

Chapter 3

Risk Minimization and Rademacher Complexity II

Ernest K. Ryu
Seoul National University

Mathematical Machine Learning Theory
Spring 2024

Outline

Calibration

$$\mathcal{R}_{\Phi^{0-1}} - \mathcal{R}_{\Phi^{0-1}}^* \stackrel{?}{\leq} H(\mathcal{R}_{\Phi} - \mathcal{R}_{\Phi}^*)$$

Calibration function with convex losses

Ridge least squares regression

Binary classification

Consider the binary classification problem, where $\tilde{\mathcal{Y}} = \mathcal{Y} = \{-1, +1\}$ and $\ell(y', y) = \mathbf{1}_{\{y' \neq y\}}$. So

$$\mathcal{R}[f] = \mathbb{E}_{(X,Y) \sim P} [\ell(f(X), Y)] = \mathbb{P}_{(X,Y) \sim P} (f(X) \neq Y).$$

Define

$$\eta(X) = \mathbb{P}(Y = +1 | X).$$

Assume $\eta(X) \neq 1/2$ with probability 1. Then,

$$f^*(X) = \begin{cases} -1 & \text{if } \eta(X) < 1/2 \\ +1 & \text{if } \eta(X) > 1/2 \end{cases}$$

is a Bayes predictor, and

$$\mathcal{R}^* = \mathbb{E}_{X \sim P_X} [\min\{1 - \eta(X), \eta(X)\}].$$

Surrogate loss

We replace $\Phi^{0-1}(u)$ with a surrogate loss such as

$$\Phi^{\text{hinge}}(u) = \max\{1 - u, 0\}$$

$$\Phi^{\text{logistic}}(u) = \log(1 + e^{-u})$$

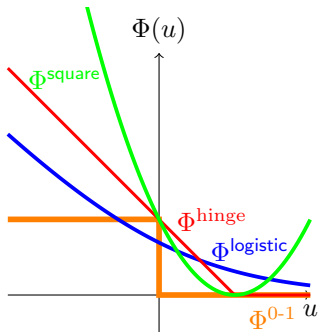
$$\Phi^{\text{square}}(u) = (1 - u)^2,$$

which are nice continuous, convex functions, and solve the continuous convex optimization problem

$$\underset{g}{\text{minimize}} \quad \underbrace{\mathbb{E}_{(X,Y) \sim P} [\Phi(Yg(X))]}_{=\mathcal{R}_{\Phi}[g]}$$

or its approximation

$$\underset{g}{\text{minimize}} \quad \underbrace{\frac{1}{N} \sum_{i=1}^N \Phi(Y_i g(X_i))}_{=\hat{\mathcal{R}}_{\Phi}[g]}$$



Binary classification with square loss

Consider the square surrogate loss

$$\mathcal{R}_{\Phi^{\text{square}}}[g] = \mathbb{E}_{(X,Y) \sim P} [(1 - Yg(X))^2] = \mathbb{E}_{(X,Y) \sim P} [(g(X) - Y)^2].$$

Bayes predictor has a simple analytic form:

$$\begin{aligned} g^*(X) &= \mathbb{E}[Y | X] = -1 \cdot \mathbb{P}(Y = -1 | X) + 1 \cdot \mathbb{P}(Y = +1 | X) \\ &= 2\eta(X) - 1. \end{aligned}$$

Also,

$$\begin{aligned} \mathcal{R}_{\Phi^{\text{square}}}[g] - \mathcal{R}_{\Phi^{\text{square}}}[g^*] &= \mathbb{E}_{(X,Y) \sim P} [(g(X) - Y)^2 - (g^*(X) - Y)^2] \\ &= \mathbb{E}_{(X,Y) \sim P} [g(X)^2 - 2g(X)Y - g^*(X)^2 + 2g^*(X)Y] \\ &= \mathbb{E}_X \left[\mathbb{E}_Y [g(X)^2 - 2g(X)Y - g^*(X)^2 + 2g^*(X)Y | X] \right] \\ &= \mathbb{E}_X [g(X)^2 - 2g(X) \mathbb{E}_Y [Y | X] - g^*(X)^2 + 2g^*(X) \mathbb{E}_Y [Y | X]] \\ &= \mathbb{E}_X [g(X)^2 - 2g(X) + g^*(X)^2] \\ \text{Calibration} \quad &= \mathbb{E}_X [(g(X) - g^*(X))^2]. \end{aligned}$$

Minimize surrogate loss $\stackrel{?}{\Rightarrow}$ Minimize original loss

However, we should not forget that we have changed the optimization problem from minimizing $\mathcal{R}_{\Phi^{0-1}}$ to \mathcal{R}_{Φ} .

Is this valid? Does the following implication hold?

$$\mathcal{R}_{\Phi}[g] - \mathcal{R}_{\Phi}^* = 0 \quad \stackrel{?}{\Rightarrow} \quad \mathcal{R}_{\Phi^{0-1}}[g] - \mathcal{R}_{\Phi^{0-1}}^* = 0$$

In general, no.

Since $\Phi^{0-1} \leq \gamma\Phi$ for some $\gamma > 0$, if $\mathcal{R}_{\Phi}^* = 0$, then $\mathcal{R}_{\Phi^{0-1}}^* = 0$ and

$$\mathcal{R}_{\Phi}[g] = 0 \quad \Rightarrow \quad \mathcal{R}_{\Phi^{0-1}}[g] = 0.$$

However, if $\mathcal{R}_{\Phi}^* > 0$, the desired implication does not hold in general.

When is minimizing \mathcal{R}_Φ valid?

We shall now study conditions that ensure:

$$\operatorname{argmin}_g \mathcal{R}_\Phi[g] \subseteq \operatorname{argmin}_g \mathcal{R}_{\Phi^{0-1}}[g].$$

If so, then (exactly) minimizing \mathcal{R}_Φ provides a minimizer to $\mathcal{R}_{\Phi^{0-1}}$, the actual risk that we care about, i.e.,

$$\mathcal{R}_\Phi[g] - \mathcal{R}_\Phi^* = 0 \quad \Rightarrow \quad \mathcal{R}_{\Phi^{0-1}}[g] - \mathcal{R}_{\Phi^{0-1}}^* = 0.$$

Conditional Φ -risk

For any $g: \mathcal{X} \rightarrow \mathbb{R}$, define the *conditional Φ -risk* as

$$\begin{aligned}\mathcal{R}_\Phi[g | X] &= \mathbb{E}_{Y \sim P_{Y|X}} [\Phi(Yg(X)) | X] \\ &= \eta(X)\Phi(g(X)) + (1 - \eta(X))\Phi(-g(X)).\end{aligned}$$

(Of course, $\mathbb{E}_X[\mathcal{R}_\Phi[g | X]] = \mathcal{R}_\Phi[g]$.)

Let

$$C_\Phi(\alpha; \eta) = \eta\Phi(\alpha) + (1 - \eta)\Phi(-\alpha).$$

Then,

$$\mathcal{R}_\Phi[g | X] = C_\Phi(g(X); \eta(X)).$$

Bayes predictor from conditional Φ -risk

Recall that the Bayes predictor was obtained by

$$g_{\Phi}^*(X) \in \operatorname{argmin}_{\alpha \in \mathbb{R}} \mathbb{E}_{Y \sim P_{Y|X}} [\Phi(Y\alpha) | X] = \operatorname{argmin}_{\alpha \in \mathbb{R}} C_{\Phi}(\alpha; \eta(X)).$$

For the true 0-1 loss, we have

$$\begin{aligned} \operatorname{argmin}_{\alpha \in \mathbb{R}} C_{\Phi^{0-1}}(\alpha; \eta(X)) &= \operatorname{argmin}_{\alpha \in \mathbb{R}} \{ \eta(X) \mathbf{1}_{\{\alpha \leq 0\}} + (1 - \eta(X)) \mathbf{1}_{\{\alpha \geq 0\}} \} \\ &= \begin{cases} \alpha > 0 & \text{if } \eta(X) > 1/2 \\ \alpha < 0 & \text{if } \eta(X) < 1/2. \end{cases} \end{aligned}$$

(For simplicity, assume $\eta(X) \neq 1/2$ with probability 1.) I.e., it is optimal to output $\alpha > 0$ if $Y = +1$ is more likely and $\alpha < 0$ if $Y = -1$ is more likely. Does this hold for the surrogate loss function?

Calibrated surrogate loss

We say a surrogate loss Φ is *classification calibrated* or *calibrated* if

$$\operatorname{argmin}_{\alpha \in \mathbb{R}} C_{\Phi}(\alpha; \eta(X)) \subseteq \operatorname{argmin}_{\alpha \in \mathbb{R}} C_{\Phi^{0-1}}(\alpha; \eta(X)) = \begin{cases} \alpha > 0 & \text{if } \eta(X) > 1/2 \\ \alpha < 0 & \text{if } \eta(X) < 1/2. \end{cases}$$

Lemma

Let Φ be classification calibrated. Then,

$$\operatorname{argmin}_g \mathcal{R}_{\Phi}[g] \subseteq \operatorname{argmin}_g \mathcal{R}_{\Phi^{0-1}}[g].$$

Proof. Let $g_{\Phi}^* \in \operatorname{argmin}_g \mathcal{R}_{\Phi}[g]$. Then,

$$g_{\Phi}^*(X) \in \operatorname{argmin}_{\alpha \in \mathbb{R}} C_{\Phi}(\alpha; \eta(X))$$

for P -almost all X . Then,

$$g_{\Phi}^*(X) \in \operatorname{argmin}_{\alpha \in \mathbb{R}} C_{\Phi^{0-1}}(\alpha; \eta(X))$$

for P -almost all X , and we conclude

$$g_{\Phi}^* \in \operatorname{argmin}_g \mathcal{R}_{\Phi^{0-1}}[g].$$



Bayes predictor for square loss is optimal for 0-1 loss

Recall that

$$g_{\Phi^{\text{square}}}^*(X) = 2\eta(X) - 1.$$

Since $g_{\Phi^{\text{square}}}^*(X) > 0$ if $\eta(X) > 1/2$ and vice versa,

$$g_{\Phi^{\text{square}}}^*(X) \in \underset{\alpha \in \mathbb{R}}{\operatorname{argmin}} C_{\Phi^{0-1}}(\alpha; \eta(X)).$$

Therefore,

$$g_{\Phi^{\text{square}}}^* \in \underset{g}{\operatorname{argmin}} \mathcal{R}_{\Phi^{0-1}}[g].$$

How about

$$g_{\Phi^{\text{logistic}}}^* \stackrel{?}{\in} \underset{g}{\operatorname{argmin}} \mathcal{R}_{\Phi^{0-1}}, \quad g_{\Phi^{\text{hinge}}}^* \stackrel{?}{\in} \underset{g}{\operatorname{argmin}} \mathcal{R}_{\Phi^{0-1}}$$

Calibrated surrogate loss

Theorem

Let $\Phi: \mathbb{R} \rightarrow \mathbb{R}$ be convex. If Φ is differentiable at 0 and $\Phi'(0) < 0$, then Φ is classification-calibrated.

Proof. Convexity of Φ implies $C_{\Phi}(\alpha; \eta)$ is convex in α for any fixed $\eta \in [0, 1]$. If $\eta > 1/2$, then

$$\frac{d}{d\alpha} C_{\Phi}(\alpha; \eta) \Big|_{\alpha=0} = \eta\Phi'(0) - (1 - \eta)\Phi'(0) < 0.$$

Therefore, $\operatorname{argmin}_{\alpha \in \mathbb{R}} C_{\Phi}(\alpha; \eta) \subseteq (0, \infty)$ by convexity.

If $\eta < 1/2$, then

$$\frac{d}{d\alpha} C_{\Phi}(\alpha; \eta) \Big|_{\alpha=0} = \eta\Phi'(0) - (1 - \eta)\Phi'(0) > 0.$$

Therefore, $\operatorname{argmin}_{\alpha \in \mathbb{R}} C_{\Phi}(\alpha; \eta) \subseteq (-\infty, 0)$ by convexity. □

Calibrated surrogate loss

Therefore, all three surrogate losses are calibrated, and

$$g_{\Phi^{\text{logistic}}}^* \in \underset{g}{\operatorname{argmin}} \mathcal{R}_{\Phi^{0-1}}[g]$$

$$g_{\Phi^{\text{hinge}}}^* \in \underset{g}{\operatorname{argmin}} \mathcal{R}_{\Phi^{0-1}}[g]$$

$$g_{\Phi^{\text{square}}}^* \in \underset{g}{\operatorname{argmin}} \mathcal{R}_{\Phi^{0-1}}[g].$$

Outline

Calibration

$$\mathcal{R}_{\Phi^{0-1}} - \mathcal{R}_{\Phi^{0-1}}^* \stackrel{?}{\leq} H(\mathcal{R}_{\Phi} - \mathcal{R}_{\Phi}^*)$$

Calibration function with convex losses

Ridge least squares regression

$$\mathcal{R}_{\Phi^{0-1}} - \mathcal{R}_{\Phi^{0-1}}^* \stackrel{?}{\leq} H(\mathcal{R}_{\Phi} - \mathcal{R}_{\Phi}^*)$$

When is approximately minimizing \mathcal{R}_Φ valid?

If Φ is calibrated, then

$$\mathcal{R}_\Phi[g] - \mathcal{R}_\Phi^* = 0 \quad \Rightarrow \quad \mathcal{R}_{\Phi^{0-1}}[g] - \mathcal{R}_{\Phi^{0-1}}^* = 0.$$

However, do we have?

$$\mathcal{R}_\Phi[g] - \mathcal{R}_\Phi^* < \text{small} \quad \Rightarrow \quad \mathcal{R}_{\Phi^{0-1}}[g] - \mathcal{R}_{\Phi^{0-1}}^* < \text{small}$$

After all, we can only hope to approximately minimize \mathcal{R}_Φ .

$$\mathcal{R}_{\Phi^{0-1}} - \mathcal{R}_{\Phi^{0-1}}^* \stackrel{?}{\leq} H(\mathcal{R}_\Phi - \mathcal{R}_\Phi^*)$$

$$\mathcal{R}_{\Phi^{0-1}}[g] - \mathcal{R}_{\Phi^{0-1}}^* \leq \mathcal{R}_{\Phi^{\text{hinge}}}[g] - \mathcal{R}_{\Phi^{\text{hinge}}}^*$$

For the hinge loss, we can carry out the analysis with direct arguments.

Recall,

$$\begin{aligned} C_{\Phi^{0-1}}(\alpha; \eta) &= \eta \mathbf{1}_{\{\alpha \leq 0\}} + (1 - \eta) \mathbf{1}_{\{\alpha \geq 0\}} \\ C_{\Phi^{\text{hinge}}}(\alpha; \eta) &= \eta(1 - \alpha)_+ + (1 - \eta)(1 + \alpha)_+. \end{aligned}$$

With direct calculations, we get

$$\inf_{\alpha \in \mathbb{R}} C_{\Phi^{0-1}}(\alpha; \eta) = \min\{\eta, 1 - \eta\}, \quad \inf_{\alpha \in \mathbb{R}} C_{\Phi^{\text{hinge}}}(\alpha; \eta) = 2 \min\{\eta, 1 - \eta\}.$$

With direct (albeit tedious) arguments, we can show

$$C_{\Phi^{0-1}}(\alpha; \eta) - \inf_{\alpha \in \mathbb{R}} C_{\Phi^{0-1}}(\alpha; \eta) \leq C_{\Phi^{\text{hinge}}}(\alpha; \eta) - \inf_{\alpha \in \mathbb{R}} C_{\Phi^{\text{hinge}}}(\alpha; \eta)$$

for all $\alpha \in \mathbb{R}$ and $\eta \in [0, 1]$, which implies

$$\mathcal{R}_{\Phi^{0-1}}[g] - \mathcal{R}_{\Phi^{0-1}}^* \leq \mathcal{R}_{\Phi^{\text{hinge}}}[g] - \mathcal{R}_{\Phi^{\text{hinge}}}^*.$$

$$\mathcal{R}_{\Phi^{0-1}} - \mathcal{R}_{\Phi^{0-1}}^* \stackrel{?}{\leq} H(\mathcal{R}_{\Phi} - \mathcal{R}_{\Phi}^*)$$

$$\mathcal{R}_{\Phi^{0-1}}[g] - \mathcal{R}_{\Phi^{0-1}}^* \not\leq \mathcal{R}_{\Phi^{\text{logistic}}}[g] - \mathcal{R}_{\Phi^{\text{logistic}}}^*$$

For the logistic loss, we have

$$C_{\Phi^{0-1}}(\alpha; \eta) \leq \frac{1}{\log 2} C_{\Phi^{\text{logistic}}}(\alpha; \eta)$$

However,

$$C_{\Phi^{0-1}}(\alpha; \eta) - \inf_{\alpha \in \mathbb{R}} C_{\Phi^{0-1}}(\alpha; \eta) \not\leq \gamma (C_{\Phi^{\text{logistic}}}(\alpha; \eta) - \inf_{\alpha \in \mathbb{R}} C_{\Phi^{\text{logistic}}}(\alpha; \eta))$$

for any constant $\gamma > 0$, and we cannot proceed with the same argument.

The same problem arises with the square loss.

$$\mathcal{R}_{\Phi^{0-1}} - \mathcal{R}_{\Phi^{0-1}}^* \stackrel{?}{\leq} H(\mathcal{R}_{\Phi} - \mathcal{R}_{\Phi}^*)$$

Lemma

Let $g^* \in \operatorname{argmin}_g \mathcal{R}_{\Phi^{0-1}}[g]$. Then,

$$\begin{aligned}\mathcal{R}_{\Phi^{0-1}}[g] - \mathcal{R}_{\Phi^{0-1}}[g^*] &= \mathbb{E}[\mathbf{1}_{\{g(X)g^*(X) < 0\}} |2\eta(X) - 1|] \\ &\leq \mathbb{E}[\mathbf{1}_{\{g(X)g^*(X) < 0\}} |2\eta(X) - 1 - b(g(X))|]\end{aligned}$$

for any $b: \mathbb{R} \rightarrow \mathbb{R}$ such that $\operatorname{sign}(b(x)) = \operatorname{sign}(x)$ for all $x \in \mathbb{R}$.

$$\mathcal{R}_{\Phi^{0-1}} - \mathcal{R}_{\Phi^{0-1}}^* \stackrel{?}{\leq} H(\mathcal{R}_{\Phi} - \mathcal{R}_{\Phi}^*)$$

Proof. The first claim follows from

$$\begin{aligned}
& \mathcal{R}_{\Phi^{0-1}}[g] - \mathcal{R}_{\Phi^{0-1}}[g^*] \\
&= \mathbb{E} \left[\mathbb{E} \left[\mathbf{1}_{\{\text{sign}(g(X)) \neq Y\}} - \mathbf{1}_{\{\text{sign}(g^*(X)) \neq Y\}} \mid X \right] \right] \\
&= \mathbb{E} \left[\mathbb{E} \left[-\mathbf{1}_{\{g(X) > 0, g^*(X) < 0\}} \mathbf{1}_{\{Y = +1\}} + \mathbf{1}_{\{g(X) > 0, g^*(X) < 0\}} \mathbf{1}_{\{Y = -1\}} \right. \right. \\
&\quad \left. \left. + \mathbf{1}_{\{g(X) < 0, g^*(X) > 0\}} \mathbf{1}_{\{Y = +1\}} - \mathbf{1}_{\{g(X) < 0, g^*(X) > 0\}} \mathbf{1}_{\{Y = -1\}} \mid X \right] \right] \\
&= \mathbb{E} \left[-\mathbf{1}_{\{g(X) > 0, g^*(X) < 0\}} \eta(X) + \mathbf{1}_{\{g(X) > 0, g^*(X) < 0\}} (1 - \eta(X)) \right. \\
&\quad \left. + \mathbf{1}_{\{g(X) < 0, g^*(X) > 0\}} \eta(X) - \mathbf{1}_{\{g(X) < 0, g^*(X) > 0\}} (1 - \eta(X)) \right] \\
&= \mathbb{E} \left[\mathbf{1}_{\{g(X) > 0, g^*(X) < 0\}} (1 - 2\eta(X)) - \mathbf{1}_{\{g(X) < 0, g^*(X) > 0\}} (1 - 2\eta(X)) \right] \\
&= \mathbb{E} \left[\mathbf{1}_{\{g(X) > 0, g^*(X) < 0\}} |1 - 2\eta(X)| + \mathbf{1}_{\{g(X) < 0, g^*(X) > 0\}} |1 - 2\eta(X)| \right] \\
&= \mathbb{E} \left[\mathbf{1}_{\{g(X)g^*(X) < 0\}} |1 - 2\eta(X)| \right],
\end{aligned}$$

where we use the fact that $g^*(X) < 0$ implies $\eta(X) < 1/2$.

$$\mathcal{R}_{\Phi^{0-1}} - \mathcal{R}_{\Phi^{0-1}}^* \stackrel{?}{\leq} H(\mathcal{R}_{\Phi} - \mathcal{R}_{\Phi}^*)$$

For the second claim,

$$\begin{aligned}
& \mathbb{E}[\mathbf{1}_{\{g(X)g^*(X) < 0\}} |2\eta(X) - 1|] \\
&= \mathbb{E}[\mathbf{1}_{\{g(X)g^*(X) < 0, g^*(X) > 0, \eta(X) > 1/2\}} (2\eta(X) - 1)] \\
&\quad + \mathbb{E}[\mathbf{1}_{\{g(X)g^*(X) < 0, g^*(X) < 0, \eta(X) < 1/2\}} (-2\eta(X) + 1)] \\
&\leq \mathbb{E}[\mathbf{1}_{\{g(X)g^*(X) < 0, g^*(X) > 0, \eta(X) > 1/2\}} (2\eta(X) - 1 - b(g(X)))] \\
&\quad + \mathbb{E}[\mathbf{1}_{\{g(X)g^*(X) < 0, g^*(X) < 0, \eta(X) < 1/2\}} (-2\eta(X) + 1 + b(g(X)))] \\
&= \mathbb{E}[\mathbf{1}_{\{g(X)g^*(X) < 0, g^*(X) > 0, \eta(X) > 1/2\}} |2\eta(X) - 1 - b(g(X))|] \\
&\quad + \mathbb{E}[\mathbf{1}_{\{g(X)g^*(X) < 0, g^*(X) < 0, \eta(X) < 1/2\}} |2\eta(X) - 1 - b(g(X))|] \\
&= \mathbb{E}[\mathbf{1}_{\{g(X)g^*(X) < 0\}} |2\eta(X) - 1 - b(g(X))|].
\end{aligned}$$

□

$$\mathcal{R}_{\Phi^{0-1}} - \mathcal{R}_{\Phi^{0-1}}^* \stackrel{?}{\leq} H(\mathcal{R}_{\Phi} - \mathcal{R}_{\Phi}^*)$$

Square loss

Equipped with this lemma, we can now analyze the relationship between $\mathcal{R}_{\Phi^{0-1}}[g] - \mathcal{R}_{\Phi^{0-1}}^*[g^*]$ and $\mathcal{R}_{\Phi^{\text{square}}}[g] - \mathcal{R}_{\Phi^{\text{square}}}^*[g^*]$:

$$\begin{aligned}\mathcal{R}_{\Phi^{0-1}}[g] - \mathcal{R}_{\Phi^{0-1}}[g^*] &\leq \mathbb{E}[\mathbf{1}_{\{g(X)g^*(X) < 0\}} |2\eta(X) - 1 - g(X)|] \\ &\leq \left(\mathbb{E}[\mathbf{1}_{\{g(X)g^*(X) < 0\}} \underbrace{|2\eta(X) - 1 - g(X)|^2}_{=g^*(X)}] \right)^{1/2} \\ &\leq \left(\mathbb{E}[|g^*(X) - g(X)|^2] \right)^{1/2} \\ &= \left(\mathcal{R}_{\Phi^{\text{square}}}[g] - \mathcal{R}_{\Phi^{\text{square}}}[g^*] \right)^{1/2},\end{aligned}$$

where the second inequality follows from Jensen.

Therefore,

$$\mathcal{R}_{\Phi^{\text{square}}} - \mathcal{R}_{\Phi^{\text{square}}}^* < \text{small} \quad \Rightarrow \quad \mathcal{R}_{\Phi^{0-1}} - \mathcal{R}_{\Phi^{0-1}}^* < \sqrt{\text{small}}.$$

$$\mathcal{R}_{\Phi^{0-1}} - \mathcal{R}_{\Phi^{0-1}}^* \stackrel{?}{\leq} H(\mathcal{R}_{\Phi} - \mathcal{R}_{\Phi}^*)$$

Logistic loss

Lemma

For any $x, u \in \mathbb{R}$

$$\log(e^{-x/2} + e^{x/2}) - ux - \inf_{x \in \mathbb{R}} \{\log(e^{-x/2} + e^{x/2}) - ux\} \geq 2 \left(u - \frac{e^x - 1}{2(e^x + 1)} \right)^2.$$

Proof. A brute-force proof:

$$\begin{aligned} \inf_x \{\log(e^{-x/2} + e^{x/2}) - ux\} &= \begin{cases} \frac{1}{2}(1 - 2u) \log \frac{1+2u}{1-2u} + \log \frac{2}{1+2u} & \text{if } -2 < u < 2 \\ -\infty & \text{otherwise.} \end{cases} \\ &\leq \log(e^{-x/2} + e^{x/2}) - ux - 2 \left(u - \frac{e^x - 1}{2(e^x + 1)} \right)^2 \end{aligned}$$

with a Taylor expansion argument. Rearrange the inequality to conclude the statement. (Better proof later.) \square

$$\mathcal{R}_{\Phi^{0-1}} - \mathcal{R}_{\Phi^{0-1}}^* \stackrel{?}{\leq} H(\mathcal{R}_{\Phi} - \mathcal{R}_{\Phi}^*)$$

Logistic loss

Recall that

$$\Phi^{\text{logistic}}(u) = \log(1 + e^{-u}).$$

Then,

$$\begin{aligned} C_{\Phi^{\text{logistic}}}(\alpha; \eta) &= \eta \log(1 + e^{-\alpha}) + (1 - \eta) \log(1 + e^{\alpha}) \\ &= \log(e^{-\alpha/2} + e^{\alpha/2}) - \frac{2\eta - 1}{2} \alpha \end{aligned}$$

Appealing to the previous lemma, we have

$$C_{\Phi^{\text{logistic}}}(\alpha; \eta) - \inf_{\alpha \in \mathbb{R}} C_{\Phi^{\text{logistic}}}(\alpha; \eta) \geq \frac{1}{2} \left(2\eta - 1 - \frac{e^{\alpha} - 1}{e^{\alpha} + 1} \right)^2.$$

$$\mathcal{R}_{\Phi^{0-1}} - \mathcal{R}_{\Phi^{0-1}}^* \stackrel{?}{\leq} H(\mathcal{R}_{\Phi} - \mathcal{R}_{\Phi}^*)$$

Logistic loss

Plug in $\alpha \leftarrow g(X)$ and $\eta \leftarrow \eta(X)$, and take the expectation to get

$$\begin{aligned}\mathcal{R}_{\Phi^{\text{logistic}}} [g] - \mathcal{R}_{\Phi^{\text{logistic}}}^* &\geq \frac{1}{2} \mathbb{E} \left[\left(2\eta(X) - 1 - \frac{e^{g(X)} - 1}{e^{g(X)} + 1} \right)^2 \right] \\ &\geq \frac{1}{2} \left(\mathbb{E} \left[\left| 2\eta(X) - 1 - \frac{e^{g(X)} - 1}{e^{g(X)} + 1} \right| \right] \right)^2 \\ &\geq \frac{1}{2} \left(\mathcal{R}_{\Phi^{0-1}} [g] - \mathcal{R}_{\Phi^{0-1}}^* \right)^2.\end{aligned}$$

Therefore, we conclude

$$\mathcal{R}_{\Phi^{0-1}} [g] - \mathcal{R}_{\Phi^{0-1}}^* \leq \sqrt{2} \left(\mathcal{R}_{\Phi^{\text{logistic}}} [g] - \mathcal{R}_{\Phi^{\text{logistic}}}^* \right)^{1/2}$$

the same (up to constant) guarantee as for the square loss.

$$\mathcal{R}_{\Phi^{0-1}} - \mathcal{R}_{\Phi^{0-1}}^* \stackrel{?}{\leq} H(\mathcal{R}_{\Phi} - \mathcal{R}_{\Phi}^*)$$

Calibration function

We established guarantees of the form

$$\mathcal{R}_{\Phi^{0-1}} - \mathcal{R}_{\Phi^{0-1}}^* \stackrel{?}{\leq} H(\mathcal{R}_{\Phi} - \mathcal{R}_{\Phi}^*),$$

where H is a monotonically increasing function. H is called the *calibration function*.

The guarantee for the hinge loss is better than the guarantee for the square or logistic loss. However, we will later see that the hinge loss is harder to optimize due to its non-differentiability. So there is a trade-off.

$$\mathcal{R}_{\Phi^{0-1}} - \mathcal{R}_{\Phi^{0-1}}^* \stackrel{?}{\leq} H(\mathcal{R}_{\Phi} - \mathcal{R}_{\Phi}^*)$$

Impact on approximation errors

So far, our analysis was carried out without any restriction on the set of functions.

In practice, however, we use a restricted function class \mathcal{F} (often with a controlled Rademacher complexity). The choice of the surrogate loss Φ affects the Bayes predictor (even though the set of Bayes predictor for Φ^{0-1} is always the same), so the approximation error is affected by the choice of Φ .

In particular,

$$\begin{aligned}g_{\Phi^{\text{hinge}}}^*(X) &= \text{sign}(2\eta(X) - 1) \\g_{\Phi^{\text{logistic}}}^*(X) &= \text{atanh}(2\eta(X) - 1) \\g_{\Phi^{\text{square}}}^*(X) &= 2\eta(X) - 1.\end{aligned}$$

If Φ admits a g_{Φ}^* that is well approximated by \mathcal{F} , that is a reason to favor Φ . (Having a favorable calibration function and the ease of optimization are two other reasons to favor a choice of Φ .)

$$\mathcal{R}_{\Phi^{0-1}} - \mathcal{R}_{\Phi^{0-1}}^* \stackrel{?}{\leq} H(\mathcal{R}_{\Phi} - \mathcal{R}_{\Phi}^*)$$

Outline

Calibration

$$\mathcal{R}_{\Phi^{0-1}} - \mathcal{R}_{\Phi^{0-1}}^* \stackrel{?}{\leq} H(\mathcal{R}_{\Phi} - \mathcal{R}_{\Phi}^*)$$

Calibration function with convex losses

Ridge least squares regression

Logistic loss

Lemma

For any $x, u \in \mathbb{R}$

$$\log(e^{-x/2} + e^{x/2}) - ux - \inf_{x \in \mathbb{R}} \{\log(e^{-x/2} + e^{x/2}) - ux\} \geq 2(u - b(x))^2,$$

where $b: \mathbb{R} \rightarrow \mathbb{R}$ is a sign-preserving function.

Better Proof. Note that

$$f(x) = \log(e^{-x/2} + e^{x/2})$$

is a convex L -smooth function with $L = 1/4$. (Easy to check that $0 \leq f''(x) \leq 1/4$ for all $x \in \mathbb{R}$.) Then, by the Fenchel–Young inequality for smooth convex functions, we have

$$f(x) + f^*(u) - ux \geq \frac{1}{2L} \|y - f'(x)\|^2.$$

Finally, it is straightforward to verify $f'(0) = 0$ and f' is strictly increasing. □

Calibration functions for square and logistic losses

Assume

$$\Phi(u) = a(u) - \gamma u + \beta,$$

where $a(0) = 0$, a is convex L -smooth, a is even, $\gamma > 0$, and $\beta \in \mathbb{R}$.

Recall

$$\Phi^{\text{square}}(u) = (1 - u)^2 = u^2 - 2u + 1 \quad (2\text{-smooth})$$

$$\Phi^{\text{logistic}}(u) = \log(1 + e^{-u}) = \log(e^{-u/2} + e^{u/2}) - \frac{1}{2}u \quad (\frac{1}{4}\text{-smooth})$$

Then,

$$\begin{aligned} C_{\Phi}(\alpha; \eta) &= \eta\Phi(\alpha) + (1 - \eta)\Phi(-\alpha) + \beta \\ &= a(\alpha) - \gamma(2\eta - 1)\alpha + \beta \end{aligned}$$

Using Fenchel–Young for smooth convex functions, we have

$$C_{\Phi}(\alpha; \eta) - \inf_{\alpha \in \mathbb{R}} C_{\Phi}(\alpha; \eta) \geq \frac{\gamma^2}{2L} \left(2\eta - 1 - \frac{1}{\gamma} a'(\alpha) \right)^2.$$

Calibration functions for square and logistic losses

Plug in $\alpha \leftarrow g(X)$ and $\eta \leftarrow \eta(X)$, and take the expectation to get

$$\begin{aligned}\mathcal{R}_{\Phi}[g] - \mathcal{R}_{\Phi}^* &\geq \frac{\gamma^2}{2L} \mathbb{E} \left[\left(2\eta(X) - 1 - \frac{1}{\gamma} a'(g(X)) \right)^2 \right] \\ &\geq \frac{\gamma^2}{2L} \left(\mathbb{E} \left[\left| 2\eta(X) - 1 - \frac{1}{\gamma} a'(g(X)) \right| \right] \right)^2 \\ &\geq \frac{\gamma^2}{2L} \left(\mathcal{R}_{\Phi^{0-1}}[g] - \mathcal{R}_{\Phi^{0-1}}^* \right)^2.\end{aligned}$$

Therefore, we conclude

$$\mathcal{R}_{\Phi^{0-1}}[g] - \mathcal{R}_{\Phi^{0-1}}^* \leq \frac{\sqrt{2L}}{\gamma} \left(\mathcal{R}_{\Phi}[g] - \mathcal{R}_{\Phi}^* \right)^{1/2}.$$

Outline

Calibration

$$\mathcal{R}_{\Phi^{0-1}} - \mathcal{R}_{\Phi^{0-1}}^* \stackrel{?}{\leq} H(\mathcal{R}_{\Phi} - \mathcal{R}_{\Phi}^*)$$

Calibration function with convex losses

Ridge least squares regression

Ridge regression

Consider the *ridge regression* problem with $\mu > 0$

$$\underset{\theta \in \mathbb{R}^d}{\text{minimize}} \quad \frac{1}{N} \|Y - \Phi\theta\|^2 + \mu \|\theta\|_2^2,$$

which has the minimizer

$$\hat{\theta}_\mu = \frac{1}{N} (\hat{\Sigma} + \mu I)^{-1} \Phi^\top y = (\Phi^\top \Phi + N\mu I)^{-1} \Phi^\top y = \Phi^\top (\Phi \Phi^\top + N\mu I)^{-1} y.$$

(Proof involving matrix inversion lemma in homework.)

Recall

$$\hat{\Sigma} = \frac{1}{N} \Phi^\top \Phi \in \mathbb{R}^{d \times d}.$$

Notably, we will no longer assume that $\hat{\Sigma}$ is invertible. Not assuming invertibility will be important when we consider kernel methods, where $d = \infty$ and $N < \infty$.

Ridge regression

Even though we consider a regularized optimization problem to obtain $\hat{\theta}_\mu$, we still consider the same (unregularized) risk

$$\mathcal{R}(\theta) = \mathbb{E}_{\varepsilon_1, \dots, \varepsilon_N} \left[\frac{1}{N} \|\Phi\theta - Y\|^2 \right].$$

Theorem

For the fixed design setting, with $\hat{\theta}_\mu = \frac{1}{N}(\hat{\Sigma} + \mu I)^{-1}\Phi^\top Y$ has expected excess risk

$$\mathbb{E}[\mathcal{R}(\hat{\theta})] - \mathcal{R}^* = \underbrace{\mu^2 \theta_\star^\top (\hat{\Sigma} + \mu I)^{-2} \hat{\Sigma} \theta_\star}_{\text{bias}} + \underbrace{\frac{\sigma^2}{N} \text{Tr}(\hat{\Sigma}^2 (\hat{\Sigma} + \mu I)^{-2})}_{\text{variance}}.$$

Proof. Recall that we had shown

$$\mathbb{E}[\mathcal{R}(\hat{\theta})] - \mathcal{R}_\star = \underbrace{\|\mathbb{E}[\hat{\theta}] - \theta_\star\|_{\hat{\Sigma}}^2}_{\text{bias}} + \underbrace{\mathbb{E}[\|\hat{\theta} - \mathbb{E}[\hat{\theta}]\|_{\hat{\Sigma}}^2]}_{\text{variance}},$$

First, we have

$$\begin{aligned}\mathbb{E}[\hat{\theta}_\mu] &= \frac{1}{N} \mathbb{E}[(\hat{\Sigma} + \mu I)^{-1} \Phi^\top (\Phi \theta_\star + \varepsilon)] \\ &= (\hat{\Sigma} + \mu I)^{-1} \hat{\Sigma} \theta_\star = (\hat{\Sigma} + \mu I)^{-1} (\hat{\Sigma} + \mu I - \mu I) \theta_\star \\ &= \theta_\star - \mu (\hat{\Sigma} + \mu I)^{-1} \theta_\star.\end{aligned}$$

So

$$\begin{aligned}\text{bias} &= \|\mu (\hat{\Sigma} + \mu I)^{-1} \theta_\star\|_{\hat{\Sigma}}^2 = \mu^2 \theta_\star^\top (\hat{\Sigma} + \mu I)^{-1} \hat{\Sigma} (\hat{\Sigma} + \mu I)^{-1} \theta_\star \\ &= \mu^2 \theta_\star^\top (\hat{\Sigma} + \mu I)^{-2} \hat{\Sigma} \theta_\star,\end{aligned}$$

which accounts for the first term.

Next, we have

$$\hat{\theta} - \mathbb{E}[\hat{\theta}] = \frac{1}{N}(\hat{\Sigma} + \mu I)^{-1}\Phi^T \varepsilon.$$

So,

$$\begin{aligned}\text{variance} &= \mathbb{E}\left[\|\hat{\theta} - \mathbb{E}[\hat{\theta}]\|_{\hat{\Sigma}}^2\right] \\ &= \frac{1}{N^2}\mathbb{E}\left[\text{Tr}(\varepsilon^T \Phi(\hat{\Sigma} + \mu I)^{-1}\hat{\Sigma}(\hat{\Sigma} + \mu I)^{-1}\Phi^T \varepsilon)\right] \\ &= \frac{\sigma^2}{N^2}\text{Tr}(\Phi(\hat{\Sigma} + \mu I)^{-1}\hat{\Sigma}(\hat{\Sigma} + \mu I)^{-1}\Phi^T) \\ &= \frac{\sigma^2}{N}\text{Tr}((\hat{\Sigma} + \mu I)^{-1}\hat{\Sigma}(\hat{\Sigma} + \mu I)^{-1}\hat{\Sigma}) \\ &= \frac{\sigma^2}{N}\text{Tr}(\hat{\Sigma}^2(\hat{\Sigma} + \mu I)^{-2}).\end{aligned}$$

□

Should we use $\mu > 0$?

For small $\mu > 0$, we have

$$\begin{aligned}\mathbb{E}[\mathcal{R}(\hat{\theta}_\mu)] - \mathcal{R}^* &= \mu^2 \theta_*^\top (\hat{\Sigma} + \mu I)^{-2} \hat{\Sigma} \theta_* + \frac{\sigma^2}{N} \text{Tr}(\hat{\Sigma}^2 (\hat{\Sigma} + \mu I)^{-2}) \\ &= \mathcal{O}(\mu^2) + \frac{\sigma^2}{N} \sum_{i=1}^{\min\{d, N\}} \frac{\lambda_i^2}{(\lambda_i + \mu)^2} \\ &= \mathcal{O}(\mu^2) + \frac{\sigma^2}{N} \sum_{i=1}^{\min\{d, N\}} \frac{1}{(1 + \mu/\lambda_i)^2} \\ &= \frac{\sigma^2}{N} \sum_{i=1}^{\min\{d, N\}} (1 - 2\mu\lambda_i) + \mathcal{O}(\mu^2) \\ &= \mathcal{O}(1) - \frac{2\sigma^2 \text{Tr}(\hat{\Sigma})}{N} \mu + \mathcal{O}(\mu^2).\end{aligned}$$

So the optimal value of μ is positive.

Optimizing regularization parameter

Theorem

Assume $\theta_\star \neq 0$. With

$$\mu_o = \frac{\sigma \text{Tr}(\widehat{\Sigma})^{1/2}}{\|\theta_\star\|_2 \sqrt{N}}$$

we have

$$\mathbb{E}[\mathcal{R}(\hat{\theta}_{\mu_o})] - \mathcal{R}^\star \leq \frac{\sigma \text{Tr}(\widehat{\Sigma})^{1/2} \|\theta_\star\|_2}{\sqrt{N}}.$$

(As we will see from the proof, μ_o is not the exact optimum, but rather a choice that optimizes an upper bound.)

Proof. Previously, we had shown that $\mathbb{E}[\mathcal{R}(\hat{\theta}_\mu)] - \mathcal{R}^* = \text{bias} + \text{variance}$.

First, bound the bias:

$$\text{bias} = \mu^2 \theta_*^\top (\hat{\Sigma} + \mu I)^{-2} \hat{\Sigma} \theta_* = \mu \theta_*^\top \underbrace{(\hat{\Sigma} + \mu I)^{-2} \mu \hat{\Sigma}}_{\preceq \frac{1}{2} I} \theta_* \leq \frac{\mu}{2} \|\theta_*\|^2,$$

where we use the fact that

$$\frac{\mu \lambda}{(\lambda + \mu)^2} \leq \frac{1}{2} \quad \forall \mu > 0, \lambda > 0.$$

Next, bound the variance:

$$\text{variance} = \frac{\sigma^2}{N} \text{Tr}(\hat{\Sigma}^2 (\hat{\Sigma} + \mu I)^{-2}) = \frac{\sigma^2}{\mu N} \text{Tr}(\hat{\Sigma} \underbrace{\mu \hat{\Sigma} (\hat{\Sigma} + \mu I)^{-2}}_{\preceq \frac{1}{2} I}) \leq \frac{\sigma^2}{2\mu N} \text{Tr} \hat{\Sigma}.$$

Finally, plugging in $\mu \leftarrow \mu_o$ (which minimizes the upper bounds on bias + variance), we conclude the statement. □

Compared to the expected excess risk of the least squares estimator without regularization

$$\mathbb{E}[\mathcal{R}(\hat{\theta}_0)] - \mathcal{R}^* = \frac{\sigma^2 d}{N}$$

the bound with regularization

$$\mathbb{E}[\mathcal{R}(\hat{\theta}_{\mu_o})] - \mathcal{R}^* \leq \frac{\sigma \text{Tr}(\hat{\Sigma})^{1/2} \|\theta_\star\|_2}{\sqrt{N}}$$

does not have an explicit dependence on d . Such bounds are said to be *dimension-independent*.

If $\|\varphi(x)\| \leq R$ for all x , then

$$\text{Tr}(\hat{\Sigma}) = \frac{1}{N} \sum_{i=1}^N \|\varphi(X_i)\|_2^2 \leq R^2,$$

and the only remaining (implicit) dependence on d is in $\|\theta_\star\|_2$.

However, the $\mathcal{O}(1/\sqrt{N})$ -rate is slower than the $\mathcal{O}(1/N)$ -rate. This is a common tradeoff in machine learning theory: a “fast rate” with bad constants vs. “slow rate” with good constants.