

Lecture 2 - Concentration of Measure; Generalization in ML

Lecturer: Chicheng Zhang

Scribe: Junfeng Xu

1 Concentration of Measure

- Concentration of Measure: basically provides a way to quantify how close is the sample mean to the population

- Factors: the distribution of the original random variable, sample size, unlucky sample draw
- Example of the concentration quality reduces all three factors above:

Theorem 1. Suppose $X_1 \dots X_n$ are iid, $\mathbb{E}[X_i] = \mu$, if X_i 's are b^2 -SG (all random variables are sub-Gaussian), then

$$P(|\bar{X}_n - \mu| \geq \epsilon) \leq 2\exp\left(-\frac{n\epsilon^2}{2b^2}\right) \quad (1)$$

- Let ϵ to be such that $2\exp\left(-\frac{2n\epsilon^2}{b^2}\right) = \delta \Rightarrow \epsilon = b\sqrt{\frac{\ln \frac{2}{\delta}}{2n}}$
- To prove equation 1: $LHS = P(\bar{X}_n - \mu \geq \epsilon) + P(\bar{X}_n - \mu \leq -\epsilon)$
 $(\square) = P(\bar{X}_n - \mu \geq \epsilon), (\triangle) = P(\bar{X}_n - \mu \leq -\epsilon)$

$$\begin{aligned} (\square) &= P(\bar{X}_n - \mu \geq \epsilon) \\ &= P(\bar{X}_n - \mu \geq \epsilon) \\ &= P\left(\sum_{i=1}^n X_i - n\mu \geq n\epsilon\right) \end{aligned}$$

First we choose a free parameter λ greater than zero and scale both sides by the factor of λ

$$= P\left(\lambda\left(\sum_{i=1}^n X_i - n\mu\right) \geq \lambda n\epsilon, \lambda > 0\right)$$

We need to use sub-Gaussian distribution property here. The sub-gaussianness is about the deviation of the random variable to its mean, but exponentiated. It is nature that we can exponentiate both sides of the equation.

$$= P\left(e^{\lambda\left(\sum_{i=1}^n X_i - n\mu\right)} \geq e^{\lambda n\epsilon}\right)$$

Denote $e^{\lambda\left(\sum_{i=1}^n X_i - n\mu\right)}$ as Z , this random variable Z has non-negativity property, based on Markov inequality, the probability that it deviates is greater than the threshold w cannot be too large if w is already very large. The tail probability will be smaller if my original random variable has a

smaller expectation or my threshold is chosen to be large. i.e. $P(Z \geq w) \leq \frac{Z}{w}$.

$$\begin{aligned} &\leq e^{-\lambda n \epsilon} \mathbb{E}[e^{\lambda(\sum_{i=1}^n X_i - n\mu)}] \\ &= e^{-\lambda n \epsilon} \phi_{\sum_i X_i - n\mu}(\lambda) \\ &= e^{-\lambda n \epsilon} \phi_{\sum_i X_i - \mathbb{E}[\sum_i X_i]}(\lambda) \end{aligned}$$

Is $\sum_i X_i$ a SG random variable? Yes. if so, what is various proxy? nb^2 . $n = 2, x_1 + x_2$ is $2b^2 - SG$.

Given a sub-gaussian random variable, it is scaled by a constant factor. The result is still a sub-Gaussian random variable. These various proxy will be scaled by a factor of the scaling square. If we have two independent random variables, both of which are sub-gaussian, their summation must be sub-gaussian. The variance proxy of the new random variable is the summation over each individual random variable's variance proxy. Then applying the definition of the sub-gaussian:

$$\leq e^{-\lambda n \epsilon + \frac{\lambda^2 n b^2}{2}}$$

for any $\lambda > 0$, choosing λ that minimizes that bound $\Rightarrow \lambda = \frac{\epsilon^2}{b^2}$

$$\Rightarrow (\square) \leq e^{-\frac{n\epsilon^2}{2b^2}}$$

Similarly

$$(\Delta) \leq e^{-\frac{n\epsilon^2}{2b^2}}$$

\Rightarrow Theorem 1

Theorem 2: (Bernstein Inequality) Let $X_1 \dots X_n$ be iid random variables, $|X_i - \mathbb{E}X_i| \leq R$, $\mu = \mathbb{E}X_i$, $\sigma^2 = \text{var}(X_i)$, then for any $\epsilon > 0$:

$$P(|\overline{X}_n - \mu| \geq \epsilon) \leq 2\exp\left(-\frac{n\epsilon^2}{2\sigma^2 + \frac{2}{3}R\epsilon}\right) \quad (2)$$

(no worse, can be much better than hoeffding equality. Because the term $\frac{2}{3}R\epsilon$ can be ignored and the denominator has the actual variance of the random variable rather than the variance proxy of the random variable. Generally for a random variable, the variance of it can be much smaller than the variance proxy which makes this bound significantly better.)

Choosing ϵ so that RHS = δ , choosing a tight enough ϵ , such that RHS $\leq \delta$

$$2\exp\left(-\frac{n\epsilon^2}{2\sigma^2 + \frac{2}{3}R\epsilon}\right) \leq \delta$$

Dividing by two on both sides of the equation and taking a logarithmic on both sides, then moving the denominator.

$$\Leftrightarrow n\epsilon^2 \geq (2\sigma^2 + \frac{2}{3}R\epsilon) \ln \frac{2}{\delta}$$

$X > A + B \Leftrightarrow X \geq 2A$ and $x \geq 2B$. X must be greater than the average of the two, which is $A+B$. Then applying this:

$$\begin{aligned} \Leftrightarrow n\epsilon^2 &\geq 2 * 2\sigma^2 \ln \frac{2}{\delta} \text{ and } n\epsilon^2 \geq 2 * \frac{2}{3} R\epsilon \ln \frac{2}{\delta} \\ \Leftrightarrow \epsilon &\geq \sqrt{\frac{4\sigma^2 \ln \frac{2}{\delta}}{n}} \text{ and } \epsilon \geq \frac{4R \ln \frac{2}{\delta}}{3n} \\ \Leftrightarrow \epsilon &\geq \sqrt{\frac{4\sigma^2 \ln \frac{2}{\delta}}{n}} + \frac{4R \ln \frac{2}{\delta}}{3n} \end{aligned}$$

In summary:

$$P(|\bar{X}_n - \mu| \geq \sqrt{\frac{4\sigma^2 \ln \frac{2}{\delta}}{n}} + \frac{4R \ln \frac{2}{\delta}}{3n}) \leq \delta \quad (3)$$

2 Generalization in supervised machine learning

- Instance space X (e.g. $[0, 1]^{W \times H}$ camera images in pixel representation)
- Label space Y (e.g. $Y = \{L, R\}$)
- Loss function, $\ell(\hat{y}, g) \in [0, B]$, \hat{y} : prediction, g : ground truth label (e.g. $\ell(\hat{y}, g) = I(\hat{y} \neq g)$)
- distribution D cover X, Y ($X, Y \sim D$) (e.g. Camera capture demonstrated by expert steering)
- prediction rule $f : X \rightarrow Y$, quality measure: generalization loss:

$$L_D(f) = \mathbb{E}_{(X,Y) \sim D}[\ell(f(X), Y)] \quad (\text{smaller the better}) \quad (4)$$

- \mathcal{F} : a predictor class to learn \hat{f} from. (e.g. neural network with a fixed architecture)
- can we design a general approach to find a good \hat{f} for any F ? \hat{f} approximately minimizes $L_D(f)$

Idea:

$$\forall f : L_n(f) = \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i) \rightarrow L_D(f) \quad (5)$$

(Concentration of measure)

Algorithm: Empirical risk minimization (ERM)

$$\text{return } \hat{f} = \arg \min_{f \in \mathcal{F}} L_n(f) \quad (6)$$

Theorem: (ERM) Suppose $|F| \leq \infty$, then $\forall \delta > 0$, with probability $1 - \delta$:

$$\forall f \in F \text{ simultaneously, } |L_n(f) - L_D(f)| \leq B \sqrt{\frac{\ln |F|}{2n}} =: \epsilon_n \quad (7)$$

and therefore

$$L_D(\hat{f}) \leq L_D(f^*) + 2\epsilon_n \quad (8)$$

(notation: $f^* = \operatorname{argmin}_{f \in F} L_D(f)$)

Interpretation:

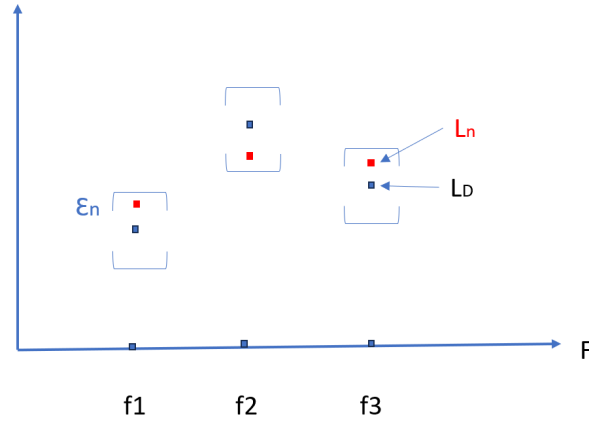


Figure 1: Thm 3 figure

- as long as $n \gg \ln |F|$, \hat{f} 's performance is competitive with the best model in F . To get a sense of how large is $\ln |F|$, suppose F has m parameters, each taking V values, $|F| = V^m \Rightarrow \ln |F| = m \ln V$
- instance space X can be enormous, the training set only has a tiny fraction of all possible instances, yet ERM has strong guarantee \Rightarrow achieves generalization.
- $L_D(f^*)$: approximation error
- $2\epsilon_n$: estimation error
- Larger $F \Rightarrow$ estimation error increases, approximation error decreases.
- F can encode learner's inductive bias, which can help the learning in application-specific ways. For example, for image classification, $X = \text{images}$, $F = \text{convolutional neural network}$, then classifiers in F satisfies translational invariance, that is, $\forall f \in F, f(X) = f(X')$, if X is a translation of X' .
- In modern deep learning regime, $\ln |F| \gg n$ is more common. The guarantees provided by this theorem is vacuous. Nevertheless, researchers found that popular learners (e.g. stochastic gradient-based learners) manage to converge to "simple" predictors, which implicitly uses a much smaller \mathcal{F} . Note: "Implicit Regularization" by Nati Srebro is a great video to watch.

Proof: (7) \Rightarrow (8) refer to figure 1.

$$\begin{aligned}
 L_D(\hat{f}) &\leq L_n(\hat{f}) + \epsilon_n && \hat{f}'\text{s training and test loss are within } \epsilon_n \\
 &\leq L_n(f^*) + \epsilon_n && \hat{f} \text{ is ERM} \\
 &\leq (L_D(f^*) + \epsilon_n) + \epsilon_n && f^* \text{ is training in the test loss are within } \epsilon_n
 \end{aligned}$$

Proving (7): want to show: $P(E) \geq 1 - \delta$, equivalently, we want to show a statement like:

$$\begin{aligned}
 &P(\forall f \in F : |L_n(f) - L_D(f)| \leq \epsilon) \geq 1 - \delta \\
 \Leftrightarrow &P(\exists f \in F : |L_n(f) - L_D(f)| > \epsilon) \leq \delta
 \end{aligned}$$

$$\begin{aligned}
LHS &= P\left(\bigcup_{f \in F} |L_n(f) - L_D(f)| > \epsilon\right) \\
&\leq \sum_{f \in F} P(|L_n(f) - L_D(f)| > \epsilon) \\
&\leq \exp\left(-\frac{2n\epsilon^2}{B^2}\right) \\
&= |F| \exp\left(-\frac{2n\epsilon^2}{B^2}\right)
\end{aligned}$$

Setting ϵ such that this bounds $\delta \Rightarrow 2|F|\exp\left(-\frac{2n\epsilon^2}{B^2}\right) = \delta \Rightarrow \epsilon = B\sqrt{\frac{\ln \frac{2|F|}{\delta}}{2n}} = \epsilon_n$

$$\Rightarrow P(\exists f : |L_n(f) - L_D(f)| > \epsilon_n) \leq \delta \tag{9}$$

Follow up question: what would we get for ϵ_n if we instead using Chebyshev's inequality for bounding the deviation probability? (Hint: $\ln \frac{|F|}{\delta}$ will become without $\frac{|F|}{\delta}$, can you see why?)