STAT 460/560 Class 11: Statistical Decision Theory

Ben Bloem-Reddy

Reading: Chapter 12, [Was04].

At a high level, statistical decision theory is a framework for assessing (and comparing the performance of statistical procedures). It connects statistical tasks like parameter estimation to real-world/practical quantities such as cost, risk, etc.

As usual, $\theta \in \Theta$ is a parameter, and $\hat{\theta}$ is an estimator. In decision theoretic terms, $\hat{\theta}$ is a **decision rule** (i.e., a function that maps data to a decision, or **action**). Note that we are not limited to parameter estimation. As an example more closely related to the language, a decision rule could actually be a data-driven decision on which action to take (e.g., which experiment to run next).

For simplicity, we will look for now only at parameter estimation. The discrepancy of $\hat{\theta}$ and θ is measured by a **loss function**, $L \colon \Theta \times \Theta \to \mathbb{R}$. A common example is **squared error loss**,

$$L(\theta, \hat{\theta}) = (\hat{\theta} - \theta)^2 . \tag{11.1}$$

The **risk** of an estimator $\hat{\theta}$ is a function of θ ,

$$R(\theta, \hat{\theta}) = \mathbb{E}_{\theta}[L(\theta, \hat{\theta}(X))] = \int L(\theta, \hat{\theta}(x)) f(x; \theta) \ dx \ . \tag{11.2}$$

When the loss is squared error, the resulting risk is just the MSE,

$$R(\theta, \hat{\theta}) = \mathbb{E}_{\theta}[(\hat{\theta}(X) - \theta)^2] = \operatorname{Var}_{\theta}(\hat{\theta}(X)) + \operatorname{bias}_{\theta}(\hat{\theta}(X))^2.$$
 (11.3)

1. Comparing risk functions

How do we compare two functions? Let's start with some examples.

Example 11.1. Let $X \sim \mathcal{N}(\theta, 1)$, and consider two estimators under squared error loss: $\hat{\theta}_1 = X$ and $\hat{\theta}_2 = 3$. The risk functions are $R(\theta, \hat{\theta}_1) = 1$ and $R(\theta, \hat{\theta}_2) = (3 - \theta)^2$. If $2 < \theta < 4$ then $\hat{\theta}_2$ has lower risk; otherwise $\hat{\theta}_1$ has lower risk. So neither estimator is uniformly better than the other.

Activity 11.1. Let $X_1, \ldots, X_n \sim_{\text{\tiny IID}} \text{Bern}(p)$, and consider squared error loss. Let $\hat{p}_1 = \bar{X}$. Compute the risk. (To make things easy, recall that this estimator is unbiased.)

Now suppose that we take a Bayesian approach, and use a Beta(α, β) prior, with $\hat{p}_2 = \mathbb{E}[p|X^n]$. Write down \hat{p}_2 and calculate its risk for $\alpha = \beta = \sqrt{n}/2$.

Plot or sketch the risk functions of the two estimators. Is either one uniformly better? Do you notice anything interesting about \hat{p}_2 ?

Solution: Since the variance of \bar{X}_n is p(1-p)/n, the risk of \hat{p}_1 is p(1-p)/n.

The posterior here is $Beta(\alpha + n\bar{X}_n, \beta + n - n\bar{X}_n)$, so

$$\hat{p}_2 = \mathbb{E}[p|X^n] = \frac{\alpha + n\bar{X}_n}{\alpha + \beta + n} = \frac{\sqrt{n}/2 + n\bar{X}_n}{\sqrt{n} + n} = \frac{1/2 + \sqrt{n}\bar{X}_n}{1 + \sqrt{n}}$$
.

The bias of \hat{p}_2 is

$$\frac{1/2 + \sqrt{n}p}{1 + \sqrt{n}} - p = \frac{1/2 - p}{1 + \sqrt{n}} ,$$

and the variance is

$$\frac{n}{(1+\sqrt{n})^2}np(1-p) = \frac{p(1-p)}{(1+\sqrt{n})^2}.$$

Putting these together, we find that the risk of \hat{p}_2 is

$$R(p, \hat{p}_2) = \frac{(1/2 - p)^2}{(1 + \sqrt{n})^2} + \frac{p(1 - p)}{(1 + \sqrt{n})^2} = \frac{1}{4(1 + \sqrt{n})^2}.$$

This is constant with respect to p, and is lower than the risk of \hat{p}_1 when $p(1-p) < \frac{n}{4(1+\sqrt{n})^2}$.

In both of these examples, if we don't know the true value of θ (which we don't) then how should we choose between the two estimators? Much of statistical decision theory deals with this seemingly simple question.

There are two widely used ways to turn the risk function into a number. The maximum risk is

$$\bar{R}(\hat{\theta}) = \sup_{\theta} R(\theta, \hat{\theta}) . \tag{11.4}$$

The Bayes risk is

$$r(f,\hat{\theta}) = \int R(\theta,\hat{\theta})f(\theta) \ d\theta \ , \tag{11.5}$$

where $f(\theta)$ is a prior PDF/PMF.

Activity 11.2. Find the maximum risk for the two estimators from Activity 11.1. Which is lower?

Find the Bayes risk for the two estimators under a uniform prior.

Solution: The risk of \hat{p}_1 is maximized at p = 1/2, in which case it equals $\frac{1}{4n}$.

The risk of \hat{p}_2 is constant, equal to $\frac{1}{4(1+\sqrt{n})^2}$. This is lower than the maximum risk of \hat{p}_1 .

Under a uniform prior on p, the Bayes risk of \hat{p}_1 is

$$\frac{1}{n} \int_0^1 p(1-p)dp = \frac{1}{n} B(2,2) = \frac{1}{n} \frac{\Gamma(2)\Gamma(2)}{\Gamma(4)} = \frac{1}{6n} ,$$

where B is the beta function, and we used the gamma function identity $\Gamma(n) = (n-1)!$.

Since the risk of \hat{p}_2 is constant, its Bayes risk is equal to its maximum risk.

Exercise 11.1. Find the Bayes risk for the two estimators under a Beta $(\sqrt{n}/2, \sqrt{n}/2)$ prior. Show that the ratio of Bayes risks can be written as

$$\frac{r(\text{Beta}(\sqrt{n}/2, \sqrt{n}/2), \hat{p}_1)}{r(\text{Beta}(\sqrt{n}/2, \sqrt{n}/2), \hat{p}_2)} = 1 + \frac{1}{\sqrt{n}}.$$
 (11.6)

Solution: Under this prior, the Bayes risk of \hat{p}_1 is

$$\frac{1}{n} \int_0^1 p(1-p) \frac{p^{\sqrt{n}/2-1}(1-p)^{\sqrt{n}/2-1}}{B(\sqrt{n}/2, \sqrt{n}/2)} dp = \frac{1}{n} \frac{B(\sqrt{n}/2+1, \sqrt{n}/2+1)}{B(\sqrt{n}/2, \sqrt{n}/2)} = \frac{1}{4\sqrt{n}(\sqrt{n}+1)}.$$

Taking the ratio with the Bayes risk of \hat{p}_2 yields the desired result.

Note that the Beta $(\sqrt{n}/2, \sqrt{n}/2)$ prior isn't likely to be motivated by any prior belief or knowledge, except in the case n = 1. We'll see later where the prior comes from.

Based on these one-number summaries of risk, we can devise general strategies for finding "good" estimators: we might try to find $\hat{\theta}$ to minimize the maximum risk; alternatively, we might find $\hat{\theta}$ to minimize the Bayes risk. A **minimax estimator** $\hat{\theta}$ is any estimator (it may not be unique) such that

$$\sup_{\theta} R(\theta, \hat{\theta}) = \inf_{\tilde{\theta}} \sup_{\theta} R(\theta, \tilde{\theta}) , \qquad (11.7)$$

where the infimum is taken over all possible estimators $\hat{\theta}$. The **Bayes estimator** with respect to prior f is any estimator $\hat{\theta}$ such that

$$r(f,\hat{\theta}) = \inf_{\tilde{\theta}} r(f,\tilde{\theta}) . \tag{11.8}$$

Bayes estimators are often easier to find. To see why, define the **posterior risk** as

$$r(\hat{\theta}|x) = \int L(\theta, \hat{\theta}(x)) f(\theta|x) \ d\theta \ . \tag{11.9}$$

Theorem 11.1. The Bayes risk satisfies

$$r(f,\hat{\theta}) = \int r(\hat{\theta}|x)m(x) dx, \qquad (11.10)$$

where $m(x) = \int f(x|\theta)f(\theta) d\theta$. Therefore, under squared error loss, the Bayes estimator is the posterior expectation,

$$\hat{\theta}(x) = \int \theta f(\theta|x) \ d\theta = \mathbb{E}[\theta|X = x] \ . \tag{11.11}$$

Exercise 11.2. Prove Theorem 11.1.

Exercise 11.3. For $X_1, \ldots, X_n \sim_{\text{\tiny IID}} \operatorname{Bern}(p)$ and a $\operatorname{Beta}(\alpha, \beta)$ prior, what is the Bayes estimator under squared error loss? Under $L(\theta, \hat{\theta}) = |\theta - \hat{\theta}|$?

2. Minimax decision rules

Finding minimax decision rules is in general hard, and there is no one-size-fits-all method for doing so. In some situations, a Bayes estimator can be minimax.

Theorem 11.2. Let $\hat{\theta}^f$ be a Bayes estimator for some prior f, and therefore

$$r(f, \hat{\theta}^f) = \inf_{\hat{\theta}} r(f, \hat{\theta}) . \tag{11.12}$$

Suppose that

$$R(\theta, \hat{\theta}^f) \le r(f, \hat{\theta}^f) \quad \text{for all } \theta \in \Theta \ .$$
 (11.13)

Then $\hat{\theta}^f$ is minimax and f is called a **least favorable prior**.

 $^{^{1}}$ Why would our prior belief/knowledge about p depend on n?

Proof. Suppose that $\hat{\theta}^f$ is not minimax. Then there is some estimator $\hat{\theta}_0$ such that

$$\sup_{\theta} R(\theta, \hat{\theta}_0) < \sup_{\theta} R(\theta, \hat{\theta}^f) . \tag{11.14}$$

Since the average of a function is always less than or equal to its maximum, $r(f, \hat{\theta}_0) \leq \sup_{\theta} R(\theta, \hat{\theta}_0)$. Then

$$r(f, \hat{\theta}_0) \le \sup_{\theta} R(\theta, \hat{\theta}_0) < \sup_{\theta} R(\theta, \hat{\theta}^f) \le r(f, \hat{\theta}^f) . \tag{11.15}$$

This is a contradiction, because $\hat{\theta}^f$ is the Bayes estimator under f. Hence, $\hat{\theta}^f$ is minimax.

Corollary 11.3. Suppose that $\hat{\theta}$ is a Bayes estimator under some prior f, and further that $\hat{\theta}$ has constant risk c. Then $\hat{\theta}$ is minimax.

Exercise 11.4. Prove Corollary 11.3.

Solution: If $\hat{\theta}$ has constant risk then its Bayes risk is $r(f, \hat{\theta}) = \int cf(\theta) = c$. Therefore, $R(\theta, \hat{\theta}) = c = r(f, \hat{\theta})$, and so Theorem 11.2 applies.

Activity 11.3. Recall that in the previous class, with $X_1, \ldots, X_n \sim_{\text{\tiny IID}} \text{Bern}(p)$ and squared error loss, we analyzed $\hat{p}_2 = \mathbb{E}[p|X^n]$ when the prior is $\text{Beta}(\sqrt{n}/2, \sqrt{n}/2)$. Argue that \hat{p}_2 is minimax.

Solution: We saw that \hat{p}_2 has constant risk and that it is the Bayes estimator under a Beta $(\sqrt{n}/2, \sqrt{n}/2)$ prior. By Corollary 11.3, it is minimax.

3. Optional: Prediction and empirical risk minimization

When the task is parameter estimation, risk is typically a purely theoretical construct: the risk usually depends on an unknown parameter and therefore cannot be estimated from data. However, consider the problem of predicting Y from X with an estimator $\hat{Y} = \hat{\theta}(X)$. The estimator will be fit to pairs $((X_1, Y_1), \dots, (X_n, Y_n)) \sim_{\text{IID}} P$, and the risk is

$$R(P,\hat{\theta}) = \mathbb{E}_P[L(Y,\hat{\theta}(X))] = \int L(y,\hat{\theta}(x))dP(x,y). \qquad (11.16)$$

This depends on the unknown distribution P. However, it can be estimated by the **empirical risk**

$$\hat{R}(\hat{\theta}) = R(\hat{P}_n, \hat{\theta}) = \frac{1}{n} \sum_{i=1}^n L(Y_i, \hat{\theta}(X_i)) . \tag{11.17}$$

Minimizing this type of objective is called **empirical risk minimization**. This is the default procedure for most machine learning methods (though its implementation for different models and massive datasets is fairly complicated).

Activity 11.4. Suppose we have \mathbb{R} -valued pairs $((X_1, Y_1), \dots, (X_n, Y_n)) \sim_{\text{IID}} P$, and we're using squared error loss. Show that an estimator that achieves the smallest risk is

$$\hat{\theta}(x) = \mathbb{E}_P[Y|X=x] . \tag{11.18}$$

Solution:

$$R(P,\hat{\theta}) = \mathbb{E}_P[L(Y,\hat{\theta}(X))] \tag{11.19}$$

$$= \mathbb{E}_P[\mathbb{E}_P[L(Y, \hat{\theta}(X))|X]] \tag{11.20}$$

$$= \mathbb{E}_{P}[\mathbb{E}_{P}[(Y - \hat{\theta}(X))^{2}|X]]. \tag{11.21}$$

As in the case of finding the Bayes estimator of θ under squared error loss, if we can find $\hat{\theta}$ that minimizes $\mathbb{E}_P[(Y - \hat{\theta}(X))^2 | X = x]$ pointwise in x then we know that the overall risk is minimized by $\hat{\theta}$. For fixed x, the minimizer is

$$a^*(x) = \arg\inf_{a} \mathbb{E}_P[(Y-a)^2 | X = x] = \mathbb{E}_P[Y | X = x],$$
 (11.22)

and therefore $\hat{\theta}(x) = \mathbb{E}_P[Y|X=x]$ achieves the smallest risk.

Exercise 11.5. Suppose that $Y \in \{0,1\}, ((X_1,Y_1),\ldots,(X_n,Y_n)) \sim_{\text{\tiny IID}} P$, and we're using 0-1 loss,

$$L(Y, \hat{\theta}(X)) = \mathbf{I}\{Y \neq \hat{\theta}(X)\}. \tag{11.23}$$

Show that an estimator that achieves the smallest risk is

$$\hat{\theta}(x) = \begin{cases} 1 & \text{if } P(Y = 1|X = x) \ge \frac{1}{2} \\ 0 & \text{o.w.} \end{cases}$$
 (11.24)

This is called the **Bayes classifier**.

Solution: This is similar to the previous activity. The conditional (on X) risk is

$$\mathbb{E}_{P}[\mathbf{I}\{Y \neq \hat{\theta}(X)\}|X] = P(Y \neq \hat{\theta}(X)\}|X) \tag{11.25}$$

$$= P(Y = 1|X)\mathbf{I}\{\hat{\theta}(X) = 0\} + P(Y = 0|X)\mathbf{I}\{\hat{\theta}(X) = 1\}$$
(11.26)

$$= P(Y = 1|X)\mathbf{I}\{\hat{\theta}(X) = 0\} + (1 - P(Y = 1|X))\mathbf{I}\{\hat{\theta}(X) = 1\}.$$
 (11.27)

This is minimized by

$$\hat{\theta}(x) = \begin{cases} 1 & \text{if } P(Y = 1|X = x) \ge 1 - P(Y = 1|X = x) \\ 0 & \mathbb{P}(Y = 1|X = x) < 1 - P(Y = 1|X = x) \end{cases}$$
(11.28)

For more on the statistical aspects of machine learning, see Hastie, Tibshirani, and Friedman [HTF09], which is now somewhat outdated but still a great place to start. Wainwright [Wai19] has more advanced material.

References

- [HTF09] T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning*. Springer New York, 2009. URL: https://hastie.su.domains/ElemStatLearn/.
- [Wai19] M. J. Wainwright. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge University Press, 2019.
- [Was04] L. Wasserman. All of Statistics. Springer New York, 2004.