

# PAC-Bayesian Bounds based on the Rényi Divergence

## CLASSIFICATION SETTING AND PAC-BAYESIAN BASICS

**Training set:** We draw  $m$  examples *i.i.d.* from a distribution  $D$  on  $\mathcal{X} \times \mathcal{Y}$ :

$$S = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\} \sim D^m.$$

**PAC-Bayesian learning:** Given a set  $\mathcal{H}$  of voters  $\mathcal{X} \rightarrow \{-1, 1\}$  and a training set  $S$ , we consider

1. The *prior* distribution  $P$  on  $\mathcal{H}$  encodes prior knowledge.

**Learner:** Learn a classifier  $h : \mathcal{X} \rightarrow \mathcal{Y}$  that has a low generalization risk on  $D$ :

$$R_D(h) \stackrel{\text{def}}{=} \mathbf{E}_{(x,y) \sim D} I[h(x) \neq y], \text{ where } I(a) = 1 \text{ if } a \text{ is true and } 0 \text{ otherwise.}$$

2. The *posterior* distribution  $Q$  on  $\mathcal{H}$  is obtained by learning from  $S$ .

PAC-Bayesian theory traditionally bounds the *Gibbs risk*  $R_D(G_Q) \stackrel{\text{def}}{=} \mathbf{E}_{h \sim Q} R_D(h)$  using its **empirical value**  $R_S(G_Q) \stackrel{\text{def}}{=} \mathbf{E}_{h \sim Q} R_S(h)$  and a the **Kullback-Leibler** divergence between  $Q$  and  $P$ . General PAC-Bayesian theorems are tools to derive various PAC-Bayesian bounds using **any convex function**  $\Delta : [0, 1] \times [0, 1] \rightarrow \mathbb{R}$ .

## CLASSICAL PAC-BAYESIAN THEORY

The *Kullback-Leibler divergence* between distributions  $Q$  and  $P$  is given by

$$\text{KL}(Q\|P) \stackrel{\text{def}}{=} \mathbf{E}_{h \sim Q} \ln \frac{Q(h)}{P(h)}.$$

**Lemma 3** (*Kullback-Leibler change of measure*) For any set  $\mathcal{H}$ , for any distributions  $P$  and  $Q$  on  $\mathcal{H}$ , and for any measurable function  $\phi : \mathcal{H} \rightarrow \mathbb{R}$ , we have

$$\mathbf{E}_{h \sim Q} \phi(h) \leq \text{KL}(Q\|P) + \ln \left( \mathbf{E}_{h \sim P} e^{\phi(h)} \right).$$

**Theorem 4** For any distribution  $D$ , for any set  $\mathcal{H}$  of voters  $\mathcal{X} \rightarrow \{-1, 1\}$ , for any prior  $P$  on  $\mathcal{H}$ , for any  $\delta \in (0, 1]$ , for any  $m' > 0$ , and for any convex function  $\Delta$ , with probability at least  $1 - \delta$  over the choice of  $S \sim D^m$ , we have

$$\forall Q \text{ on } \mathcal{H}: \Delta(R_S(G_Q), R_D(G_Q)) \leq \frac{1}{m'} \left[ \text{KL}(Q\|P) + \ln \frac{\mathcal{I}_{\Delta}^K(m, m')}{\delta} \right],$$

where  $\mathcal{I}_{\Delta}^K(m, m') \stackrel{\text{def}}{=} \sup_{r \in [0, 1]} \left[ \sum_{k=0}^m \text{Bin}_k^m(r) e^{m' \Delta(\frac{k}{m}, r)} \right]$ ,  $\text{Bin}_k^m(r) \stackrel{\text{def}}{=} \binom{m}{k} r^k (1-r)^{m-k}$ .

Two commons  $\Delta$  are  $\Delta_{\text{KL}}(q, p) \stackrel{\text{def}}{=} q \ln \frac{q}{p} + (1-q) \ln \frac{1-q}{1-p} \geq \Delta_{V^2}(q, p) \stackrel{\text{def}}{=} 2(q-p)^2$ :

**Corollary 6** (*Seeger, 2002; McAllester, 2003*) With probability at least  $1 - \delta$ ,

$$\forall Q \text{ on } \mathcal{H}: \quad \text{a) } \Delta_{\text{KL}}(R_S(G_Q), R_D(G_Q)) \leq \frac{1}{m} \left[ \text{KL}(Q\|P) + \ln \frac{2\sqrt{m}}{\delta} \right],$$

$$\quad \text{b) } R_D(G_Q) \leq R_S(G_Q) + \sqrt{\frac{1}{2m} \left[ \text{KL}(Q\|P) + \ln \frac{2\sqrt{m}}{\delta} \right]}.$$

## RÉNYI PAC-BAYESIAN THEORY

For any  $\alpha > 1$ , the *Rényi divergence* between distributions  $Q$  and  $P$  is given by

$$D_{\alpha}(Q\|P) \stackrel{\text{def}}{=} \frac{1}{\alpha - 1} \ln \left[ \mathbf{E}_{h \sim P} \left( \frac{Q(h)}{P(h)} \right)^{\alpha} \right]. \quad D_{\alpha}(Q\|P) = \text{KL}(Q\|P) \text{ when } \alpha \rightarrow 1.$$

**Theorem 8** (*Rényi change of measure*) For any set  $\mathcal{H}$ , for any distributions  $P$  and  $Q$  on  $\mathcal{H}$ , for any  $\alpha > 1$ , and for any measurable function  $\phi : \mathcal{H} \rightarrow \mathbb{R}$ , we have

$$\frac{\alpha}{\alpha - 1} \ln \mathbf{E}_{h \sim Q} \phi(h) \leq D_{\alpha}(Q\|P) + \ln \left( \mathbf{E}_{h \sim P} \phi(h)^{\frac{\alpha}{\alpha - 1}} \right).$$

**Theorem 9** For any distribution  $D$ , for any set  $\mathcal{H}$  of voters  $\mathcal{X} \rightarrow \{-1, 1\}$ , for any prior  $P$  on  $\mathcal{H}$ , for any  $\delta \in (0, 1]$ , for any  $\alpha > 1$ , and for any convex function  $\Delta$ , with probability at least  $1 - \delta$  over the choice of  $S \sim D^m$ , we have

$$\forall Q \text{ on } \mathcal{H}: \quad \ln \Delta(R_S(G_Q), R_D(G_Q)) \leq \frac{1}{\alpha'} \left[ D_{\alpha}(Q\|P) + \ln \frac{\mathcal{I}_{\Delta}^R(m, \alpha')}{\delta} \right],$$

where  $\alpha' = \frac{\alpha}{\alpha - 1}$ ,  $\mathcal{I}_{\Delta}^R(m, \alpha') \stackrel{\text{def}}{=} \sup_{r \in [0, 1]} \left[ \sum_{k=0}^m \text{Bin}_k^m(r) \Delta(\frac{k}{m}, r)^{\alpha'} \right]$ .

By choosing  $\alpha = 2$ , and  $\Delta(q, p) = p - q$ , we obtain as a special case:

**Corollary 10** (*≈ Honorio and Jaakkola, 2014*) With probability at least  $1 - \delta$ ,

$$\forall Q \text{ on } \mathcal{H}: \quad R_D(G_Q) \leq R_S(G_Q) + \sqrt{\frac{\chi^2(Q\|P) + 1}{4m\delta}},$$

where  $\chi^2(Q\|P) \stackrel{\text{def}}{=} \mathbf{E}_{h \sim P} \left[ \left( \frac{Q(h)}{P(h)} \right)^2 - 1 \right]$  is the *chi-squared divergence*.

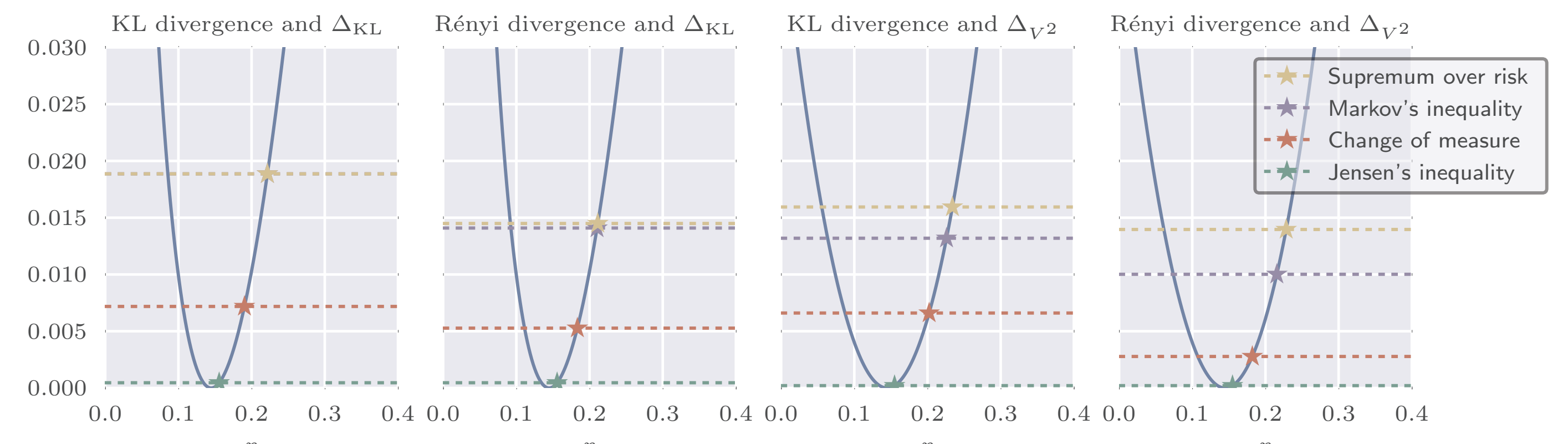
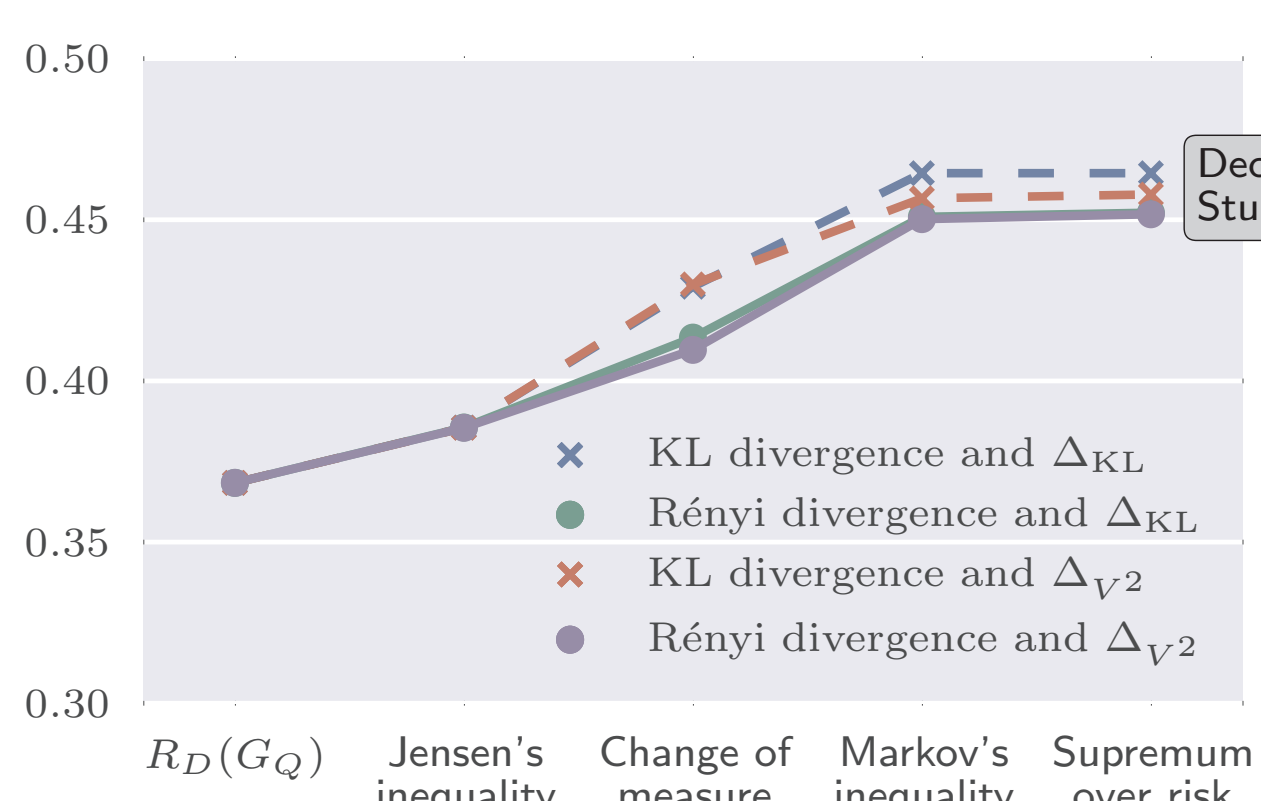
## STREAMLINED AND CUSTOMIZABLE PAC-BAYESIAN PROOFS

	KL-divergence	Rényi divergence with $\alpha' = \frac{\alpha}{\alpha - 1}$
Jensen's inequality	$\Delta \left( \mathbf{E}_{h \sim Q} R_S(h), \mathbf{E}_{h \sim Q} R_D(h) \right)$ $\leq \mathbf{E}_{h \sim Q} \Delta(R_S(h), R_D(h))$	$\ln \Delta \left( \mathbf{E}_{h \sim Q} R_S(h), \mathbf{E}_{h \sim Q} R_D(h) \right)$ $\leq \ln \mathbf{E}_{h \sim Q} \Delta(R_S(h), R_D(h))$
Change of measure	$\leq \frac{1}{m'} \left[ \text{KL}(Q\ P) + \ln \mathbf{E}_{h \sim P} e^{m' \Delta(R_S(h), R_D(h))} \right]$	$\leq \frac{1}{\alpha'} \left[ D_{\alpha}(Q\ P) + \ln \mathbf{E}_{h \sim P} \Delta(R_S(h), R_D(h))^{\alpha'} \right]$
Markov's inequality	$\leq \frac{1}{1-\delta} \frac{1}{m'} \left[ \text{KL}(Q\ P) + \ln \frac{1}{\delta} \mathbf{E}_{S' \sim D^m} \mathbf{E}_{h \sim P} e^{m' \Delta(R_{S'}(h), R_D(h))} \right]$	$\leq \frac{1}{1-\delta} \frac{1}{\alpha'} \left[ D_{\alpha}(Q\ P) + \ln \frac{1}{\delta} \mathbf{E}_{S' \sim D^m} \mathbf{E}_{h \sim P} \Delta(R_{S'}(h), R_D(h))^{\alpha'} \right]$
Expectations swap	$= \frac{1}{m'} \left[ \text{KL}(Q\ P) + \ln \frac{1}{\delta} \mathbf{E}_{h \sim P} \mathbf{E}_{S' \sim D^m} e^{m' \Delta(R_{S'}(h), R_D(h))} \right]$	$= \frac{1}{\alpha'} \left[ D_{\alpha}(Q\ P) + \ln \frac{1}{\delta} \mathbf{E}_{h \sim P} \mathbf{E}_{S' \sim D^m} \Delta(R_{S'}(h), R_D(h))^{\alpha'} \right]$
Binomial law	$= \frac{1}{m'} \left[ \text{KL}(Q\ P) + \ln \frac{1}{\delta} \mathbf{E}_{h \sim P} \sum_{k=0}^m \text{Bin}_k^m(R_D(h)) e^{m' \Delta(\frac{k}{m}, R_D(h))} \right]$	$= \frac{1}{\alpha'} \left[ D_{\alpha}(Q\ P) + \ln \frac{1}{\delta} \mathbf{E}_{h \sim P} \sum_{k=0}^m \text{Bin}_k^m(R_D(h)) \Delta(\frac{k}{m}, R_D(h))^{\alpha'} \right]$
Supremum over risk	$\leq \frac{1}{m'} \left[ \text{KL}(Q\ P) + \ln \frac{1}{\delta} \sup_{r \in [0, 1]} \left\{ \sum_{k=0}^m \text{Bin}_k^m(r) e^{m' \Delta(\frac{k}{m}, r)} \right\} \right]$	$\leq \frac{1}{\alpha'} \left[ D_{\alpha}(Q\ P) + \ln \frac{1}{\delta} \sup_{r \in [0, 1]} \left\{ \sum_{k=0}^m \text{Bin}_k^m(r) \Delta(\frac{k}{m}, r)^{\alpha'} \right\} \right]$

## EMPIRICAL STUDY

(HOW THE BOUND VALUES ARE IMPACTED BY EACH INEQUALITY OF OUR PROOF?)

Each example generated by  $D$  is a random draw among the 8124 examples of the *mushroom* dataset. That is, the training set  $S \sim D^m$  contains  $m$  examples drawn *with replacement* and *uniform probability* from the full dataset. From training set  $S$ , we learn a majority vote using AdaBoost. We compare three different kinds of voters.



Values for each inequality computed with the three kinds of voters. The dashed lines correspond to the traditional bounds with the Kullback-Leibler divergence. The full lines correspond to the bounds considering the Rényi divergence.

Alternate representation of the quantities obtained using the weak decision trees. The blue curve corresponds to the function  $\Delta(R_D(G_Q), r)$ . Each dashed horizontal line corresponds to the value given by the right-hand side of the bound after each inequality.