

Generalization Bounds for Averaged Classifiers

Yoav Freund, Yishay Mansour and Robert Schapire, *Annals of Statistics*, 2004

Feb. 6th, 2014
presented by Chicheng Zhang chichengzhang@ucsd.edu

Outline

Generalization Bounds for Averaged Classifiers

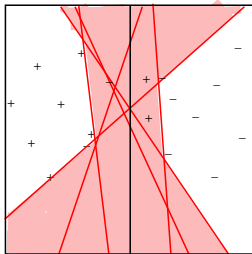
- Classification with Reject Option: an Introduction
- Algorithm
- Analysis

Concentration Bounds for Sequential Prediction

- Online Learning with Stochastic Data: an Introduction
- General Techniques
- Relative Entropy Inequalities
- Bennett Inequalities
- Applications to Exponential Weight Algorithm

Introduction

- ▶ Overfitting
 - ▶ A classification problem w/ insufficient training data
 - ▶ ERM: w.p. $1 - \delta$, $\epsilon(\hat{h}) \leq \epsilon(h^*) + \sqrt{\epsilon(h^*) \frac{\ln |\mathcal{H}| + \ln \frac{1}{\delta}}{m}} + \frac{\ln |\mathcal{H}| + \ln \frac{1}{\delta}}{m}$
- ▶ Avoid: Saying Don't Know on part of test examples
 - ▶ Agnostic selective classification [EYW11]: Version Space, Disagreement-Based
 - ▶ This work: weighted average over \mathcal{H} , taking voting "margin" into account



- ▶ Generally, $\Pr(\text{Don't Know}) \uparrow$, $\Pr(\text{Mistake}) \downarrow$.
- ▶ Goal: find an alg. such that both $\Pr(\text{Don't Know})$ and $\Pr(\text{Mistake})$ small.

Preliminaries

- ▶ Batch Learning(Rather than online learning!), distribution \mathcal{D} defined over $\mathcal{X} \times \{-1, +1\}$, $(x_1, y_1), \dots, (x_m, y_m) \sim \mathcal{D}$.
- ▶ a hypothesis class \mathcal{H} , Each classifier $h \in \mathcal{H}$,
 $h : \mathcal{X} \rightarrow \{-1, +1\}$
- ▶ True error $\epsilon(h) = \mathbb{P}(h(X) \neq Y)$, error of the optimal classifier $\epsilon = \epsilon(h^*) = \inf_{h \in \mathcal{H}} \epsilon(h)$, Empirical error $\hat{\epsilon}(h) = \frac{1}{m} \sum_{i=1}^m l(h(x_i) \neq y_i)$.
- ▶ $\mathcal{H}_x^+ = \{h \in \mathcal{H} : h(x) = +1\}$, $\mathcal{H}_x^- = \{h \in \mathcal{H} : h(x) = -1\}$

Algorithm: Intuition

- ▶ Inspired by Exponential Weight algorithm(Hedge, say?)
"Bayesian" $w_{i,t+1} \propto e^{-\eta L_{i,t}}$.
- ▶ Translated into binary classification in batch case:
 $w(h) \propto e^{-\eta \hat{\epsilon}(h)}$
- ▶ "Softly" Put higher weight over classifiers performing well.
- ▶ Algorithm:

- ▶ $\hat{\ell}(x) = \frac{1}{\eta} \ln\left(\frac{\sum_{h(x)=+} e^{-\eta \hat{\epsilon}(h)}}{\sum_{h(x)=-} e^{-\eta \hat{\epsilon}(h)}}\right)$
- ▶ if $|\hat{\ell}(x)| \leq \Delta$, then predict 0 (Saying Don't Know).
- ▶ otherwise, predict w/ $\text{sign}(\hat{\ell}(x))$.

- ▶ How to choose Δ ? Will analyze a related quantity

$$\ell(x) = \frac{1}{\eta} \ln\left(\frac{\sum_{h(x)=+} e^{-\eta \epsilon(h)}}{\sum_{h(x)=-} e^{-\eta \epsilon(h)}}\right) \text{ first.}$$

Original Theorem and Proof

Theorem

Let $\eta > 0$, $\Delta \geq 0$, $\Delta\eta \leq 1/2$. Then $\forall \gamma \geq \frac{\ln 8|\mathcal{H}|}{\eta}$,

$$\Pr(y_\ell(x) \leq 0) \leq 2(1 + 2|\mathcal{H}|e^{-\eta\gamma})(\epsilon + \gamma)$$

$$\begin{aligned}\Pr(y_\ell(x) \leq 2\Delta) &\leq (1 + e^{2\Delta\eta})(1 + 2|\mathcal{H}|e^{\eta(2\Delta-\gamma)})(\epsilon + \gamma) \\ &\leq 4(1 + 2|\mathcal{H}|e^{\eta(2\Delta-\gamma)})(\epsilon + \gamma)\end{aligned}$$

$\hat{\ell}(x)$ converges to $\ell(x)$

Theorem

For any \mathcal{D} , any $x \in \mathcal{X}$, any $\lambda, \eta > 0$:

$$\Pr_{S \sim \mathcal{D}^m} (|\ell(x) - \hat{\ell}(x)| \geq 2\lambda + \frac{\eta}{8m}) \leq 4e^{-2m\lambda^2}$$

- ▶ Intuitively, $\eta \uparrow$, the convergence become worse.
- ▶ Note the convergence rate does not depend on $|\mathcal{H}|$!
- ▶ Define

$$\hat{R}_\eta(\mathcal{K}) = \frac{1}{\eta} \ln \left(\sum_{h \in \mathcal{K}} e^{-\eta \hat{\epsilon}(h)} \right), R_\eta(\mathcal{K}) = \frac{1}{\eta} \ln \left(\sum_{h \in \mathcal{K}} e^{-\eta \epsilon(h)} \right)$$

- ▶ Note $\hat{\ell}(x) = \hat{R}_\eta(\mathcal{H}_x^+) - \hat{R}_\eta(\mathcal{H}_x^-)$, $\ell(x) = R_\eta(\mathcal{H}_x^+) - R_\eta(\mathcal{H}_x^-)$
- ▶ We will prove it based on convergence of $\hat{R}_\eta(\mathcal{H}_x^+)$ to $R_\eta(\mathcal{H}_x^+)$, and $\hat{R}_\eta(\mathcal{H}_x^-)$ to $R_\eta(\mathcal{H}_x^-)$

Proof I

- ▶ Define "weak" classifiers: $\text{weak} = \{h \in \mathcal{H} : \epsilon \geq \epsilon + \gamma\}$, otherwise, call them "strong".
- ▶ Fix (x, y) , weight of each group of classifiers:

$$W_s^\vee(x, y) = \frac{\sum_{h(x)=y, \epsilon(h) < \epsilon + \gamma} e^{-\eta \epsilon(h)}}{\sum_{\mathcal{H}} e^{-\eta \epsilon(h)}}, \quad W_s^X(x, y) = \frac{\sum_{h(x) \neq y, \epsilon(h) < \epsilon + \gamma} e^{-\eta \epsilon(h)}}{\sum_{\mathcal{H}} e^{-\eta \epsilon(h)}}$$

$$W_w = \frac{\sum_{\epsilon(h) \geq \epsilon + \gamma} e^{-\eta \epsilon(h)}}{\sum_{\mathcal{H}} e^{-\eta \epsilon(h)}} \leq \frac{|\mathcal{H}| e^{-\eta(\gamma + \epsilon)}}{e^{-\eta \gamma}} = |\mathcal{H}| e^{-\eta \gamma} \leq \frac{1}{8}$$

- ▶ When $y \ell(x) \leq 2\Delta$,

$$2\Delta \geq \frac{1}{\eta} \ln \frac{W_s^\vee(x, y)}{W_s^X(x, y)} \geq \frac{1}{\eta} \ln \frac{W_s^\vee(x, y)}{W_s^X(x, y) + W_w}$$

Hence

$$W_s^X(x, y) + W_w \geq \frac{1}{1 + e^{2\Delta \eta}} =: c \Rightarrow \frac{W_s^X(x, y)}{W_s^X(x, y) + W_s^\vee(x, y)} \geq \frac{c - W_w}{1 - W_w}$$

Proof II

► Analogously, we have

$$\begin{aligned}\Pr(y \ell(x) \leq 2\Delta) &\leq \Pr\left(\frac{W_s^x(x, y)}{W_s^x(x, y) + W_s^y(x, y)} \geq \frac{c - W_w}{1 - W_w}\right) \\&\leq \Pr_{(x, y) \sim \mathcal{D}} \left(\Pr_{h \sim w|s} (h(x) \neq y) \geq \frac{c - W_w}{1 - W_w} \right) \\&\leq \mathbb{E}_{(x, y) \sim \mathcal{D}} \Pr_{h \sim w|s} (h(x) \neq y) \frac{1 - W_w}{c - W_w} \\&= \mathbb{E}_{(x, y) \sim \mathcal{D}} \mathbb{E}_{h \sim w|s} I(h(x) \neq y) \frac{1 - W_w}{c - W_w} \\&= \mathbb{E}_{h \sim w|s} \mathbb{E}_{(x, y) \sim \mathcal{D}} I(h(x) \neq y) \frac{1 - W_w}{c - W_w} \\&\leq (\epsilon + \gamma) \frac{1 - W_w}{c - W_w} \\&\leq (\epsilon + \gamma)(1 + 2W_w e^{2\Delta\eta})(1 + e^{2\Delta\eta}) \\&\leq (1 + e^{2\Delta\eta})(1 + 2|\mathcal{H}|e^{\eta(2\Delta - \gamma)})(\epsilon + \gamma)\end{aligned}$$

Performance of $\ell(x)$

Intuitively, $\eta \uparrow$, the performance of $\ell(x)$ becomes better.

Theorem

Gets closer to the performance of the best h in H .

For any $\Delta \geq 0$, we have:

Is this theorem in the paper? What is the intuition?

$$\Pr(y\ell(x) \leq 2\Delta) \leq (1 + e^{2\eta\Delta})(\epsilon + \frac{\ln |\mathcal{H}|}{\eta})$$

In particular, let $\Delta = 0$, we have

$$\Pr(y\ell(x) \leq 0) \leq 2(\epsilon + \frac{\ln |\mathcal{H}|}{\eta})$$

The factor 2 is unavoidable for such voting methods.

Proof I

- Fix an example (x, y) , "correct" and "incorrect" weight

$$W_s^\vee(x, y) = \frac{\sum_{h(x)=y} e^{-\eta\epsilon(h)}}{\sum_{\mathcal{H}} e^{-\eta\epsilon(h)}}, W_s^X(x, y) = \frac{\sum_{h(x)\neq y} e^{-\eta\epsilon(h)}}{\sum_{\mathcal{H}} e^{-\eta\epsilon(h)}}$$

Markov's Inequality:

$$\begin{aligned} I(y \ell(x) \leq 2\Delta) &= I(W_s^\vee(x, y) \leq e^{2\eta\Delta} W_s^X(x, y)) \\ &= I(W_s^\vee(x, y) \geq \frac{1}{1 + e^{2\eta\Delta}}) \\ &\leq (1 + e^{2\eta\Delta}) W_s^\vee(x, y) \\ &= (1 + e^{2\eta\Delta}) \frac{\sum_{h(x)\neq y} e^{-\eta\epsilon(h)}}{\sum_h e^{-\eta\epsilon(h)}} \end{aligned}$$

Proof II

Taking expectations on both sides:

$$Pr(y\ell(x) \leq 2\Delta) \leq (1 + e^{2\eta\Delta}) \frac{\sum_h \epsilon(h) e^{-\eta\epsilon(h)}}{\sum_h e^{-\eta\epsilon(h)}}$$

- Convexity of $x \ln x$, $x > 0$:

$$\mathbb{E}X \ln X \geq \mathbb{E}X \ln(\mathbb{E}X)$$

- Take $X(h) = e^{-\eta\epsilon(h)}$, uniform distribution over \mathcal{H} :

$$\frac{\sum_h \epsilon(h) e^{-\eta\epsilon(h)}}{\sum_h e^{-\eta\epsilon(h)}} \leq -\frac{1}{\eta} \ln\left(\frac{1}{|\mathcal{H}|} \sum_h e^{-\eta\epsilon(h)}\right)$$

- Familiar "singleton" bound: $\leq \epsilon + \frac{\ln |\mathcal{H}|}{\eta}$.

Proof (1): $\hat{R}_\eta(\mathcal{K})$ converges to $\mathbb{E}\hat{R}_\eta(\mathcal{K})$

► Note

$$\hat{R}_\eta(\mathcal{K})((x_1, y_1), \dots, (x_m, y_m)) = \frac{1}{\eta} \ln \left(\sum_{h \in \mathcal{K}} e^{-\eta \hat{\epsilon}(h)} \right)$$

satisfies bounded difference

- Suppose we have modified training set $S' = (S \setminus \{(x_i, y_1)\}) \cup \{(x'_i, y'_i)\}$. denote $\hat{\epsilon}'(h)$ the empirical error of h in S' . Then $\sup_{h \in \mathcal{K}} |\hat{\epsilon}'(h) - \hat{\epsilon}(h)| \leq \frac{1}{m}$.

$$-\frac{1}{m} \leq \frac{1}{\eta} \inf_{h \in \mathcal{K}} \ln \left(\frac{e^{-\eta \hat{\epsilon}(h)}}{e^{-\eta \hat{\epsilon}'(h)}} \right) \leq \frac{1}{\eta} \ln \left(\frac{\sum_{h \in \mathcal{K}} e^{-\eta \hat{\epsilon}(h)}}{\sum_{h \in \mathcal{K}} e^{-\eta \hat{\epsilon}'(h)}} \right) \leq \frac{1}{\eta} \sup_{h \in \mathcal{K}} \ln \left(\frac{e^{-\eta \hat{\epsilon}(h)}}{e^{-\eta \hat{\epsilon}'(h)}} \right) \leq \frac{1}{m}$$

- By McDiarmid's Lemma, w.p. $1 - 2e^{-2m\lambda^2}$

$$|\hat{R}_\eta(\mathcal{K}) - \mathbb{E}\hat{R}_\eta(\mathcal{K})| \leq \lambda$$

- How does $\mathbb{E}\hat{R}_\eta(\mathcal{K})$ relate to $R_\eta(\mathcal{K})$?

Proof (2): $\mathbb{E}\hat{R}_\eta(\mathcal{K})$ converges to $R_\eta(\mathcal{K})$

Lemma

$$R_\eta(\mathcal{K}) \leq \mathbb{E}\hat{R}_\eta(\mathcal{K}) \leq R_\eta(\mathcal{K}) + \frac{\eta}{8m}$$

- ▶ The first inequality directly follows from convexity of $f(x) = \ln \sum_i e^{x_i}$
- ▶ The second uses Hoeffding's Inequality:
 $X \in [a, b] \Rightarrow \mathbb{E}e^X \leq e^{\mathbb{E}X} e^{(b-a)^2/8}.$

$$\begin{aligned}\mathbb{E}\hat{R}_\eta(\mathcal{K}) &= \mathbb{E} \frac{1}{\eta} \ln \left(\sum_{h \in \mathcal{K}} e^{-\eta \hat{\epsilon}(h)} \right) \\ &\leq \frac{1}{\eta} \ln \left(\sum_{h \in \mathcal{K}} \mathbb{E} e^{-\eta \hat{\epsilon}(h)} \right) \\ &\leq \frac{1}{\eta} \ln \left(\sum_{h \in \mathcal{K}} \left(e^{-\frac{\eta}{m} \epsilon(h)} e^{\frac{\eta^2}{8m^2}} \right)^m \right) \\ &\leq \frac{1}{\eta} \ln \left(\sum_{h \in \mathcal{K}} e^{-\eta \epsilon(h)} \right) + \frac{\eta}{8m}\end{aligned}$$

Proof (3): Combine \mathcal{H}_x^+ and \mathcal{H}_x^-

- ▶ w.p. $1 - 4e^{-2m\lambda^2}$ the following hold simultaneously:

$$\hat{R}_\eta(\mathcal{H}_x^+) \leq \mathbb{E}\hat{R}_\eta(\mathcal{H}_x^+) + \lambda \leq R_\eta(\mathcal{H}_x^+) + \lambda + \frac{\eta}{8m}$$

$$\hat{R}_\eta(\mathcal{H}_x^-) \geq \mathbb{E}\hat{R}_\eta(\mathcal{H}_x^-) - \lambda \geq R_\eta(\mathcal{H}_x^-) - \lambda$$

- ▶ Hence $\hat{\ell}(x) \leq \ell(x) + 2\lambda + \frac{\eta}{8m}$
- ▶ Analogously, $-\hat{\ell}(x) \leq -\ell(x) + 2\lambda + \frac{\eta}{8m}$
- ▶ Proof generalized into uncountably infinite \mathcal{H} , with technicalities resolved in paper

Bounding the fraction of "atypical" test examples

Theorem

For any $\delta > 0$ and $\eta > 0$, if we set $\Delta = 2\sqrt{\frac{\ln(2/\delta)}{m}} + \frac{\eta}{8m}$, then w.p. $1 - \delta$ over choice of S ,

$$\Pr_{(x,y) \sim \mathcal{D}} (|\ell(x) - \hat{\ell}(x)| \geq \Delta) \leq \delta$$

Note that setting e.g. $\delta = O(m^{-10})$ won't affect much of the bound

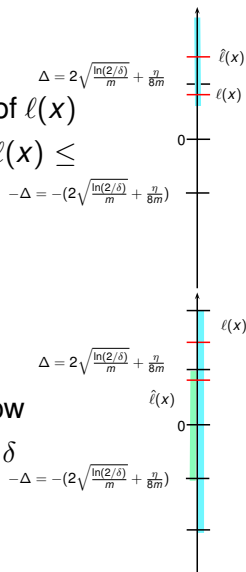
Proof.

Taking $\lambda = \ln(2/\delta)$ using the previous theorem,

$$\begin{aligned} & \mathbb{E}_{S \sim \mathcal{D}^m} \Pr_{(x,y) \sim \mathcal{D}} (|\ell(x) - \hat{\ell}(x)| \geq \Delta) \\ & \leq \mathbb{E}_{S \sim \mathcal{D}^m} \mathbb{E}_{(x,y) \sim \mathcal{D}} I(|\ell(x) - \hat{\ell}(x)| \geq \Delta) \\ & \leq \mathbb{E}_{(x,y) \sim \mathcal{D}} \mathbb{E}_{S \sim \mathcal{D}^m} I(|\ell(x) - \hat{\ell}(x)| \geq \Delta) \\ & \leq \delta^2 \end{aligned}$$

Implication of Deviation

- ▶ Relate the error of $\hat{\ell}(x)$ to the error of $\ell(x)$
- ▶ $\Pr(|\hat{\ell}(x)| > \Delta \wedge y\hat{\ell}(x) \leq 0) \leq \Pr(y\ell(x) \leq 0) + \delta$



- ▶ The probability of saying Don't Know
- ▶ $\Pr(|\hat{\ell}(x)| \leq \Delta) \leq \Pr(|\ell(x)| \leq 2\Delta) + \delta$

Putting them together

► Mistake Bound

$$\begin{aligned}\Pr(\text{Mistake}) &= \Pr(|\hat{\ell}(x)| > \Delta \wedge y\hat{\ell}(x) \leq 0) \\ &\leq \Pr(y\ell(x) \leq 0) + \delta \\ &\leq 2\left(\epsilon + \frac{\ln |\mathcal{H}|}{\eta}\right) + \delta\end{aligned}$$

► Don't know Bound

$$\begin{aligned}\Pr(\text{Don't Know}) &= \Pr(|\hat{\ell}(x)| \leq \Delta) \\ &\leq \Pr(|\ell(x)| \leq 2\Delta) + \delta \\ &\leq \Pr(y\ell(x) \leq 2\Delta) + \delta \\ &\leq (1 + e^{2\Delta\eta})\left(\epsilon + \frac{\ln |\mathcal{H}|}{\eta}\right) + \delta\end{aligned}$$

Putting them together

- ▶ Simple tuning($\eta = \ln |\mathcal{H}| m^{1/2}$):

$$\Pr(\text{Mistake}) \leq 2(\epsilon + m^{-1/2}) + \delta$$

$$\Pr(\text{Don't Know}) \leq e^{\sqrt{\ln \frac{1}{\delta}} \ln |\mathcal{H}| + (\ln |\mathcal{H}|)^2} (\epsilon + m^{-1/2}) + \delta$$

- ▶ In region of prediction, the probability of making a mistake is independent of complexity of \mathcal{H} any more!
- ▶ The upper bound of probability of saying Don't Know might be loose; should be estimated by unlabelled data in practice.

Putting them together(2)

- ▶ A subtler tuning($\eta = \ln |\mathcal{H}| m^{1/2-\theta}$):

$$\Pr(\text{Mistake}) \leq 2(\epsilon + m^{-1/2+\theta}) + \delta$$

$$\Pr(\text{Don't Know}) \leq (1 + e^{2\frac{\sqrt{\ln 1/\delta} \ln |\mathcal{H}|}{m^\theta} + \frac{(\ln |\mathcal{H}|)^2}{m^{2\theta}}})(\epsilon + m^{-1/2+\theta}) + \delta$$

- ▶ When $m \leq O((\ln |\mathcal{H}| + \ln(1/\delta))^{1/\theta})$, the mistake bound improves over ERM.
- ▶ OTOH, $m \geq \Omega((\ln |\mathcal{H}| \ln(1/\delta))^{1/\theta})$,

$$2\frac{\sqrt{\ln 1/\delta} \ln |\mathcal{H}|}{m^\theta} + \frac{(\ln |\mathcal{H}|)^2}{m^{2\theta}} \leq 1$$

$\Pr(\text{Don't Know}) \leq 5(\epsilon + m^{-1/2+\theta})$, almost as small as the error guarantee for ERM.

Generalization to uncountably infinite hypothesis class

- ▶ Have a "prior" μ (hopefully) puts more weights for "good" classifiers
- ▶ define $\ell(x)$ slightly differently: $\ell(x) = \frac{\int_{h(x)=+} e^{-\eta\epsilon(h)} d\mu}{\int_{h(x)=-} e^{-\eta\epsilon(h)} d\mu}$. Then same argument goes:

$$Pr(y\ell(x) \leq \Delta) \leq (1 + e^{\eta\Delta}) \leq (1 + e^{\eta\Delta}) \left(-\frac{1}{\eta} \ln \int e^{-\eta\epsilon(h)} d\mu(h) \right)$$

- ▶ Applying Compression Lemma:

$$-\frac{1}{\eta} \ln \int e^{-\eta\epsilon(h)} d\mu(h) \leq \int \epsilon(h) d\nu(h) + \frac{D(\nu||\mu)}{\eta}$$

Where $D(\nu||\mu) = \int \ln \frac{d\nu}{d\mu} d\nu$ is the relative entropy.

- ▶ Taking $V_\epsilon = \int_{\epsilon(h) \leq \epsilon} d\mu(h)$, $d\nu(h) = \frac{I(\epsilon(h) \leq \epsilon)}{V_\epsilon} d\mu(h)$:

$$-\frac{1}{\eta} \ln \int e^{-\eta\epsilon(h)} d\mu(h) \leq \epsilon + \frac{\ln 1/V_\epsilon}{\eta}$$

Aside: Proof of Compression Lemma

Lemma

If μ, ν are two probability measures, $\nu \ll \mu$, then

$$\int f(h) d\nu(h) \leq D(\nu || \mu) + \ln\left(\int e^{f(h)} d\mu(h)\right)$$

Proof.

Define a new probability measure $d\hat{\mu}(h) = e^{f(h)} d\mu(h) / Z$, $Z = \int e^{f(h)} d\mu(h)$. Since $D(\nu || \hat{\mu}) \geq 0$, expanding,

$$\int \ln\left(\frac{d\nu Z}{d\mu e^{f(h)}}\right) d\nu \geq 0$$

i.e.

$$\ln Z + D(\nu || \mu) \geq \int f(h) d\nu(h)$$



Online Learning with Stochastic Data

- Perceptron Algorithm:

- $w_1 = 0$;

For $t = 1, 2, \dots, m$:

Observing x_t , predicts $\hat{y}_t = \text{sign}(w_t \cdot x_t) =: h_t(x_t)$.

Receive y_t , incur loss $l(y_t \neq \text{sign}(w_t \cdot x_t))$.

Update w_{t+1} based on $w_t, (x_t, y_t)$.

- Mistake Bound: Suppose $\|X\|_2 \leq X$,

$$\begin{aligned} & \sum_{t=1}^m l(y_t \neq \text{sign}(w_t \cdot x_t)) \\ & \leq \inf_{u, \gamma > 0, \lambda > 0} \left(\left(1 + \frac{1}{\lambda}\right) \left(\sum_{t=1}^m \left(1 - \frac{y_t u^T x_t}{\gamma}\right)_+\right) + (1 + \lambda) \frac{X^2 \|u\|^2}{\gamma^2} \right) \end{aligned}$$

- it holds for arbitrary sequence, hence for stochastic sequence as well
- can we relate $\{\epsilon(h_t)\}_{t=1}^m$ to $\{l(h_t(x_t) \neq y_t)\}_{t=1}^m$? Also answered in [CBG05, CBCG04].

A Basic Inequality

- ▶ Consider a sequence of iid random variables Z_1, \dots, Z_m , and functions $\xi_1(z_1), \xi_2(z_1, z_2), \dots, \xi_m(z_1, z_2, \dots, z_m)$.
- ▶ Specifically in our context: $z_t = (x_t, y_t)$,
 $\xi_t(z_1, z_2, \dots, z_t) = I(h_t(x_t) \neq y_t)$, h_t depends only on $(x_1, y_1), \dots, (x_{t-1}, y_{t-1})$
- ▶ Notice $\mathbb{E}_{Z_t} \xi_t(Z_1, Z_2, \dots, Z_t) = \epsilon(h_t)$. Find the relationship between $\{\mathbb{E}_{Z_t} \xi_t\}_{t=1}^m$ and $\{\xi_t\}_{t=1}^m$. Specifically,
 $\mu_m = \frac{1}{m} \sum_{t=1}^m \mathbb{E}_{Z_t} \xi_t$, $s_m = \frac{1}{m} \sum_{t=1}^m \xi_t$

Lemma

For any functions $\zeta_1(x_1), \dots, \zeta_n(x_1, \dots, x_n)$

$$\mathbb{E} \frac{e^{\zeta_1(Z_1)}}{\mathbb{E}_{Z_1} e^{\zeta_1(Z_1)}} \frac{e^{\zeta_2(Z_1, Z_2)}}{\mathbb{E}_{Z_2} e^{\zeta_2(Z_1, Z_2)}} \cdots \frac{e^{\zeta_n(Z_1, \dots, Z_n)}}{\mathbb{E}_{Z_n} e^{\zeta_n(Z_1, \dots, Z_n)}} = 1$$

Hence by Markov's Inequality, w.p. $1 - \delta$,

$$\begin{aligned} & \zeta_1(Z_1) + \zeta_2(Z_1, Z_2) + \dots + \zeta_n(Z_1, \dots, Z_n) \\ \leq & \ln \mathbb{E}_{Z_1} e^{\zeta_1(Z_1)} + \ln \mathbb{E}_{Z_2} e^{\zeta_2(Z_1, Z_2)} + \dots + \ln \mathbb{E}_{Z_n} e^{\zeta_n(Z_1, \dots, Z_n)} + \ln \frac{1}{\delta} \end{aligned}$$

Relative Entropy Inequalities I

- ▶ Assume $\xi_i \in [0, 1]$
- ▶ Taking $\zeta_i = -\rho\xi_i$, rearranging,

$$\begin{aligned} & \rho\xi_1 + \rho\xi_2 + \dots + \rho\xi_n + \ln \frac{1}{\delta} \\ \geq & -\ln \mathbb{E}_{Z_1} e^{-\rho\xi_1} - \ln \mathbb{E}_{Z_2} e^{-\rho\xi_2} - \dots - \ln \mathbb{E}_{Z_n} e^{-\rho\xi_n} \\ \geq & \sum_{t=1}^n -\ln(1 - (1 - e^{-\rho})\mathbb{E}_{Z_t}\xi_t) \\ \geq & -n\ln(1 - (1 - e^{-\rho})\mu_n) \end{aligned}$$

Equivalent to $-\rho s_n - \ln(1 - (1 - e^{-\rho})\mu_n) \leq \frac{\ln 1/\delta}{n}$

- ▶ Since $D(q||p) = \sup_{\rho \geq 0} (-\rho q - \ln(1 - (1 - e^{-\rho})q))$, $p \geq q$, after some manipulations: $\forall \alpha \in [0, 1], t \geq 0$,

$$\Pr(\mu_n \geq D^{-1}(\alpha, t), s_n \leq \alpha) \leq e^{-nt}$$

Relative Entropy Inequalities II

- ▶ Note: ρ cannot be tuned apriori!

- ▶ Taking union bound over all

$$(\alpha, t) \in \{(0, t_0 + \frac{2 \ln 2}{n}), (\frac{1}{n}, t_0 + \frac{2 \ln(3)}{n}), \dots, (1, t_0 + \frac{2 \ln(n+2)}{n})\}:$$

$$\Pr(\mu_n \geq D^{-1}(\frac{\lceil ns_n \rceil}{n}, \frac{2 \ln(\lceil ns_n \rceil + 2)}{n} + t)) \leq e^{-nt}$$

- ▶ Since $D^{-1}(\alpha, t) \leq \alpha + 2t + \sqrt{2\alpha t}$, we get w.p. $1 - \delta$

$$\mu_n \leq \frac{\lceil ns_n \rceil}{n} + 2 \frac{2 \ln(\lceil ns_n \rceil + 2) + \ln \frac{1}{\delta}}{n} + \sqrt{2(\frac{\lceil ns_n \rceil}{n} \frac{2 \ln(\lceil ns_n \rceil + 2) + \ln \frac{1}{\delta}}{n})}$$

- ▶ Simplifying:

$$\mu_n \leq s_n + O(\frac{\ln n + \ln(1/\delta)}{n} + \sqrt{s_n \frac{\ln n + \ln(1/\delta)}{n}})$$

Relative Entropy Inequalities III

- Specifically for perceptron:

$$\mu_n \leq \frac{L_{\gamma,n}}{n\gamma} + O\left(\frac{X^2\|u\|^2}{n\gamma^2} + \frac{\ln n + \ln(1/\delta)}{n}\right) + \sqrt{\frac{L_{\gamma,n}}{n\gamma} \frac{X^2\|u\|^2}{n\gamma^2} + \frac{\ln n + \ln(1/\delta)}{n}}$$

Bennett Inequalities I

- ▶ Assume $\xi_i - \mathbb{E}_{Z_i}\xi_i \geq -1$
- ▶ Taking $\zeta_i = -\rho(\xi_i - \mathbb{E}_{Z_i}\xi_i)$, rearranging,

$$\begin{aligned} & \rho(\xi_1 - \mathbb{E}_{Z_1}\xi_1) + \rho(\xi_2 - \mathbb{E}_{Z_2}\xi_2) + \dots + \rho(\xi_n - \mathbb{E}_{Z_n}\xi_n) + \ln \frac{1}{\delta} \\ & \geq -\ln \mathbb{E}_{Z_1} e^{-\rho(\xi_1 - \mathbb{E}_{Z_1}\xi_1)} - \dots - \ln \mathbb{E}_{Z_n} e^{-\rho(\xi_n - \mathbb{E}_{Z_n}\xi_n)} \\ & \geq -(e^\rho - \rho - 1)(\mathbb{E}_{Z_1}(\xi_1 - \mathbb{E}_{Z_1}\xi_1)^2 + \dots + \mathbb{E}_{Z_n}(\xi_n - \mathbb{E}_{Z_n}\xi_n)^2) \end{aligned}$$

- ▶ Imposing the assumption that $\mathbb{E}_{Z_t}(\xi_t - \mathbb{E}_{Z_t}\xi_t)^2 \leq b\mathbb{E}_{Z_t}\xi_t$ (e.g. $\mathbb{E}_{(x_t, y_t)}(I(h(x_t) \neq y_t) - \epsilon(h_t))^2 \leq \epsilon(h_t) - \epsilon(h_t)^2$)

$$\begin{aligned} & (\mathbb{E}_{Z_1}\xi_1 + \mathbb{E}_{Z_2}\xi_2 + \dots + \mathbb{E}_{Z_n}\xi_n) - (\xi_1 + \xi_2 + \dots + \xi_n) \\ \leq & \frac{\ln(1/\delta)}{\rho} + \frac{\rho}{2(1 - \rho/3)}(\mathbb{E}_{Z_1}(\xi_1 - \mathbb{E}_{Z_1}\xi_1)^2 + \dots + \mathbb{E}_{Z_n}(\xi_n - \mathbb{E}_{Z_n}\xi_n)^2) \\ \leq & \frac{\ln(1/\delta)}{\rho} + \frac{b\rho}{2(1 - \rho/3)}(\mathbb{E}_{Z_1}\xi_1 + \mathbb{E}_{Z_2}\xi_2 + \dots + \mathbb{E}_{Z_n}\xi_n) \end{aligned}$$

Bennett Inequalities II

- ▶ after some manipulations, $\forall \alpha > 0, \exists c_b$

$$\Pr(\mu_n \geq s_n + \sqrt{2b\alpha t} + c_b t, s_n \leq \alpha) \leq e^{-nt}$$

- ▶ Same as before, taking union bound over all $(\alpha, t) \in \{(0, t_0 + \frac{2\ln 2}{n}), (\frac{1}{n}, t_0 + \frac{2\ln(3)}{n}), \dots, (1, t_0 + \frac{2\ln(n+2)}{n})\}$, w.p. $1 - \delta$:

$$\mu_n \leq s_n + c_b \frac{2\ln(\lceil ns_n \rceil + 2) + \ln \frac{1}{\delta}}{n} + \sqrt{2b \left(\frac{\lceil ns_n \rceil}{n} \frac{2\ln(\lceil ns_n \rceil + 2) + \ln \frac{1}{\delta}}{n} \right)}$$

- ▶ Similar bounds obtained for perceptron

Applications to Exponential Weight Algorithm I

► Reminder:

$$\begin{aligned} & \rho\xi_1 + \rho\xi_2 + \dots + \rho\xi_n + \ln \frac{1}{\delta} \\ \geq & -\ln \mathbb{E}_{Z_1} e^{-\rho\xi_1} - \ln \mathbb{E}_{Z_2} e^{-\rho\xi_2} - \dots - \ln \mathbb{E}_{Z_n} e^{-\rho\xi_n} \end{aligned}$$

- Consider applying Hedge(η) to iid sequence $\{(x_t, y_t)\}_{t=1}^n$,
Denote $\mu_n(\eta) = \frac{1}{n} \sum_{t=1}^n \mathbb{E}_{h \sim w_t} \epsilon(h_t)$, let
 $\xi_t = \ln \mathbb{E}_{h \sim w_t} e^{-\eta l(h(x_t) \neq y_t)}$, set $\rho = 1$:

$$\begin{aligned} & -\ln(\mathbb{E}_{h \sim w_1} e^{-\eta n \hat{\epsilon}(h)}) + \ln(1/\delta) \\ \geq & -\ln(1 - (1 - e^{-\eta}) \mathbb{E}_{h \sim w_1} \epsilon(h)) - \dots - \ln(1 - (1 - e^{-\eta}) \mathbb{E}_{h \sim w_n} \epsilon(h)) \\ \geq & -n \ln(1 - (1 - e^{-\eta}) \mu_n(\eta)) \end{aligned}$$

Applications to Exponential Weight Algorithm II

- By Compression Lemma,

$$\Pr(\exists \pi, -\ln(1 - (1 - e^{-\eta})\mu_n(\eta)) \geq \eta \int \hat{\epsilon}(h) d\nu + \frac{D(\nu||\mu)}{n} + t) \leq e^{-nt}$$

- Fix $\alpha \in [0, 1]$, $t \geq 0$, $d \geq 0$, $\eta = \eta(\alpha, t + d)$, where $\eta(\alpha, u) = \ln(D^{-1}(\alpha, u)(1 - \alpha)) - \ln(\alpha(1 - D^{-1}(\alpha, u)))$.

$$\Pr(\exists \pi, \mu_n(\eta) \geq D^{-1}(\alpha, t + d), \int \hat{\epsilon}(h) d\nu \leq \alpha, D(\nu||\mu) \leq nd) \leq e^{-nt}$$

- Taking union bound over all $(\alpha, \delta, t) \in \{(\frac{p}{n}, \frac{q}{n}, t_0 + 2\ln(p+2) + 2\ln(q+2)), p, q = 0, \dots, n\}$,
Let

$$d_n^\pi = \frac{D(\nu||\mu) + 2\ln(\lceil n \int \hat{\epsilon}(h) d\nu \rceil + 2) + 2\ln(D(\nu||\mu) + 2) + \ln(1/\delta)}{n}$$

w.p. $1 - \delta$:





$$\forall \pi, \mu_n(\eta(\frac{\lceil n \int \hat{\epsilon}(h) d\nu \rceil}{n}, d_n^\pi)) \leq D^{-1}(\frac{\lceil n \int \hat{\epsilon}(h) d\nu \rceil}{n}, d_n^\pi)$$

Applications to Exponential Weight Algorithm III

- In particular, taking $\mu = \text{uniform over } \mathcal{H}$, $\nu = \delta_{\hat{h}}$, with η tuned optimally after observing the data,

$$\mu_n(\eta) \leq \hat{\epsilon}(\hat{h}) + 2 \frac{\ln(\frac{|\mathcal{H}|n}{\delta})}{n} + \sqrt{2\hat{\epsilon}(\hat{h}) \frac{\ln(\frac{|\mathcal{H}|n}{\delta})}{n}}$$

References I

-  Nicolò Cesa-Bianchi, Alex Conconi, and Claudio Gentile.
On the generalization ability of on-line learning algorithms.
IEEE Transactions on Information Theory,
50(9):2050–2057, 2004.
-  Nicolò Cesa-Bianchi and Claudio Gentile.
Improved risk tail bounds for on-line algorithms.
In *NIPS*, 2005.
-  R. El-Yaniv and Y. Wiener.
Agnostic selective classification.
In *NIPS*, 2011.
-  Yoav Freund, Yishay Mansour, and Robert E. Schapire.
Generalization bounds for averaged classifiers.
The Annals of Statistics, 32(4):1698–1722, 08 2004.

References II



Tong Zhang.

Data dependent concentration bounds for sequential prediction algorithms.

In *COLT*, pages 173–187, 2005.