
Notes on PAC Bayes

Renjie Liao
Department of Computer Science
University of Toronto
rjliao@cs.toronto.edu
8th July, 2020

Abstract

In this note, we review the PAC-Bayesian approaches following [4, 3].

1 Prerequisites

We follow the literature to use KL^+ to denote the KL divergence between two Bernoulli distributions p and q as,

$$\text{KL}^+(p\|q) = p \ln \frac{p}{q} + (1-p) \ln \frac{1-p}{1-q}.$$

Recall the definition of strongly convex functions as follows.

Definition 1.1. A twice differentiable function f is called strongly convex with parameter $m > 0$ if for any x in the domain, we have

$$\nabla^2 f(x) \succeq mI.$$

Lemma 1.1. For any $p, q \in [0, 1]$, we have

$$\text{KL}^+(p\|q) \geq 2(p-q)^2,$$

Proof. Let $f(p) = \text{KL}^+(p\|q) - 2(p-q)^2$, we have

$$\begin{aligned} f'(p) &= \ln \left(\frac{p}{1-p} \right) - \ln \left(\frac{q}{1-q} \right) - 4(p-q) \\ f''(p) &= \frac{1}{p(1-p)} - 4 \end{aligned}$$

Since $p(1-p)$ achieves its maximum $1/4$ with $p = 1/2$, we have $f''(p) \geq 0, \forall p \in [0, 1]$. Note that $f'(q) = 0$. Hence, $f(p)$ decreases when $p \leq q$, increases when $p \geq q$, and achieves the minimum value $f(q) = 0$. \square

Lemma 1.2. For any $p, q \in [0, 1]$ with $p \leq q$, we have

$$\text{KL}^+(p\|q) \geq \frac{(p-q)^2}{2q},$$

Proof. Let $f(p) = \text{KL}^+(p\|q) - \frac{(p-q)^2}{2q}$, we have

$$\begin{aligned} f'(p) &= \ln \left(\frac{p}{1-p} \right) - \ln \left(\frac{q}{1-q} \right) - \frac{p-q}{q} \\ f''(p) &= \frac{1}{p(1-p)} - \frac{1}{q} = \frac{q-p+p^2}{qp(1-p)} \end{aligned}$$

Since $p \leq q$, then we have $f''(p) \geq 0$. Note that $f'(q) = 0$. Hence, $f(p)$ decreases when $p \leq q$. Therefore $f(p) \geq f(q) = 0$ which proves the claim. \square

Lemma 1.3. $\text{KL}^+(p||q)$ is 4-strongly convex w.r.t. argument p and convex w.r.t. argument q .

Proof. Similar to the proof of Lemma 1.1, denoting $f(p) = \text{KL}^+(p||q)$, we have

$$\begin{aligned} f'(p) &= \ln\left(\frac{p}{1-p}\right) - \ln\left(\frac{q}{1-q}\right) \\ f''(p) &= \frac{1}{p(1-p)} \geq 4. \end{aligned}$$

Therefore, according to the definition 1.1, $f(p)$ is 4-strongly convex.

Denoting $g(q) = \text{KL}^+(p||q)$, we have

$$\begin{aligned} g'(q) &= -\frac{p}{q} + \frac{1-p}{1-q} \\ g''(q) &= \frac{p}{q^2} + \frac{1-p}{(1-q)^2} > 0. \end{aligned}$$

Therefore $g(q)$ is convex. \square

Lemma 1.4. (Hoeffding's Lemma [2]) For bounded random variables X_1, \dots, X_m where $X_i \in [0, 1]$ and is i.i.d. with mean μ , let $\bar{X}_m = \frac{1}{m} \sum_{i=1}^m X_i$, then for any $\epsilon > 0$ with $\mu + \epsilon < 1$ and $\mu - \epsilon > 0$,

$$\begin{aligned} \mathbb{P}(\bar{X}_m \geq \mu + \epsilon) &\leq e^{-m \text{KL}^+(\mu + \epsilon || \mu)} \leq e^{-2m\epsilon^2} \\ \mathbb{P}(\bar{X}_m \leq \mu - \epsilon) &\leq e^{-m \text{KL}^+(\mu - \epsilon || \mu)} \leq e^{-2m\epsilon^2} \end{aligned}$$

Proof. Let $t > 0$, $S_m = \sum_{i=1}^m X_i$, $p = \mu + \epsilon$, $\bar{p} = 1 - p$, and $\bar{\mu} = 1 - \mu$, we have

$$\begin{aligned} \mathbb{P}(\bar{X}_m - \mu \geq \epsilon) &= \mathbb{P}(S_m - m\mu \geq m\epsilon) \\ &= \mathbb{P}\left(e^{tS_m} \geq e^{tm(\mu + \epsilon)}\right) \\ &\leq \frac{\mathbb{E}[e^{tS_m}]}{e^{tm(\mu + \epsilon)}} \quad (\text{Markov Inequality}) \\ &= \frac{\mathbb{E}[e^{tX_1}]^m}{e^{tm(\mu + \epsilon)}} \\ &\leq \frac{\mathbb{E}[X_1 e^t + (1 - X_1)e^0]^m}{e^{tm(\mu + \epsilon)}} \quad (\text{Convexity of } e^{tX}) \\ &= \frac{(\mu e^t + 1 - \mu)^m}{e^{tm(\mu + \epsilon)}} \\ &= \left(\frac{\mu e^t + 1 - \mu}{e^{t(\mu + \epsilon)}}\right)^m \\ &= (\mu e^{t\bar{p}} + \bar{\mu} e^{-tp})^m \end{aligned} \tag{1}$$

Let $f(t) = m \ln g(t) = m \ln(\mu e^{t\bar{p}} + \bar{\mu} e^{-tp})$, we have

$$\begin{aligned} f'(t) &= m \frac{g'(t)}{g(t)} \\ f''(t) &= m \frac{m(g(t)g''(t) - g'(t)^2)}{g(t)^2} \\ g'(t) &= \mu \bar{p} e^{t\bar{p}} - \bar{\mu} p e^{-tp} \\ g''(t) &= \mu \bar{p}^2 e^{t\bar{p}} + \bar{\mu} p^2 e^{-tp} \\ g(t)g''(t) - g'(t)^2 &= \mu \bar{\mu} e^{t(\bar{p}-p)} (\bar{p} - p)^2. \end{aligned}$$

Since $f''(t) \geq 0$, f achieves its minimum while $e^t = \frac{\bar{\mu}p}{\mu\bar{p}}$,

$$\begin{aligned}
\min_t f(t) &= \min_t m \ln \left(\mu e^{t(1-p)} + \bar{\mu} e^{-tp} \right) \\
&= \min_t m \ln \left(e^{-tp} (\mu e^t + \bar{\mu}) \right) \\
&= m \ln \left(\left(\frac{\mu\bar{p}}{\bar{\mu}p} \right)^p \left(\frac{\bar{\mu}p}{\bar{p}} + \bar{\mu} \right) \right) \\
&= m \ln \left(\frac{\mu^p \bar{p}^{(p-1)}}{\bar{\mu}^{(p-1)} p^p} \right) \\
&= m \ln \left(\frac{\mu^p \bar{\mu}^{\bar{p}}}{\bar{p}^{\bar{p}} p^p} \right) \tag{2}
\end{aligned}$$

Combine Eq. (1) and Eq. (2), we have

$$\begin{aligned}
\mathbb{P}(\bar{X}_m - \mu \geq \epsilon) &\leq \min_t e^{f(t)} \\
&= \left(\frac{\mu^p \bar{\mu}^{\bar{p}}}{\bar{p}^{\bar{p}} p^p} \right)^m \\
&= e^{-mp \ln(\frac{p}{\mu}) - m\bar{p} \ln(\frac{\bar{p}}{\bar{\mu}})} \\
&= e^{-m \text{KL}^+(\mu + \epsilon \| \mu)} \\
&\leq e^{-2m\epsilon^2} \quad (\text{Lemma 1.1})
\end{aligned}$$

Similarly, we can prove the case for the other side. □

Remark. This proof is based on the one from [8].

Lemma 1.5. For non-negative continuous random variables X , we have

$$\mathbb{E}[X] = \int_0^\infty \mathbb{P}(X \geq \nu) d\nu.$$

Proof.

$$\begin{aligned}
\mathbb{E}[X] &= \int_0^\infty X \mathbb{P}(X) dX \\
&= \int_0^\infty \int_0^X \mathbf{1} d\nu \mathbb{P}(X) dX \\
&= \int_0^\infty \int_0^X \mathbb{P}(X) d\nu dX \\
&= \int_0^\infty \int_\nu^\infty \mathbb{P}(X) dX d\nu \quad (\text{region of the integral is the same}) \\
&= \int_0^\infty \mathbb{P}(X \geq \nu) d\nu
\end{aligned}$$

□

Remark. This lemma can also be obtained via **Integration by Parts**. Specifically, let $f(\nu) = \int_\nu^\infty \mathbb{P}(X) dX = \mathbb{P}(X \geq \nu)$ and $g(\nu) = \int_0^\nu \mathbf{1} dX = \nu$, we have

$$\begin{aligned}
\int_0^\infty \mathbb{P}(X \geq \nu) d\nu &= \int_0^\infty f(\nu) g'(\nu) d\nu = f(\nu) g(\nu) \Big|_0^\infty - \int_0^\infty f'(\nu) g(\nu) d\nu \\
&= - \int_0^\infty -\mathbb{P}(\nu) \nu d\nu \\
&= \mathbb{E}[\nu] = \mathbb{E}[X]
\end{aligned}$$

Similar result holds for discrete nonnegative random variables.

Lemma 1.6. [2-side] Let X be a random variable satisfying $\mathbb{P}(X \geq \epsilon) \leq e^{-2m\epsilon^2}$ and $\mathbb{P}(X \leq -\epsilon) \leq e^{-2m\epsilon^2}$ where $m \geq 1$ and $\epsilon > 0$, we have

$$\mathbb{E}[e^{2(m-1)X^2}] \leq 2m.$$

Proof.

$$\begin{aligned} \mathbb{E}[e^{2(m-1)X^2}] &= \int_0^\infty \mathbb{P}\left(e^{2(m-1)X^2} \geq \nu\right) d\nu \quad (\text{Lemma 1.5}) \\ &= \int_0^\infty \mathbb{P}\left(X^2 \geq \frac{\ln \nu}{2(m-1)}\right) d\nu \\ &= \int_0^\infty \mathbb{P}\left(X \geq \sqrt{\frac{\ln \nu}{2(m-1)}}\right) d\nu + \int_0^\infty \mathbb{P}\left(X \leq -\sqrt{\frac{\ln \nu}{2(m-1)}}\right) d\nu \quad (3) \end{aligned}$$

$$\begin{aligned} \int_0^\infty \mathbb{P}\left(X \geq \sqrt{\frac{\ln \nu}{2(m-1)}}\right) d\nu &= \int_0^1 \mathbb{P}\left(X \geq \sqrt{\frac{\ln \nu}{2(m-1)}}\right) d\nu + \int_1^\infty \mathbb{P}\left(X \geq \sqrt{\frac{\ln \nu}{2(m-1)}}\right) d\nu \\ &\leq 1 + \int_1^\infty \mathbb{P}\left(X \geq \sqrt{\frac{\ln \nu}{2(m-1)}}\right) d\nu \\ &\leq 1 + \int_1^\infty e^{-2m \frac{\ln \nu}{2(m-1)}} d\nu \\ &= 1 + \left(-m \nu^{-\frac{1}{m-1}} \Big|_1^\infty\right) \\ &= m \quad (4) \end{aligned}$$

Similarly, we can show that

$$\int_0^\infty \mathbb{P}\left(X \leq -\sqrt{\frac{\ln \nu}{2(m-1)}}\right) d\nu \leq m \quad (5)$$

Combining Eq. (3) and Eq. (4), we finish the proof. \square

Remark. One can also obtain the 1-side version $\mathbb{E}[e^{2(m-1)X^2}] \leq m$ by, e.g., adding the assumption that $X \geq 0$ and following the same proof.

Lemma 1.7. [2-side] Let X be a random variable defined on $(0, 1)$ with mean μ satisfying $\mathbb{P}(X \geq \mu + \epsilon) \leq e^{-m \text{KL}^+(\mu + \epsilon \| \mu)}$ and $\mathbb{P}(X \leq \mu - \epsilon) \leq e^{-m \text{KL}^+(\mu - \epsilon \| \mu)}$ where $m \geq 1$, $\epsilon > 0$, $\mu + \epsilon < 1$ and $\mu - \epsilon > 0$, we have

$$\mathbb{E}[e^{(m-1) \text{KL}^+(X \| \mu)}] \leq 2m.$$

Proof. Denoting $f(p) = \text{KL}^+(p \| \mu)$, we have

$$\begin{aligned} f'(p) &= \ln\left(\frac{p}{1-p}\right) - \ln\left(\frac{\mu}{1-\mu}\right) \\ f''(p) &= \frac{1}{p(1-p)}. \end{aligned}$$

Since $f''(p) \geq 0$ for all $p \in (0, 1)$, $f(p)$ decreases when $0 < p \leq \mu$, increases when $\mu \leq p < 1$ and attains its minimum 0 at $p = \mu$. Therefore, based on the inverse function theorem, there exists the inverse function ϕ of $\text{KL}^+(p \| \mu)$ when $p \in (0, \mu)$, i.e., $p = \phi(\text{KL}^+(p \| \mu))$, $\forall p \in (0, \mu)$. Similarly, there exists another inverse function ψ of $\text{KL}^+(p \| \mu)$ when $p \in (\mu, 1)$, i.e., $p = \psi(\text{KL}^+(p \| \mu))$, $\forall p \in (\mu, 1)$. Note that the functional forms of ϕ and ψ may or may not be the same (i.e., whether

$\text{KL}^+(p\|\mu)$ is symmetric w.r.t. $p = \mu$ depending on the value of μ . As shown below, the exact functional form of the inverse function does not matter as long as it exists.

$$\begin{aligned}\mathbb{E}[e^{(m-1)\text{KL}^+(X\|q)}] &= \int_0^\infty \mathbb{P}\left(e^{(m-1)\text{KL}^+(X\|q)} \geq \nu\right) d\nu \quad (\text{Lemma 1.5}) \\ &= \int_0^\infty \mathbb{P}\left(\text{KL}^+(X\|q) \geq \frac{\ln \nu}{m-1}\right) d\nu \\ &= \int_0^\infty \mathbb{P}\left(X \geq \phi\left(\frac{\ln \nu}{m-1}\right)\right) d\nu + \int_0^\infty \mathbb{P}\left(X \leq \psi\left(\frac{\ln \nu}{m-1}\right)\right) d\nu \quad (6)\end{aligned}$$

We use the same trick again as in Lemma 1.6,

$$\begin{aligned}\int_0^\infty \mathbb{P}\left(X \geq \phi\left(\frac{\ln \nu}{m-1}\right)\right) d\nu &= \int_0^1 \mathbb{P}\left(X \geq \phi\left(\frac{\ln \nu}{m-1}\right)\right) d\nu + \int_1^\infty \mathbb{P}\left(X \geq \phi\left(\frac{\ln \nu}{m-1}\right)\right) d\nu \\ &\leq 1 + \int_1^\infty \mathbb{P}\left(X \geq \phi\left(\frac{\ln \nu}{m-1}\right)\right) d\nu \\ &= 1 + \int_1^\infty e^{-\frac{m \ln \nu}{m-1}} d\nu \\ &= 1 + \left(-m-1 \nu^{-\frac{1}{m-1}} \Big|_1^\infty\right) \\ &= m \quad (7)\end{aligned}$$

Similarly, we have $\int_0^\infty \mathbb{P}\left(X \leq \psi\left(\frac{\ln \nu}{m-1}\right)\right) d\nu \leq m$. Therefore, combining it with Eq. (6) and Eq. (7), we prove the claim. \square

Remark. Similarly, one can obtain the 1-side version $\mathbb{E}[e^{(m-1)\text{KL}^+(X\|\mu)}] \leq m$ by, e.g., adding the assumption that $X \geq \mu$ and following the same proof.

2 Main Result

We only consider the binary classification problem. Let us first introduce some basic setup.

- Data $z, z = (x, y)$ input $x \in \mathbb{R}^d$ and the output $y \in \{0, 1\}$
- Data Space $\mathcal{Z}, z \in \mathcal{Z}$
- Data Distribution $D, z \stackrel{iid}{\sim} D$
- Hypothesis h , model
- Hypothesis Class $\mathcal{H}, h \in \mathcal{H}$
- Training Set S with size $m, S = \{z_1, \dots, z_m\}$
- Loss $\ell, \ell : \mathcal{H} \times \mathcal{Z} \rightarrow \{0, 1\}$

As usual, we care about the *generalization error*, i.e., *error rate* or *misclassification rate* in this case,

$$L_D(h) = \mathbb{P}_{z \sim D}(h(x) \neq y), \quad (8)$$

where we use the subscript to emphasize the dependency on the data distribution. However, since we can not observe it directly, we approximate it using the empirical distribution, a.k.a., *empirical error*,

$$L_S(h) = \frac{1}{m} \sum_{i=1}^m \ell(h, z_i) = \frac{1}{m} \sum_{i=1}^m \mathbf{1}[h(x_i) \neq y_i]. \quad (9)$$

where ℓ is the 0-1 loss.

Bayesian View We are interested in bounding the generalization error using the empirical error. PAC-Bayes [4, 3] takes a Bayesian view of PAC theory [7]. In particular, it assumes that we have a prior distribution P over the hypothesis class \mathcal{H} and we gonna obtain the posterior Q after the learning process on the training set. We can define the generalization error and empirical error in terms of this Bayesian view as,

$$\begin{aligned} L_S(Q) &= \mathbb{E}_{h \sim Q}[L_S(h)] \\ L_D(Q) &= \mathbb{E}_{h \sim Q}[L_D(h)]. \end{aligned}$$

Before introducing the main results, let us prove some useful lemmas. First, note that $L_S(h)$ and $L_D(h)$ are distributions over Bernoulli random variables themselves.

Lemma 2.1. For any h and any $\epsilon > 0$ with $L_D(h) + \epsilon < 1$ and $L_D(h) - \epsilon > 0$, we have

$$\begin{aligned} \mathbb{P}(L_S(h) \geq L_D(h) + \epsilon) &\leq e^{-m \text{KL}^+(L_D(h) + \epsilon \| L_D(h))} \\ \mathbb{P}(L_S(h) \leq L_D(h) - \epsilon) &\leq e^{-m \text{KL}^+(L_D(h) - \epsilon \| L_D(h))} \end{aligned}$$

Proof. First, recall the definition of $L_S(h)$ in Eq. (9) and observe that for any h we have $\mathbb{E}_S[L_S(h)] = L_D(h)$. Then we replace μ and \bar{X}_m in Lemma 1.4 with $L_D(h)$ and $L_S(h)$, we finish the proof. \square

Lemma 2.2. For any distribution Q over \mathcal{H} , $\forall h \in \mathcal{H}$, $p(h) \in (0, 1)$ and $q(h) \in (0, 1)$, we have

$$\text{KL}^+(\mathbb{E}_{h \sim Q}[p(h)] \| \mathbb{E}_{h \sim Q}[q(h)]) \leq \mathbb{E}_{h \sim Q} [\text{KL}^+(p(h) \| q(h))]$$

Proof. Denoting $f(p) = \text{KL}^+(p \| q)$, from Lemma 1.3, we know $f(p)$ is strongly convex for any q . Therefore, based on the Jensen's inequality, we have

$$\text{KL}^+(\mathbb{E}_{h \sim Q}[p(h)] \| \mathbb{E}_{h \sim Q}[q(h)]) \leq \mathbb{E}_{h \sim Q} [\text{KL}^+(p(h) \| \mathbb{E}_{h \sim Q}[q(h)])]. \quad (10)$$

Denoting $g(q) = \text{KL}^+(p \| q)$, from Lemma 1.3, we know $g(q)$ is convex for any p ,

$$\text{KL}^+(p(h) \| \mathbb{E}_{h \sim Q}[q(h)]) \leq \mathbb{E}_{h \sim Q} [\text{KL}^+(p(h) \| q(h))]. \quad (11)$$

Combining Eq. (10) and Eq. (11), we have

$$\begin{aligned} \text{KL}^+(\mathbb{E}_{h \sim Q}[p(h)] \| \mathbb{E}_{h \sim Q}[q(h)]) &\leq \mathbb{E}_{h \sim Q} [\mathbb{E}_{h \sim Q} [\text{KL}^+(p(h) \| q(h))]] \\ &= \mathbb{E}_{h \sim Q} [\text{KL}^+(p(h) \| q(h))], \end{aligned}$$

which proves the claim. \square

2.1 Generalization Bound with KL divergence

Now we introduce the generalization bound which uses KL^+ to measure the distance between prior P and posterior Q over models.

Theorem 2.3. Let P be a prior distribution over \mathcal{H} and let $\delta \in (0, 1)$. Then, with probability $1 - \delta$ over the choice of an i.i.d. training set S according to D , for all distributions Q over \mathcal{H} , we have

$$\text{KL}^+(L_S(Q) \| L_D(Q)) \leq \frac{\text{KL}(Q \| P) + \ln \frac{2m}{\delta}}{m - 1}$$

Proof. Fix a hypothesis h , based on Lemma 2.1, we have

$$\begin{aligned} \mathbb{P}(L_S(h) \geq L_D(h) + \epsilon) &\leq e^{-m \text{KL}^+(L_D(h) + \epsilon \| L_D(h))} \\ \mathbb{P}(L_S(h) \leq L_D(h) - \epsilon) &\leq e^{-m \text{KL}^+(L_D(h) - \epsilon \| L_D(h))}. \end{aligned}$$

Then, based on Lemma 1.7, we have

$$\mathbb{E}_S[e^{(m-1) \text{KL}^+(L_S(h) \| L_D(h))}] \leq 2m. \quad (12)$$

Therefore, for all S and any $\delta > 0$, we have,

$$\mathbb{P} \left(e^{(m-1) \text{KL}^+(L_S(h)\|L_D(h))} \geq \frac{2m}{\delta} \right) \leq \frac{\mathbb{E}_S[e^{(m-1) \text{KL}^+(L_S(h)\|L_D(h))}] \delta}{2m} \leq \delta \quad (13)$$

For any function $f(h)$, we have

$$\begin{aligned} \mathbb{E}_{h \sim Q}[f(h)] &= \mathbb{E}_{h \sim Q}[\ln e^{f(h)}] \\ &= \mathbb{E}_{h \sim Q}[\ln e^{f(h)} + \ln \frac{Q}{P} + \ln \frac{P}{Q}] \\ &= \text{KL}(Q\|P) + \mathbb{E}_{h \sim Q} \left[\ln \left(\frac{P}{Q} e^{f(h)} \right) \right] \\ &\leq \text{KL}(Q\|P) + \ln \mathbb{E}_{h \sim Q} \left[\frac{P}{Q} e^{f(h)} \right] \\ &= \text{KL}(Q\|P) + \ln \mathbb{E}_{h \sim P} [e^{f(h)}]. \end{aligned} \quad (14)$$

Let $f(h) = (m-1) \text{KL}^+(L_S(h)\|L_D(h))$. From Eq. (13) and Eq. (14), we have, with at least probability $1 - \delta$,

$$\begin{aligned} (m-1) \mathbb{E}_{h \sim Q}[\text{KL}^+(L_S(h)\|L_D(h))] &\leq \text{KL}(Q\|P) + \ln \mathbb{E}_{h \sim P} [e^{(m-1) \text{KL}^+(L_S(h)\|L_D(h))}] \\ &\leq \text{KL}(Q\|P) + \ln \left(\frac{2m}{\delta} \right). \end{aligned} \quad (15)$$

From Lemma 2.2, we have

$$\begin{aligned} (m-1) \text{KL}^+(\mathbb{E}_{h \sim Q}[L_S(h)]\|\mathbb{E}_{h \sim Q}[L_D(h)]) &\leq (m-1) \mathbb{E}_{h \sim Q}[\text{KL}^+(L_S(h)\|L_D(h))] \\ &\leq \text{KL}(Q\|P) + \ln \left(\frac{2m}{\delta} \right), \end{aligned} \quad (16)$$

which proves the theorem. \square

Remark. The proof follows the original proof of [3]. Again, we can have a one-side version,

$$\text{KL}^+(L_S(Q)\|L_D(Q)) \leq \frac{\text{KL}(Q\|P) + \ln \frac{m}{\delta}}{m-1},$$

by adding the assumption $L_S(Q) > L_D(Q)$ which is reasonable in practice.

One can also prove a slight different generalization bound by using a different technique as below.

Theorem 2.4. Let P be a prior distribution over \mathcal{H} and let $\delta \in (0, 1)$. Then, with probability $1 - \delta$ over the choice of an i.i.d. training set S according to D , for all distributions Q over \mathcal{H} , we have

$$\text{KL}^+(L_S(Q)\|L_D(Q)) \leq \frac{\text{KL}(Q\|P) + \ln \frac{m+1}{\delta}}{m}$$

Proof. From Eq. (14) in the proof of Theorem 2.3, for any function $f(h)$, we have

$$\mathbb{E}_{h \sim Q}[f(h)] \leq \text{KL}(Q\|P) + \ln \mathbb{E}_{h \sim P} [e^{f(h)}]. \quad (17)$$

Based on Lemma 2.2 and let $f(h) = m \text{KL}^+(L_S(h)\|L_D(h))$, we have

$$\begin{aligned} \text{KL}^+(L_S(Q)\|L_D(Q)) &\leq \mathbb{E}_{h \sim Q}[\text{KL}^+(L_S(h)\|L_D(h))] \\ &= \mathbb{E}_{h \sim Q} \left[\frac{1}{m} f(h) \right] \\ &\leq \frac{\text{KL}(Q\|P) + \ln \mathbb{E}_{h \sim P} [e^{f(h)}]}{m}. \end{aligned} \quad (18)$$

Note that since we consider 0-1 loss and samples are i.i.d., $mL_S(h)$ can only take values from $\{0, 1, 2, \dots, m\}$ and follows the binomial distribution $B(m, L_D(h))$. Hence, we have,

$$\begin{aligned}
\mathbb{E}_S \left[\mathbb{E}_{h \sim P} \left[e^{f(h)} \right] \right] &= \mathbb{E}_S \left[\mathbb{E}_{h \sim P} \left[e^{m \text{KL}^+(L_S(h) \| L_D(h))} \right] \right] \\
&= \mathbb{E}_S \left[\mathbb{E}_{h \sim P} \left[e^{mL_S(h) \ln \frac{L_S(h)}{L_D(h)} + m(1-L_S(h)) \ln \frac{1-L_S(h)}{1-L_D(h)}} \right] \right] \\
&= \mathbb{E}_S \left[\mathbb{E}_{h \sim P} \left[\left(\frac{L_S(h)}{L_D(h)} \right)^{mL_S(h)} \left(\frac{1-L_S(h)}{1-L_D(h)} \right)^{m(1-L_S(h))} \right] \right] \\
&= \mathbb{E}_{h \sim P} \left[\mathbb{E}_S \left[\left(\frac{L_S(h)}{L_D(h)} \right)^{mL_S(h)} \left(\frac{1-L_S(h)}{1-L_D(h)} \right)^{m(1-L_S(h))} \right] \right] \\
&= \mathbb{E}_{h \sim P} \left[\sum_{k=0}^m \binom{m}{k} L_D(h)^k (1-L_D(h))^{m-k} \left(\frac{k/m}{L_D(h)} \right)^k \left(\frac{1-k/m}{1-L_D(h)} \right)^{m-k} \right] \\
&= \mathbb{E}_{h \sim P} \left[\sum_{k=0}^m \binom{m}{k} \left(\frac{k}{m} \right)^k \left(\frac{m-k}{m} \right)^{m-k} \right] \\
&\leq m+1, \tag{19}
\end{aligned}$$

where the last inequality uses the fact that $\binom{m}{k} \left(\frac{k}{m} \right)^k \left(\frac{m-k}{m} \right)^{m-k}$ is the probability of a binomial random variable (following $B(m, \frac{k}{m})$) taking the value as k , thus being no larger than 1. The second to the last line makes use of the law of the unconscious statistician (LOTUS).

Based on Markov's inequality, for any $\delta > 0$, we have

$$\mathbb{P} \left(\mathbb{E}_{h \sim P} \left[e^{f(h)} \right] \geq \frac{m+1}{\delta} \right) \leq \frac{\delta \mathbb{E}_S \left[\mathbb{E}_{h \sim P} \left[e^{f(h)} \right] \right]}{m+1} \leq \delta, \tag{20}$$

which along with Eq. (18) proves the theorem. \square

Remark. This proof largely follows the one in [1]. The observation that, for any h , $mL_S(h)$ is distributed according to the binomial distribution $B(m, L_D(h))$ is really insightful. One can further improve Eq. (19) to show $\sqrt{m} \leq \mathbb{E}_S \left[\mathbb{E}_{h \sim P} \left[e^{f(h)} \right] \right] \leq \sqrt{2m}$. How?

2.2 Canonical Generalization Bound

Theorem 2.5. Let P be a prior distribution over \mathcal{H} and let $\delta \in (0, 1)$. Then, with probability $1 - \delta$ over the choice of an i.i.d. training set S according to D , for all distributions Q over \mathcal{H} , we have

$$L_D(Q) \leq L_S(Q) + \sqrt{\frac{\text{KL}(Q \| P) + \ln \frac{2m}{\delta}}{2(m-1)}}$$

Proof. (Proof-I) From Lemma 1.1 and Theorem 2.3, we have

$$2(L_S(Q) - L_D(Q))^2 \leq \text{KL}^+(L_S(Q) \| L_D(Q)) \leq \frac{\text{KL}(Q \| P) + \ln \frac{2m}{\delta}}{m-1}.$$

Therefore, we have,

$$L_D(Q) \leq L_S(Q) + \sqrt{\frac{\text{KL}(Q \| P) + \ln \frac{2m}{\delta}}{2(m-1)}},$$

which proves the theorem. \square

Proof. (Proof-II) Let $\Delta(h) = L_D(h) - L_S(h)$. From Eq. (14), for any function $f(h)$, we have

$$\mathbb{E}_{h \sim Q} [f(h)] \leq \text{KL}(Q \| P) + \ln \mathbb{E}_{h \sim P} \left[e^{f(h)} \right]. \tag{21}$$

Let $f(h) = 2(m-1)\Delta(h)^2$. We have

$$\begin{aligned} 2(m-1)\mathbb{E}_{h\sim Q}[\Delta(h)^2] &\leq 2(m-1)\mathbb{E}_{h\sim Q}[\Delta(h)^2] \quad (\text{Jensen's inequality}) \\ &\leq \text{KL}(Q\|P) + \ln \mathbb{E}_{h\sim P} \left[e^{2(m-1)\Delta(h)^2} \right]. \end{aligned} \quad (22)$$

Since $L_D(h) \in [0, 1]$, based on Hoeffding's inequality, for any $\epsilon > 0$, we have

$$\begin{aligned} \mathbb{P}(\Delta(h) \geq \epsilon) &\leq e^{-2m\epsilon^2} \\ \mathbb{P}(\Delta(h) \leq -\epsilon) &\leq e^{-2m\epsilon^2} \end{aligned}$$

Hence, based on Lemma 1.6, we have

$$\begin{aligned} \mathbb{E}_S \left[e^{2(m-1)\Delta(h)^2} \right] \leq 2m &\Rightarrow \mathbb{E}_{h\sim P} \left[\mathbb{E}_S \left[e^{2(m-1)\Delta(h)^2} \right] \right] \leq 2m \\ &\Leftrightarrow \mathbb{E}_S \left[\mathbb{E}_{h\sim P} \left[e^{2(m-1)\Delta(h)^2} \right] \right] \leq 2m \end{aligned}$$

Based on Markov's inequality, we have

$$\mathbb{P} \left(\mathbb{E}_{h\sim P} \left[e^{2(m-1)\Delta(h)^2} \right] \geq \frac{2m}{\delta} \right) \leq \frac{\delta \mathbb{E}_S \left[\mathbb{E}_{h\sim P} \left[e^{2(m-1)\Delta(h)^2} \right] \right]}{2m} \leq \delta. \quad (23)$$

Combining Eq. (22) and Eq. (23), with probability $1 - \delta$, we have

$$\mathbb{E}_{h\sim Q}[\Delta(h)^2] \leq \frac{\text{KL}(Q\|P) + \ln \left(\frac{2m}{\delta} \right)}{2(m-1)} \quad (24)$$

which proves the theorem. \square

Remark. This theorem has a simple form and is more similar to the majority of generalization bounds. Therefore, it is frequently used in the literature. Proof-I is based on the one in [3, 1]. Proof-II is based on the one in the chapter 31 of [6]. The proof technique in Proof-II is more general in a sense that one can generalize the loss beyond 0-1 loss under this framework and derive similar results.

Theorem 2.6. Let P be a prior distribution over \mathcal{H} and let $\delta \in (0, 1)$. Then, with probability $1 - \delta$ over the choice of an i.i.d. training set S according to D , for all distributions Q over \mathcal{H} , we have

$$L_D(Q) \leq L_S(Q) + \sqrt{\frac{2L_S(Q) (\text{KL}(Q\|P) + \ln \frac{2m}{\delta})}{m-1}} + \frac{2 (\text{KL}(Q\|P) + \ln \frac{2m}{\delta})}{m-1}$$

Proof. From Lemma 1.2, we have for any $p, q \in [0, 1]$ with $p \leq q$,

$$\text{KL}^+(p\|q) \geq \frac{(p-q)^2}{2q}.$$

If $\text{KL}^+(p\|q) \leq x$, then we have

$$x \geq \frac{(p-q)^2}{2q} \Leftrightarrow 2qx \geq (p-q)^2 \Leftrightarrow q \leq p + \sqrt{2qx}. \quad (25)$$

Note that

$$\begin{aligned} q \leq p + \sqrt{2qx} &\Leftrightarrow \left(\sqrt{q} - \frac{\sqrt{2x}}{2} \right)^2 \leq p + \frac{x}{2} \\ &\Leftrightarrow \sqrt{2q} \leq \sqrt{2p+x} + \sqrt{x} \end{aligned} \quad (26)$$

Based on Eq. (25) and Eq. (26), we have

$$\begin{aligned} q &\leq p + \sqrt{2qx} \\ &\leq p + (\sqrt{2p+x} + \sqrt{x})\sqrt{x} \\ &= p + \sqrt{2px+x^2} + x \\ &\leq p + \sqrt{2px} + 2x \quad (\text{Subadditivity: } \sqrt{x+y} < \sqrt{x} + \sqrt{y}) \end{aligned} \quad (27)$$

Let $p = L_S(Q)$, $q = L_D(Q)$, and $x = \frac{\text{KL}(Q\|P) + \ln \frac{2m}{\delta}}{m-1}$. Theorem 2.3 shows $\text{KL}^+(p\|q) \leq x$. Then Eq. (27) proves the theorem. \square

Remark. Note that whether Theorem 2.5 or Theorem 2.6 provides a sharper bound depends on the actual value of $L_S(Q)$. But Theorem 2.5 has a simpler form.

2.3 Generalization Bound of Deterministic Models

Let us first review a result from [5] which generalizes the PAC-Bayes bound to a general class of deterministic models. We define the model to be $f_w \in \mathcal{H} : \mathcal{X} \rightarrow \mathbb{R}^k$ where w are the parameters of the model, \mathcal{X} is the input space, and $k \geq 1$. We also define the γ -margin loss for the k -category classification as,

$$\tilde{L}_D(f_w, \gamma) = \mathbb{P}_{z \sim D} \left(f_w(x)[y] \leq \gamma + \max_{j \neq y} f_w(x)[j] \right),$$

where $\gamma > 0$ and $f_w(x)[j]$ means the j -th output of the model. Accordingly, we can define the empirical version,

$$\tilde{L}_S(f_w, \gamma) = \frac{1}{m} \sum_{z_i \in S} \mathbf{1} \left(f_w(x)[y] \leq \gamma + \max_{j \neq y} f_w(x)[j] \right),$$

We use \mathbb{N}_m^+ to denote the first m positive integers, i.e., $\mathbb{N}_m^+ = \{1, 2, \dots, m\}$.

Theorem 2.7. Let $f_w(x) : \mathcal{X} \rightarrow \mathbb{R}^k$ be any model with parameters w , and P be any distribution on the parameters that is independent of the training data. For any w , we construct a posterior $Q(w + u)$ by adding any random perturbation u to w , s.t., $\mathbb{P}(\max_{x \in \mathcal{X}} \|f_{w+u}(x) - f_w(x)\|_\infty < \frac{\gamma}{4}) > \frac{1}{2}$. Then, for any $\gamma, \delta > 0$, with probability at least $1 - \delta$ over the size- m training set S , for any w , we have:

$$\tilde{L}_D(f_w, 0) \leq \tilde{L}_S(f_w, \gamma) + \sqrt{\frac{2 \text{KL}(Q \| P) + \ln \frac{8m}{\delta}}{2(m-1)}} \quad (28)$$

Proof. Let $\tilde{w} = w + u$. Let \mathcal{C} be the set of perturbation with the following property,

$$\mathcal{C} = \left\{ w' \mid \max_{x \in \mathcal{X}} \|f_{w'}(x) - f_w(x)\|_\infty < \frac{\gamma}{4} \right\}. \quad (29)$$

$\tilde{w} = w + u$ (w is deterministic and u is stochastic) is distributed according to $Q(\tilde{w})$. We now construct a new posterior \tilde{Q} as follows,

$$\tilde{Q}(\tilde{w}) = \begin{cases} \frac{1}{Z} Q(\tilde{w}) & \tilde{w} \in \mathcal{C} \\ 0 & \tilde{w} \in \bar{\mathcal{C}}. \end{cases} \quad (30)$$

Here $Z = \int_{\tilde{w} \in \mathcal{C}} d\tilde{Q}(\tilde{w}) = \mathbb{P}(\tilde{w} \in \mathcal{C})$. We know from the assumption that $Z > \frac{1}{2}$. $\bar{\mathcal{C}}$ is the complement set of \mathcal{C} . Therefore, for any $\tilde{w} \sim \tilde{Q}$, we have

$$\begin{aligned} & \max_{i \in \mathbb{N}_k^+, j \in \mathbb{N}_k^+, x \in \mathcal{X}} \left| |f_{\tilde{w}}(x)[i] - f_{\tilde{w}}(x)[j]| - |f_w(x)[i] - f_w(x)[j]| \right| \\ & \leq \max_{i \in \mathbb{N}_k^+, j \in \mathbb{N}_k^+, x \in \mathcal{X}} \left| f_{\tilde{w}}(x)[i] - f_{\tilde{w}}(x)[j] - f_w(x)[i] + f_w(x)[j] \right| \\ & \leq \max_{i \in \mathbb{N}_k^+, j \in \mathbb{N}_k^+, x \in \mathcal{X}} \left| f_{\tilde{w}}(x)[i] - f_w(x)[i] \right| + \left| f_{\tilde{w}}(x)[j] - f_w(x)[j] \right| \\ & \leq \max_{i \in \mathbb{N}_k^+, x \in \mathcal{X}} \left| f_{\tilde{w}}(x)[i] - f_w(x)[i] \right| + \max_{j \in \mathbb{N}_k^+, x \in \mathcal{X}} \left| f_{\tilde{w}}(x)[j] - f_w(x)[j] \right| \\ & < \frac{\gamma}{4} + \frac{\gamma}{4} = \frac{\gamma}{2} \end{aligned} \quad (31)$$

Recall that

$$\begin{aligned} \tilde{L}_D(f_w, 0) &= \mathbb{P}_{z \sim D} \left(f_w(x)[y] \leq \max_{j \neq y} f_w(x)[j] \right) \\ \tilde{L}_D(f_{\tilde{w}}, \frac{\gamma}{2}) &= \mathbb{P}_{z \sim D} \left(f_{\tilde{w}}(x)[y] \leq \frac{\gamma}{2} + \max_{j \neq y} f_{\tilde{w}}(x)[j] \right), \end{aligned}$$

Denoting $j_1^* = \arg \max_{j \neq y} f_{\bar{w}}(x)[j]$ and $j_2^* = \arg \max_{j \neq y} f_w(x)[j]$, from Eq. (31), we have

$$\begin{aligned} & \left| f_{\bar{w}}(x)[y] - f_{\bar{w}}(x)[j_2^*] - f_w(x)[y] + f_w(x)[j_2^*] \right| < \frac{\gamma}{2} \\ \Rightarrow & f_{\bar{w}}(x)[y] - f_{\bar{w}}(x)[j_2^*] < f_w(x)[y] - f_w(x)[j_2^*] + \frac{\gamma}{2} \end{aligned} \quad (32)$$

Note that since $f_{\bar{w}}(x)[j_1^*] \geq f_{\bar{w}}(x)[j_2^*]$, we have

$$\begin{aligned} f_{\bar{w}}(x)[y] - f_{\bar{w}}(x)[j_1^*] & \leq f_{\bar{w}}(x)[y] - f_{\bar{w}}(x)[j_2^*] \\ & \leq f_w(x)[y] - f_w(x)[j_2^*] + \frac{\gamma}{2} \end{aligned} \quad (\text{Eq. (32)})$$

Therefore, we have

$$f_w(x)[y] - f_w(x)[j_2^*] \leq 0 \Rightarrow f_{\bar{w}}(x)[y] - f_{\bar{w}}(x)[j_1^*] \leq \frac{\gamma}{2},$$

which indicates $\mathbb{P}_{z \sim D} (f_w(x)[y] \leq f_w(x)[j_2^*]) \leq \mathbb{P}_{z \sim D} (f_{\bar{w}}(x)[y] \leq f_{\bar{w}}(x)[j_1^*] + \frac{\gamma}{2})$, or equivalently

$$\tilde{L}_D(f_w, 0) \leq \tilde{L}_D(f_{\bar{w}}, \frac{\gamma}{2}). \quad (33)$$

Note that this holds for any perturbation $\bar{w} \sim \tilde{Q}$.

Again, recall that

$$\begin{aligned} \tilde{L}_D(f_{\bar{w}}, \frac{\gamma}{2}) & = \mathbb{P}_{z \sim D} \left(f_{\bar{w}}(x)[y] \leq \frac{\gamma}{2} + \max_{j \neq y} f_{\bar{w}}(x)[j] \right) \\ \tilde{L}_D(f_w, \gamma) & = \mathbb{P}_{z \sim D} \left(f_w(x)[y] \leq \gamma + \max_{j \neq y} f_w(x)[j] \right) \end{aligned}$$

From Eq. (31), we have

$$\begin{aligned} & \left| f_{\bar{w}}(x)[y] - f_{\bar{w}}(x)[j_1^*] - f_w(x)[y] + f_w(x)[j_1^*] \right| < \frac{\gamma}{2} \\ \Rightarrow & f_w(x)[y] - f_w(x)[j_1^*] < f_{\bar{w}}(x)[y] - f_{\bar{w}}(x)[j_1^*] + \frac{\gamma}{2} \end{aligned} \quad (34)$$

Note that since $f_w(x)[j_2^*] \geq f_w(x)[j_1^*]$, we have

$$\begin{aligned} f_w(x)[y] - f_w(x)[j_2^*] & \leq f_w(x)[y] - f_w(x)[j_1^*] \\ & \leq f_{\bar{w}}(x)[y] - f_{\bar{w}}(x)[j_1^*] + \frac{\gamma}{2} \end{aligned} \quad (\text{Eq. (34)})$$

Therefore, we have

$$f_{\bar{w}}(x)[y] - f_{\bar{w}}(x)[j_1^*] \leq \frac{\gamma}{2} \Rightarrow f_w(x)[y] - f_w(x)[j_2^*] \leq \gamma,$$

which indicates $\tilde{L}_D(f_{\bar{w}}, \frac{\gamma}{2}) \leq \tilde{L}_D(f_w, \gamma)$. Therefore, from the perspective of the empirical estimation of the probability, for any $\bar{w} \sim \tilde{Q}$, we almost surely have

$$\tilde{L}_S(f_{\bar{w}}, \frac{\gamma}{2}) \leq \tilde{L}_S(f_w, \gamma). \quad (35)$$

Now with probability at least $1 - \delta$, we have

$$\begin{aligned} \tilde{L}_D(f_w, 0) & \leq \mathbb{E}_{\bar{w} \sim \tilde{Q}} \left[\tilde{L}_D(f_{\bar{w}}, \frac{\gamma}{2}) \right] \quad (\text{Eq. (33)}) \\ & \leq \mathbb{E}_{\bar{w} \sim \tilde{Q}} \left[\tilde{L}_S(f_{\bar{w}}, \frac{\gamma}{2}) \right] + \sqrt{\frac{\text{KL}(\tilde{Q} \| P) + \ln \frac{2m}{\delta}}{2(m-1)}} \quad (\text{Theorem 2.5}) \\ & \leq \tilde{L}_S(f_w, \gamma) + \sqrt{\frac{\text{KL}(\tilde{Q} \| P) + \ln \frac{2m}{\delta}}{2(m-1)}} \quad (\text{Eq. (35)}) \end{aligned} \quad (36)$$

Note that

$$\begin{aligned}
\text{KL}(Q\|P) &= \int_{\tilde{w} \in \mathcal{C}} Q \ln \frac{Q}{P} d\tilde{w} + \int_{\tilde{w} \in \bar{\mathcal{C}}} Q \ln \frac{Q}{P} d\tilde{w} \\
&= \int_{\tilde{w} \in \mathcal{C}} \frac{QZ}{Z} \ln \frac{Q}{ZP} d\tilde{w} + \int_{\tilde{w} \in \mathcal{C}} Q \ln Z d\tilde{w} + \int_{\tilde{w} \in \bar{\mathcal{C}}} \frac{Q(1-Z)}{1-Z} \ln \frac{Q}{(1-Z)P} d\tilde{w} + \int_{\tilde{w} \in \bar{\mathcal{C}}} Q \ln(1-Z) d\tilde{w} \\
&= Z \text{KL}(\tilde{Q}\|P) + (1-Z) \text{KL}(\bar{\tilde{Q}}\|P) - H(Z), \tag{37}
\end{aligned}$$

where $\bar{\tilde{Q}}$ denotes the normalized density of Q restricted to $\bar{\mathcal{C}}$. $H(Z)$ is the entropy of a Bernoulli random variable with parameter Z . Since we know $\frac{1}{2} \leq Z \leq 1$ from the beginning, $0 \leq H(Z) \leq \ln 2$, and KL is nonnegative, from Eq. (37), we have

$$\begin{aligned}
\text{KL}(\tilde{Q}\|P) &= \frac{1}{Z} [\text{KL}(Q\|P) + H(Z) - (1-Z) \text{KL}(\bar{\tilde{Q}}\|P)] \\
&\leq \frac{1}{Z} [\text{KL}(Q\|P) + H(Z)] \\
&\leq 2 \text{KL}(Q\|P) + 2 \ln 2. \tag{38}
\end{aligned}$$

Combining Eq. (36) and Eq. (38), we have

$$\tilde{L}_D(f_w, 0) \leq \tilde{L}_S(f_w, \gamma) + \sqrt{\frac{\text{KL}(Q\|P) + \frac{1}{2} \ln \frac{8m}{\delta}}{m-1}}, \tag{39}$$

which finishes the proof. \square

Remark. Note that the constants are slightly different from the one in [5] due to the facts that we use two-side version of Theorem 2.5 and we use natural logarithm rather than the one with base 2.

References

- [1] Ofer Dekel. Pac-bayes analysis. <https://courses.cs.washington.edu/courses/cse522/11wi/scribes/lecture13.pdf>, 2011.
- [2] Wassily Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13–30, 1963.
- [3] David McAllester. Simplified pac-bayesian margin bounds. In *Learning theory and Kernel machines*, pages 203–215. Springer, 2003.
- [4] David A McAllester. Pac-bayesian model averaging. In *Proceedings of the twelfth annual conference on Computational learning theory*, pages 164–170, 1999.
- [5] Behnam Neyshabur, Srinadh Bhojanapalli, and Nathan Srebro. A pac-bayesian approach to spectrally-normalized margin bounds for neural networks. *arXiv preprint arXiv:1707.09564*, 2017.
- [6] Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.
- [7] Leslie G Valiant. A theory of the learnable. *Communications of the ACM*, 27(11):1134–1142, 1984.
- [8] Robert L Wolpert. Markov, chebychev and hoeffding inequalities. <https://www2.stat.duke.edu/courses/Spring09/sta205/lec/hoef.pdf>, 2009.