**Disclaimer**: *These notes have not been subjected to the usual scrutiny reserved for formal publications. They may be distributed outside this class only with the permission of the Instructor.*

## 6.1 Minimax Estimators and Worst-Case Optimality

In minimax estimation, we collapse our risk function by looking at the worst-case risk. Given $X \sim P_\theta$, where $\theta \in \Omega$, and a loss function $L(\theta, d)$, we want to find an estimator $\delta$ that minimizes the maximum risk:

$$\sup_{\theta \in \Omega} R(\theta, \delta).$$

Any such $\delta$ is called a minimax estimator. In Bayes estimation, we essentially had a single, general-purpose way of deriving a Ba estimator: minimize the posterior risk. The derivation of minimax estimators is often prescriptive and more problem-specific. However, we will develop a few tools which, w they apply, will identify minimax estimators.

Perhaps surprisingly, one of the most effective ways of finding minimax estimator: to restrict our attention to Bayes estimators. To understand the connection, we will n to introduce some additional notation. First we recall the definition of the Bayes risk minimum average risk) under any prior distribution $\Lambda$,

$$r_\Lambda = \inf_\delta r(\Lambda, \delta) = \inf_\delta \int_{\theta \in \Omega} R(\theta, \delta) d\Lambda(\theta).$$

**Definition 6.1** *We say that a prior $\Lambda$ is a least favorable prior if $r_\Lambda \geq r_{\Lambda'}$ for any other prior distribution $\Lambda'$. (Note that we always use the unmodified word "prior" to mean a proper prior.)*

**Theorem 6.2 (TPE 5.1.4)** *Suppose $\delta_\Lambda$ is Bayes for $\Lambda$ with*

$$r_{\Lambda \in \Omega} = \sup_\theta R(\theta, \delta_\Lambda)$$

*That is, the Bayes risk of $\delta_\Lambda$ is the maximum risk of $\delta_\Lambda$. Then,*

1. *$\delta_\Lambda$ is minimax*

2. *$\Lambda$ is a least favorable prior*

3. *If $\delta_\Lambda$ is the unique Bayes estimator for $\Lambda$ (a.s. for all $P_\theta$ ), then it is the unique minimax estimator.*

**Proof:** If $\delta$ is any other estimator, then we have that

$$\sup_{\theta \in \Omega} R(\theta, \delta) \geq \int R(\theta, \delta) d\Lambda(\theta) \geq \int R(\theta, \delta_\Lambda) d\Lambda(\theta) = \sup_{\theta \in \Omega} R(\theta, \delta_\Lambda)$$

where the first step holds because the worst-case risk of $\delta$ is greater than (or equal to) the average risk of $\delta$, the second step holds because $\delta_\Lambda$ is Bayes (and hence has an average risk no higher than that of $\delta$), and the third step holds because of our assumption that the Bayes risk of $\delta_\Lambda$ is equal to the worst-case risk. This implies that $\delta_\Lambda$ is minimax.

If $\delta_\Lambda$ is the unique Bayes estimator, then the second inequality above is strict for $\delta \neq \delta_\Lambda$, which implies that $\delta_\Lambda$ is the unique minimax. Let $\Lambda'$ be any other prior distribution. Then, we have that

$$r_{\Lambda'} = \inf_\delta \int R(\theta, \delta) d\Lambda'(\theta) \leq \int R(\theta, \delta_\Lambda) d\Lambda'(\theta) \leq \sup_\theta R(\theta, \delta_\Lambda) = r_\Lambda.$$

The first step first and second steps are by the definition of Bayes risk, and the third step holds because the worst-case risk of $\delta_\Lambda$ is no less than its average risk over the distribution $\Lambda'$. Since the worst-case risk of $\delta_\Lambda$ is its Bayes risk over $\Lambda$ (by our assumption), we can infer that $\Lambda$ is a least favorable prior distribution.   ■

Thus, we can find a minimax estimator by finding a Bayes estimator with Bayes risk equal to its maximum risk. The following corollary highlights an important special case of this strategy.

**Corollary 6.3 (TPE 5.1.5)** *If a Bayes estimator $\delta_A$ has constant risk (that is, $R(\theta, \delta_A) = R(\theta', \delta_\Lambda)$ for all $\theta$ and $\theta'$), then $\delta_\Lambda$ is minimax. Note that this is a sufficient but not necessary condition.*

It is often relatively easy to check whether an estimator has constant risk, and this is typically our first line of attack for determining whether an estimator is minimax. More generally we could find a prior support set $\omega$ such that $\Lambda(\omega) = 1$ and for which $R(\theta, \delta_\Lambda)$ is maximum for all $\theta \in \omega$.
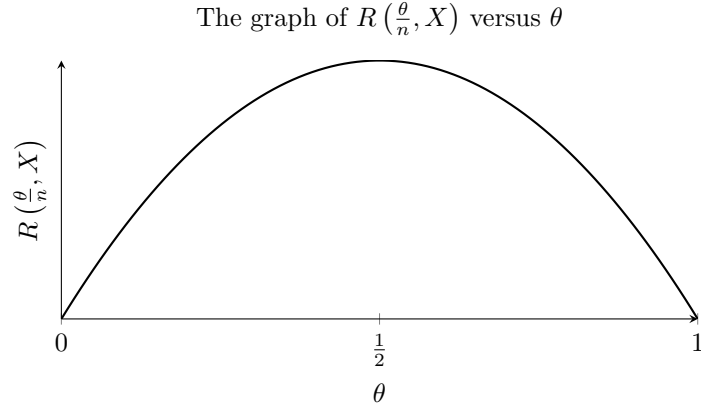
**Corollary 6.4 (TPE 5.1.6)** *Define*

$$\omega_\Lambda = \left\{ \theta : R(\theta, \delta_\Lambda) = \sup_{\theta'} R(\theta', \delta_\Lambda) \right\}$$

*Then, a Bayes estimator $\delta_\Lambda$ is minimax if $\Lambda(\omega_\Lambda) = 1$. (TPE misstates this result as if and only if, but the only if component is false. In other words, this condition is sufficient, but not necessary.)*

**Example 6.5** *Suppose $X \sim \text{Binom}(n, \theta)$ for some $\theta \in (0, 1)$ and that we use the squared error loss function. Is the sample proportion $\frac{X}{n}$ minimax? The risk of this estimator is*

$$R\left(\theta, \frac{X}{n}\right) = \frac{\theta(1 - \theta)}{n}.$$

*The graph of $R\left(\theta, \frac{X}{n}\right)$ versus $\theta$ looks like the following:*

The graph of $R\left(\frac{\theta}{n}, X\right)$ versus $\theta$



The risk has a unique maximum at $\theta = \frac{1}{2}$, so the worst-case risk is

$$\sup_{\theta \in \Omega} R\left(\theta, \frac{X}{n}\right) = R\left(\frac{1}{2}, \frac{X}{n}\right) = \frac{1}{4n}.$$

Unfortunately, we cannot apply TPE 5.1.6 directly because if $\Lambda\left(\left\{\frac{1}{2}\right\}\right) = 1$, then $\delta_\Lambda(X) = \frac{1}{2} \neq \frac{X}{n}$.

However, we can use the TPE 5.1.5 to find a minimax estimator and then compare the risk of the minimax estimator with that of $\frac{X}{n}$. To find a minimax estimator, we will search for a prior such that the Bayes estimator has constant risk.

Recall the following useful fact. Under the prior distribution Beta $(a, b)$, the Bayes estimator under the squared error loss is

$$\delta_{a,b}(X) = \frac{X + a}{n + a + b}.$$

For any $a$ and $b$,

$$R(\theta, \delta_{a,b}) = E_\theta\left[\left(\frac{X + a}{n + a + b} - \theta\right)^2\right]$$

$$= \frac{1}{(n + a + b)^2} E_\theta\left[(X + a - (n + a + b)\theta)^2\right]$$

$$= \frac{1}{(n + a + b)^2} E_\theta\left[(X - n\theta - a(\theta - 1) - \theta b)^2\right]$$

$$= \frac{1}{(n + a + b)^2}\left(n\theta(1 - \theta) + (a(\theta - 1) + \theta b)^2\right).$$

This is a quadratic function of $\theta$. To eliminate the $\theta$ dependence in $R(\theta, \delta_{a,b})$, we need the coefficients of the linear and quadratic terms to equal zero. The coefficient of $\theta^2$ is

$$-n + (a + b)^2$$

so we need $a + b = \sqrt{n}$ (since $a, b > 0$ ). The coefficient of $\theta$ is

$$n - 2a(a + b) = n - 2a\sqrt{n},$$

so we need $a = b = \frac{\sqrt{n}}{2}$. With these choices of $a$ and $b$, the risk of $R(\theta, \delta_{a,b})$ is constant, which implies that Beta $\left(\frac{\sqrt{n}}{2}, \frac{\sqrt{n}}{2}\right)$ is a least favorable prior with constant risk. Then our Bayes estimator

$$\delta_{\frac{\sqrt{\pi}}{2}, \frac{\sqrt{n}}{2}}(X) = \frac{X + \frac{\sqrt{n}}{2}}{n + \sqrt{n}},$$

is minimax with constant risk of

$$\frac{1}{4(\sqrt{n}+1)^2}.$$

Since the worst-case risk of $\frac{X}{n}$ is $\frac{1}{4n} > \frac{1}{4(\sqrt{n}+1)^2}$, we can conclude that $\frac{X}{n}$ is not minimax.

## 6.2 Minimaxity and least favorable prior sequences

In this part, we will extend our tools for deriving minimax estimators. Last time, we discovered that minimax estimators can arise from Bayes estimators under least favorable priors. However, it turns out that minimax estimators may not be Bayes estimators. Consider the following example, where our old approach fails.

**Example 6.6 (Minimax for i.i.d. Normal random variables with unknown mean $\theta$)** *Let $X_1, \ldots, X_n \overset{iid}{\sim} \mathcal{N}(\theta, \sigma^2)$, with $\sigma^2$ known. Our goal is to estimate $\theta$ under squared-error loss. For our first guess, pick the natural estimator $X$. Note that it has constant risk $\frac{\sigma^2}{n}$, which suggests minimaxity because we know that Bayes estimators with constant risk are also minimax estimators. However, $X$ is not Bayes for any prior, because under squared-error loss unbiased estimators are Bayes estimators only in the degenerate situations of zero risk (TPE Theorem 4.2.3), and $X$ is unbiased. Thus, we cannot conclude by our previous results (e.g., TPE Corollary 5.1.5) that $\bar{X}$ is minimax.*

*We might try to consider the wider class of estimators $\delta_{a,\mu_0}(X) = aX + (1-a)\mu_0$ for $a \in (0,1)$ and $\mu_0 \in R$, because many of the Bayes estimators we've encountered are convex combinations of a prior and a data mean. Note however that the worst case risk for these estimators is infinite:*

$$\sup_\theta E_\theta[\theta - \delta(X)]^2 = \sup_\theta \left\{a^2 \operatorname{Var}_\theta(\bar{X}) + (1-a)^2 (\theta - \mu_0)^2\right\}$$

$$= \frac{a^2\sigma^2}{n} + (1-a)^2 \sup_\theta (\theta - \mu_0)^2$$

$$= +\infty$$

*Since these estimators have poorer worst case risk than $\bar{X}$, they certainly cannot be minimax. We could keep trying to find Bayes estimators with better worst-case performance than $X$, but we would fail: it turns out that $\bar{X}$ is in fact minimax. To establish this, we will extend our minimax results to the limits of Bayes estimators, rather than restricting attention to Bayes estimators only.*

**Definition 6.7 (Least Favorable Sequence of Priors)** *Let $\{\Lambda_m\}$ be a sequence of priors with minimal average risk $\tau_{\Lambda_m} = \inf_\delta \int R(\theta, \delta) d\Lambda_m(\theta)$. Then, $\{\Lambda_m\}$ is a least favorable sequence of priors if there is a real number $r$ such that $r_{\Lambda_m} \to r < \infty$ and $r \geq r_{\Lambda'}$ for any prior $\Lambda'$.*

The reason for studying the limit of priors is that it may help us establish minimaxity. Since there need not exist a prior $\Lambda$ such that the associated Bayes estimator has average risk $r$, this definition is less restrictive than that of a least-favorable prior. We can prove an analogue of TPE Theorem 5.1.4 in this new setting.

**Theorem 6.8 (TPE 5.1.12)** *Suppose there is real number $r$ such that $\{\Lambda_m\}$ is a sequence of priors with $r_{\Lambda_m} \to r < \infty$. Let $\delta$ be any estimator such that $\sup_\theta R(\theta, \delta) = r$. Then,*

1. *$\delta$ is minimax,*

2. *$\{\Lambda_m\}$ is least-favorable.*

**Proof:** 1. Let $\delta'$ be any other estimator. Then, for any $m$,

$$\sup_\theta R\left(\theta, \delta'\right) \geq \int R\left(\theta, \delta'\right) d\Lambda_m(\theta) \geq r_{\Lambda_m},$$

so that sending $m \to \infty$ yields

$$\sup_\theta R\left(\theta, \delta'\right) \geq r = \sup_\theta R(\theta, \delta),$$

which means that $\delta$ is minimax.
2. Let $\Lambda'$ be any prior, then

$$r_{\Lambda'} = \int R\left(\theta, \delta_{\Lambda'}\right) d\Lambda'(\theta) \leq \int R(\theta, \delta) d\Lambda'(\theta) \leq \sup_\theta R(\theta, \delta) = r,$$

which means that $\{\Lambda_m\}$ is least favorable. ∎

Unlike Theorem 5.1.4, this result does not guarantee uniqueness, even if the Bayes estimators $\delta_{\Lambda_m}$ are unique. This is because the limiting step in the proof of (1) changes any strict inequality to nonstrict inequality. However, this result allows to check much wider class of estimators, since to check that the estimator is indeed a minimax estimator we need to find only the sequence of Bayes risks convergent to maximum risk of our candidate.

**Example 6.9 (Minimax for i.i.d. Normal random variables, continued)** *We now have the tools to confirm our suspicion that $X$ is minimax. By Theorem TPE 5.1.12 above, it suffices to find a sequence $\{\Lambda_m\}$ such that $r_{\Lambda_m} \to \frac{\sigma^2}{n} =: r$. Using the conjugate prior is a good starting point, so we let $\{\Lambda_m\}$ be the conjugate priors $\left\{N\left(0, m^2\right)\right\}$ with variance tending to $\infty$, so that $\Lambda_m$ tends to the (improper with $\pi(\theta) = 1, \forall \theta \in R$ ) uniform prior on $R$. By TPE Example 4.2.2, the posterior for $\theta$ associated with each $\Lambda_m$ is*

$$\theta \mid X_1, \ldots, X_n \sim N\left(\frac{\frac{nX}{\sigma^2}}{\frac{n}{\sigma^2} + \frac{1}{m^2}}, \frac{1}{\frac{n}{\sigma^2} + \frac{1}{m^2}}\right).$$

*In particular, the posterior variance does not depend on $X_1, \ldots, X_n$, so Lemma 1 below automatically yields the Bayes risk*

$$r_{\Lambda_m} = \frac{1}{\frac{n}{\sigma^2} + \frac{1}{m^2}} \xrightarrow{m \to \infty} \frac{\sigma^2}{n} = \sup_\theta R(\theta, X).$$

*It follows from Theorem TPE 5.1.12 that $X$ is minimax and $\{\Lambda_m\}$ is least favorable.*

**Lemma 6.10 (TPE 5.1.13)** *If the posterior variance $\mathrm{Var}_{\Theta|X}(g(\Theta) \mid X = x$ ) is constant in $x$, then under squared error loss, $r_\Lambda = \mathrm{Var}_{\Theta|X}(g(\Theta) \mid X = x)$.*

We know that the posterior mean minimizes Bayes risk, so this result can be obtained by plugging in the posterior mean of $g(\theta)$ into the average risk.

## 6.3   Minimaxity via submodel restriction

The following example illustrates the technique of deriving a minimax estimator for a general family of models by restricting attention to a subset of that family. The idea comes from simple observation that if the estimator is minimax in submodel and its risk doesn't change when we go to a larger model then estimator is minimax in this larger class.

**Example 6.11 (Minimax for i.i.d. Normal random variables, unknown mean and variance)** *Reconsider Example 1 in the case that the variance is unknown. That is, let $X_1, \ldots, X_n \overset{iid}{\sim} \mathcal{N}\left(\theta, \sigma^2\right)$, with both $\theta$ and $\sigma^2$ unknown. Note that*

$$\sup_{\theta, \sigma^2} R\left(\left(\theta, \sigma^2\right), \bar{X}\right) = \sup_{\sigma^2} \frac{\sigma^2}{n} = \infty,$$

*and in fact, the maximum risk of any estimator in this setting is infinite, so the question of minimaxity is uninteresting. Therefore, we restrict attention to the family parameterized by $\Omega = \left\{\left(\theta, \sigma^2\right) : \theta \in R, \sigma^2 \leq B\right\}$, where $B$ is a known constant. Assume $\delta$ is any other estimator. Calculating the risk of $\bar{X}$ within this family, we find*

$$\sup_{\theta \in R, \sigma^2 \leq B} R\left(\left(\theta, \sigma^2\right), \bar{X}\right) = \frac{B}{n}$$
$$= \sup_{\theta \in R, \sigma^2 = B} R\left(\left(\theta, \sigma^2\right), \bar{X}\right)$$
$$\leq \sup_{\theta \in R, \sigma^2 = B} R\left(\left(\theta, \sigma^2\right), \delta\right) \text{ [submodel minimax]}$$
$$\leq \sup_{\theta \in R, \sigma^2 \leq B} R\left(\left(\theta, \sigma^2\right), \delta\right),$$

*where the first inequality follows from the fact that $\bar{X}$ is minimax for i.i.d. normals with known $\sigma^2$, and the second inequality follows from the fact that we are taking the supremum over a larger set. Hence, we are able to show that $\bar{X}$ is minimax over $\Omega$ by focusing on the case where $\sigma^2$ is known. Notice further that the form of the estimator does not depend on the upper bound $B$, though the bound is necessary for minimaxity to be worth investigating.*