(-4, -2.1)$ .
  - (a) Assuming the 0 – 1 loss, what is the Bayes risk?
  - (b) Assuming the square error loss  $\ell(a, y) = (a - y)^2$ , what is the Bayes risk?
  - (c) Using the full hypothesis space of all (measurable) functions, what is the minimum achievable empirical risk for the square error loss?
  - (d) Using the hypothesis space of all affine functions (i.e., of the form  $f(x) = cx + d$  for some  $c, d \in \mathbb{R}$ ), what is the minimum achievable empirical risk for the square error loss?
  - (e) Using the hypothesis space of all quadratic functions (i.e., of the form  $f(x) = cx^2 + dx + e$  for some  $c, d, e \in \mathbb{R}$ ), what is the minimum achievable empirical risk for the square error loss?

*Solution.*

(a) For any decision function  $f$  the risk is given by

$$\mathbb{E}[\mathbf{1}(f(X) \neq Y)] = P(f(X) \neq Y) = 1 - P(f(X) = Y) = 1.$$

To see this note that

$$P(f(X) = Y) = \frac{1}{20\sqrt{2\pi}} \int_{-10}^{10} \int_{-\infty}^{\infty} \mathbf{1}(f(x) = y) e^{-(y-a-bx)^2/2} dy dx = \frac{1}{20\sqrt{2\pi}} \int_{-10}^{10} 0 dx = 0.$$

Thus every decision function is a Bayes decision function, and the Bayes risk is 1.

(b) By problem 2 above we know the Bayes risk is given by

$$\mathbb{E}[\text{Var}(Y|X)] = \mathbb{E}[1] = 1,$$

since  $\text{Var}(Y|X = x) = 1$ .

(c) We choose  $\hat{f}$  such that

$$\hat{f}(0) = 1.5, \hat{f}(1) = 3, \hat{f}(2.5) = 3.1, \hat{f}(-4) = 2.1,$$

and  $\hat{f}(x) = 0$  otherwise. Then we achieve the minimum empirical risk of 1/10.

(d) Letting

$$A = \begin{pmatrix} 1 & 0 \\ 1 & 0 \\ 1 & 1 \\ 1 & 2.5 \\ 1 & -4 \end{pmatrix}, \quad y = \begin{pmatrix} 1 \\ 2 \\ 3 \\ 3.1 \\ -2.1 \end{pmatrix}$$

we obtain (using a computer)

$$\hat{w} = \begin{pmatrix} \hat{d} \\ \hat{c} \end{pmatrix} = (A^T A)^{-1} A^T y = \begin{pmatrix} 1.4856 \\ 0.8556 \end{pmatrix}.$$

This gives

$$\hat{R}_5(\hat{f}) = \frac{1}{5} \|A\hat{w} - y\|_2^2 = 0.2473.$$

[Aside: In general, to solve systems like the one above on a computer you shouldn't actually invert the matrix  $A^T A$ , but use something like  $w=A \backslash y$  in Matlab which performs a QR factorization of  $A$ .]

(e) Letting

$$A = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 1 & 1 \\ 1 & 2.5 & 6.25 \\ 1 & -4 & 16 \end{pmatrix}, \quad y = \begin{pmatrix} 1 \\ 2 \\ 3 \\ 3.1 \\ -2.1 \end{pmatrix}$$

we obtain (using a computer)

$$\hat{w} = \begin{pmatrix} \hat{e} \\ \hat{d} \\ \hat{c} \end{pmatrix} = (A^T A)^{-1} A^T y = \begin{pmatrix} 1.7175 \\ 0.7545 \\ -0.0521 \end{pmatrix}.$$

This gives

$$\hat{R}_5(\hat{f}) = \frac{1}{5} \|A\hat{w} - y\|_2^2 = 0.1928.$$