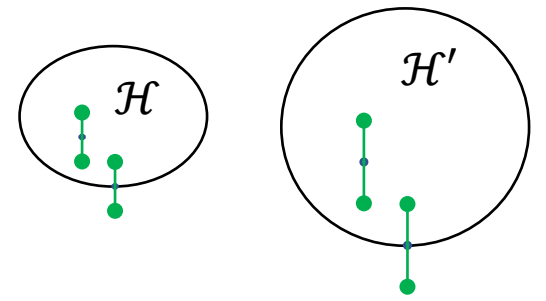


PAC-Bayesian Analysis

Yevgeny Seldin

PAC-Bayesian Analysis



- Selection \rightarrow bias
- PAC-Bayesian analysis
 - Randomized classifiers \rightarrow active avoidance of selection \rightarrow reduced bias
 - The idea: instead of committing to a particular classifier, return a distribution over classifiers (avoid commitment)
 - For example: if two classifiers have the same empirical error, do not select among them, but return a 50/50 distribution
 - Stays at the same level of approximation error, but reduces the estimation error
 - Can be applied to uncountably infinite \mathcal{H}

Randomized Classifiers

- Let ρ be a distribution on \mathcal{H}
- Randomized classification:
 1. Sample $h \sim \rho(h)$
 2. Observe X
 3. Return $h(X)$
- ρ is a *randomized classifier / Gibbs classifier*
- Expected error: $\mathbb{E}_{h \sim \rho(h)}[L(h)]$
- Empirical error: $\mathbb{E}_{h \sim \rho(h)}[\hat{L}(h, S)]$

PAC-Bayes-kl inequality

- Theorem: For any “prior” distribution π on \mathcal{H} that is independent of S

$$\mathbb{P} \left(\exists \rho: \text{kl}(\mathbb{E}_{\rho}[\hat{L}(h, S)] || \mathbb{E}_{\rho}[L(h)]) \geq \frac{\text{KL}(\rho || \pi) + \ln \frac{2\sqrt{n}}{\delta}}{n} \right) \leq \delta$$

- In other words, with probability at least $1 - \delta$, for all ρ simultaneously

$$\text{kl}(\mathbb{E}_{\rho}[\hat{L}(h, S)] || \mathbb{E}_{\rho}[L(h)]) \leq \frac{\text{KL}(\rho || \pi) + \ln \frac{2\sqrt{n}}{\delta}}{n}$$

$$\text{kl}(\mathbb{E}_\rho[\hat{L}(h, S)] || \mathbb{E}_\rho[L(h)]) \leq \frac{\text{KL}(\rho || \pi) + \ln \frac{2\sqrt{n}}{\delta}}{n}$$

Understanding the bound

- Refined Pinsker's relaxation of kl:

$$\mathbb{E}_\rho[L(h)] \leq \mathbb{E}_\rho[\hat{L}(h, S)] + \sqrt{\frac{2\mathbb{E}_\rho[\hat{L}(h, S)] \left(\text{KL}(\rho || \pi) + \ln \frac{2\sqrt{n}}{\delta} \right)}{n}} + \frac{2 \left(\text{KL}(\rho || \pi) + \ln \frac{2\sqrt{n}}{\delta} \right)}{n}$$

- Pick ρ that optimizes the trade-off between $\mathbb{E}_\rho[\hat{L}(h, S)]$ and $\text{KL}(\rho || \pi)$
 - $\mathbb{E}_\rho[\hat{L}(h, S)]$ - assign high weight on h with small $\hat{L}(h, S)$
 - $\text{KL}(\rho || \pi)$ - stay close to π in the KL sense
 - $\text{KL}(\rho || \pi) = \sum_h \rho(h) \ln \frac{\rho(h)}{\pi(h)} = \sum_h \rho(h) \ln \frac{1}{\pi(h)} + \sum_h \rho(h) \ln \rho(h) = \sum_h \rho(h) \ln \frac{1}{\pi(h)} - H(\rho)$
 - Distribute $\rho(h)$ uniformly among h with similar $\hat{L}(h, S)$ and $\pi(h)$ (maximize $H(\rho)$)
 - Extreme case: if $\rho = \pi$, then $\text{KL}(\rho || \pi) = 0$. No selection, no penalty!
 - Fast rates: small $\hat{L}(h, S)$ allows more aggressive deviation from π

PAC-Bayes vs. VC vs. VC+Occam

- PAC-Bayes (Pinsker's relaxation)

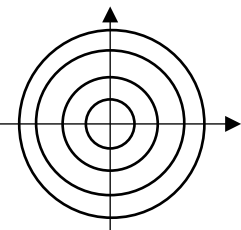
$$\mathbb{E}_\rho[L(h)] \leq \mathbb{E}_\rho[\hat{L}(h, S)] + \sqrt{\frac{\text{KL}(\rho||\pi) + \ln \frac{2\sqrt{n}}{\delta}}{2n}}$$

- VC (Empirical Risk Minimization – ERM)

$$L(h) \leq \hat{L}(h, S) + \sqrt{\frac{8 \ln \frac{2((2n)^{d_{VC}(\mathcal{H})} + 1)}{\delta}}{n}}$$

- VC+Occam (Structural Risk Minimization – SRM)

$$L(h) \leq \hat{L}(h, S) + \sqrt{\frac{8 \ln \frac{2((2n)^{d_{VC}(\mathcal{H}_{j(h)})} + 1)}{\pi(j)\delta}}{n}}$$



Special case: SVMs

$$L_{\text{FAT}}(h_{w,b}) \leq \hat{L}_{\text{FAT}}(h_{w,b}, S) + \sqrt{\frac{8 \ln \frac{2((2n)^{\lceil ||w||^2 \rceil + 1} + 1) \lceil ||w||^2 \rceil (\lceil ||w||^2 \rceil + 1)}{\delta}}{n}}$$

- PAC-Bayes

- Complexity $\pi(h)$ defined for each h individually
- Penalization by the actual selection $\text{KL}(\rho||\pi)$
 - No selection, no penalty
 - $\rho = \pi \Rightarrow \text{KL}(\rho||\pi) = 0$
- Possibility to incorporate prior knowledge via $\pi(h)$

- VC

- Complexity $d_{VC}(\mathcal{H})$ defined for all \mathcal{H}
- “Pre-paid” penalization $d_{VC}(\mathcal{H})$
 - Same cost irrespective of selection
- No possibility to incorporate prior knowledge

- VC+Occam: in-between

- Occam provides partial adaptivity across subsets of \mathcal{H}
- $d_{VC}(\mathcal{H}_j)$ dominates the complexity
- $\pi(j)$ allows to incorporate structural prior knowledge, but only limited data-dependent prior knowledge

$$\text{kl}(\mathbb{E}_\rho[\hat{L}(h, S)] || \mathbb{E}_\rho[L(h)]) \leq \frac{\text{KL}(\rho || \pi) + \ln \frac{2\sqrt{n}}{\delta}}{n}$$

PAC-Bayes vs. Bayesian learning

- PAC-Bayesian bounds

- $\pi(h)$ is an auxiliary construction in the proof (same as in Occam); the bounds always hold
- High-probability guarantees on the distance between $\mathbb{E}_\rho[\hat{L}(h, S)]$ and $\mathbb{E}_\rho[L(h)]$
- Depend on the loss function $\ell(h(X), Y)$

- Bayesian learning

- $\pi(h)$ is a prior “belief”. The Bayes rule provides a way to update it to posterior “belief” given evidence (data)
- No guarantees
- Built for the log-loss

A proof of the PAC-Bayes-kl inequality

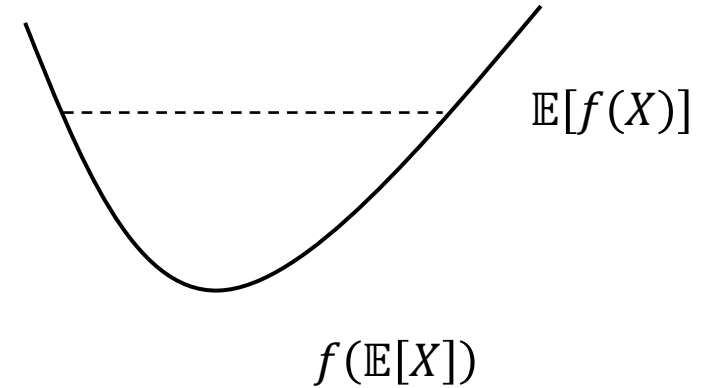
- Theorem: For any “prior” distribution π on \mathcal{H} that is independent of S

$$\mathbb{P} \left(\exists \rho: \text{kl}(\mathbb{E}_{\rho}[\hat{L}(h, S)] || \mathbb{E}_{\rho}[L(h)]) \geq \frac{\text{KL}(\rho || \pi) + \ln \frac{2\sqrt{n}}{\delta}}{n} \right) \leq \delta$$

Proof tools – Jensen's inequality

- Jensen's inequality:

For convex f : $\mathbb{E}[f(X)] \geq f(\mathbb{E}[X])$



- Example:

- $\mathbb{E}[X^2] \geq (\mathbb{E}[X])^2$

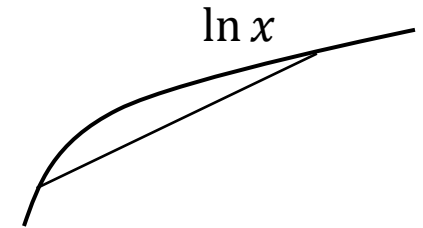
Main proof tool – Change of measure inequality

- Lemma (Change of measure inequality):

For any f, ρ , and π : $\mathbb{E}_\rho[f(h)] \leq KL(\rho||\pi) + \ln \mathbb{E}_\pi[e^{f(h)}]$

- Proof:

$$\begin{aligned}\mathbb{E}_\rho[f(h)] &= \mathbb{E}_\rho \left[\ln \left(\frac{\rho(h)}{\pi(h)} e^{f(h)} \frac{\pi(h)}{\rho(h)} \right) \right] \\ &= KL(\rho||\pi) + \mathbb{E}_\rho \left[\ln \left(e^{f(h)} \frac{\pi(h)}{\rho(h)} \right) \right] \\ &\stackrel{\text{Jensen}}{\leq} KL(\rho||\pi) + \ln \mathbb{E}_\rho \left[e^{f(h)} \frac{\pi(h)}{\rho(h)} \right] \\ &= KL(\rho||\pi) + \ln \mathbb{E}_\pi[e^{f(h)}]\end{aligned}$$



- We obtain a deterministic relation between $\mathbb{E}_\rho[f(h)]$ and $\mathbb{E}_\pi[e^{f(h)}]$, no probabilities involved.

$$\mathbb{P} \left(\exists \rho: \text{kl}(\mathbb{E}_\rho[\hat{L}(h, S)] || \mathbb{E}_\rho[L(h)]) \geq \frac{\text{KL}(\rho || \pi) + \ln \frac{2\sqrt{n}}{\delta}}{n} \right) \leq \delta$$

Proof of PAC-Bayes-kl

Markov:

$$\mathbb{P}(X \geq \varepsilon) \leq \frac{\mathbb{E}[X]}{\varepsilon} = \delta$$

$$\mathbb{P} \left(X \leq \frac{1}{\delta} \mathbb{E}[X] \right) \geq 1 - \delta$$

- Change of measure inequality

$$\mathbb{E}_\rho[f(h)] \leq \text{KL}(\rho || \pi) + \ln \mathbb{E}_\pi[e^{f(h)}]$$

Deterministic inequality relating all ρ to a single π

- PAC-Bayes Lemma

$$\mathbb{E}_\rho[f(h, S)] \leq \text{KL}(\rho || \pi) + \ln \overbrace{\mathbb{E}_\pi[e^{f(h, S)}]}$$

(deterministically $\forall \rho$)

$$\stackrel{\substack{\text{Markov} \\ \text{w.p.} \geq 1-\delta}}{\leq} \text{KL}(\rho || \pi) + \ln \frac{1}{\delta} + \ln \mathbb{E}_S \left[\overbrace{\mathbb{E}_\pi[e^{f(h, S)}]} \right]$$

(single application of Markov to $\mathbb{E}_\pi[e^{f(h, S)}]$)

$$\stackrel{\substack{\text{w.p.} \geq 1-\delta \\ \pi \text{ independent of } S}}{=} \text{KL}(\rho || \pi) + \ln \frac{1}{\delta} + \ln \mathbb{E}_\pi \left[\mathbb{E}_S[e^{f(h, S)}] \right]$$

- Choice of $f(h, S) = n \text{kl}(\hat{L}(h, S) || L(h))$

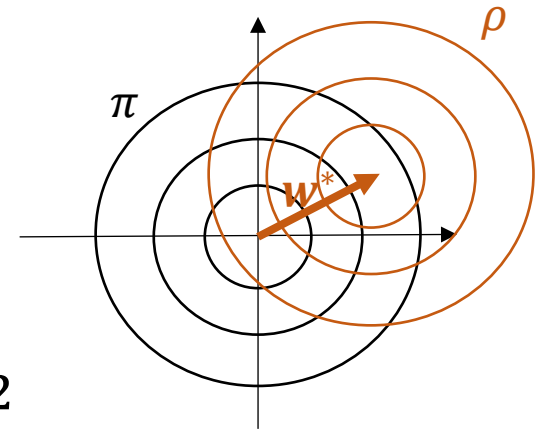
- kl-Lemma: $\mathbb{E}_S[e^{f(h, S)}] \leq 2\sqrt{n}$

- $n \text{kl}(\mathbb{E}_\rho[\hat{L}(h, S)] || \mathbb{E}_\rho[L(h)]) \stackrel{\substack{\text{convexity} \\ \text{of kl}}}{\leq} \mathbb{E}_\rho \left[n \text{kl}(\hat{L}(h, S) || L(h)) \right] \stackrel{\substack{\text{w.p.} \geq 1-\delta \\ \text{PAC-Bayes lemma}}}{\leq} \text{KL}(\rho || \pi) + \ln \frac{2\sqrt{n}}{\delta}$

Working with the bound

$$\text{kl}(\mathbb{E}_\rho[\hat{L}(h, S)] || \mathbb{E}_\rho[L(h)]) \leq \frac{\text{KL}(\rho || \pi) + \ln \frac{2\sqrt{n}}{\delta}}{n}$$

- Select π . Example: $\pi(h_w) = \mathcal{N}(0, I)$
- Select ρ . Example: $\rho(h_w) = \mathcal{N}(w^*, I)$
- Calculate $\text{KL}(\rho || \pi)$. In the example: $\text{KL}(\rho || \pi) = ||w^*||^2$
- Calculate $\mathbb{E}_\rho[\hat{L}(h, S)]$
 - The challenging part



Modularity of the bound

$$\text{kl}(\mathbb{E}_\rho[\hat{L}(h, S)] || \mathbb{E}_\rho[L(h)]) \leq \frac{\text{KL}(\rho || \pi) + \ln \frac{2\sqrt{n}}{\delta}}{n}$$

- Different choices of $f(h, S)$ give divergence measures between $\mathbb{E}_\rho[\hat{L}(h, S)]$ and $\mathbb{E}_\rho[L(h)]$
 - PAC-Bayes-kl: $f(h, S) = n \text{kl}(\hat{L}(h, S) || L(h))$
 - kl-Lemma: $\mathbb{E}_S \left[e^{n \text{kl}(\hat{L}(h, S) || L(h))} \right] \leq 2\sqrt{n}$
 - PAC-Bayes-Hoeffding: $f(h, S) = n\lambda (L(h) - \hat{L}(h, S))$
 - Hoeffding's Lemma: $\mathbb{E}_S \left[e^{n\lambda(L(h) - \hat{L}(h, S))} \right] \leq e^{\frac{n\lambda^2}{8}}$
- Different choices of ρ and π give different regularizations.
 - Gaussian prior and posterior \rightarrow regularization by $\|w\|^2 = \sum_{i=1}^d w_i^2$
 - Laplacian prior and posterior \rightarrow regularization by $\|w\|_1 = \sum_{i=1}^d |w_i|$

$$\text{kl}(\mathbb{E}_\rho[\hat{L}(h, S)] || \mathbb{E}_\rho[L(h)]) \leq \frac{\text{KL}(\rho || \pi) + \ln \frac{2\sqrt{n}}{\delta}}{n}$$

Minimization of the bound

- Relaxation: PAC-Bayes- λ (based on refined Pinsker's inequality)

$$\mathbb{E}_\rho[L(h)] \leq \frac{\mathbb{E}_\rho[\hat{L}(h, S)]}{1 - \frac{\lambda}{2}} + \frac{\text{KL}(\rho || \pi) + \ln \frac{2\sqrt{n}}{\delta}}{n\lambda \left(1 - \frac{\lambda}{2}\right)}$$

- For a fixed ρ convex in λ , for a fixed λ convex in ρ

- Apply alternating minimization

- $\rho_\lambda^*(h) = \frac{\pi(h)e^{-n\lambda\hat{L}(h,S)}}{\mathbb{E}_\pi[e^{-n\lambda\hat{L}(h,S)}]}$ (Gibbs distribution)
- Holds for any \mathcal{H} , but computationally tractable only for finite \mathcal{H}
- $\lambda_\rho^* = \frac{2}{\sqrt{\frac{2n\mathbb{E}_\rho[\hat{L}(h,S)]}{\text{KL}(\rho || \pi)} + 1} + 1} \in (0, 1]$

- The bound is *not* jointly convex in ρ and λ

- Convergence to a *local* minimum, but in many practical cases still global

Construction of an interesting finite \mathcal{H}

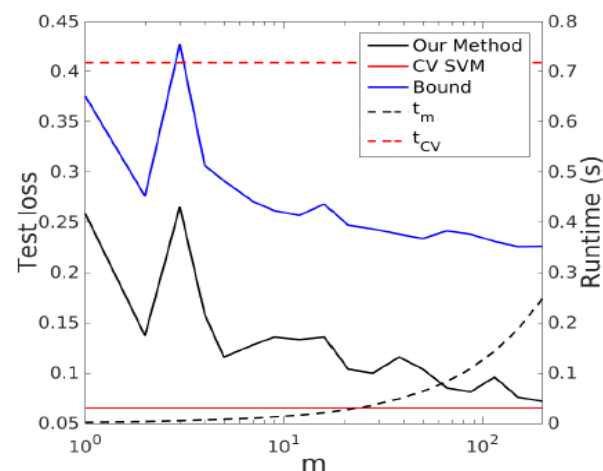


- Train m classifiers on subsamples of size r
- Validate each classifier on the corresponding remaining $n - r$ samples
- Same as in cross-validation or bagging (e.g., random forests)
- Let $\hat{L}^{\text{val}}(h, S)$ be the corresponding validation losses
- Adapted bound:

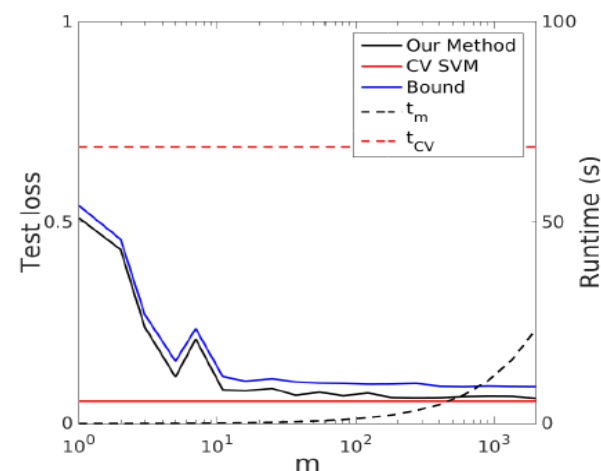
$$\mathbb{E}_{\rho}[L(h)] \leq \frac{\mathbb{E}_{\rho}[\hat{L}^{\text{val}}(h, S)]}{1 - \frac{\lambda}{2}} + \frac{\text{KL}(\rho || \pi) + \ln \frac{2\sqrt{n-r}}{\delta}}{\lambda \left(1 - \frac{\lambda}{2}\right) (n-r)}$$

- Can be used to provide generalization bounds for random forests
- In the case of kernel SVMs provides computational speed-up when r is small (“inverse cross-validation”) and $n > mr$
 - Kn^{2+} vs. $m(r^{2+} + r(n-r) + \text{BoundOptimization})$

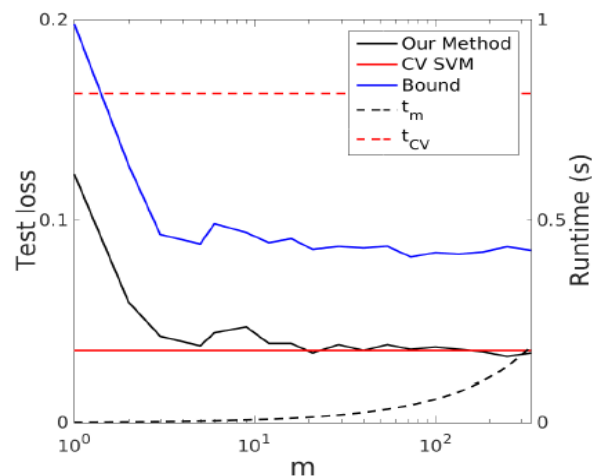
Empirical Evaluation – ensembles of small SVMs



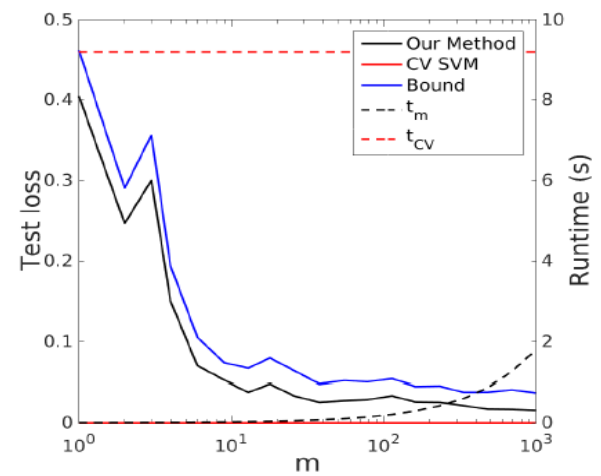
(a) Ionosphere dataset. $n = 200$,
 $r = d + 1 = 35$.



(b) Waveform dataset. $n = 2000$,
 $r = d + 1 = 41$.



(c) Breast cancer dataset. $n = 340$,
 $r = d + 1 = 11$.



(d) AvsB dataset. $n = 1000$, $r = d + 1 = 17$.

Summary

- PAC-Bayes-kl bound: with probability at least $1 - \delta$, for all ρ

$$\text{kl}(\mathbb{E}_\rho[\hat{L}(h, S)] || \mathbb{E}_\rho[L(h)]) \leq \frac{\text{KL}(\rho || \pi) + \ln \frac{2\sqrt{n}}{\delta}}{n}$$

- Refined Pinsker's relaxation for intuition

$$\mathbb{E}_\rho[L(h)] \leq \mathbb{E}_\rho[\hat{L}(h, S)] + \sqrt{\frac{2\mathbb{E}_\rho[\hat{L}(h, S)] \left(\text{KL}(\rho || \pi) + \ln \frac{2\sqrt{n}}{\delta} \right)}{n}} + \frac{2 \left(\text{KL}(\rho || \pi) + \ln \frac{2\sqrt{n}}{\delta} \right)}{n}$$

- Refined Pinsker's relaxation for alternating minimization

$$\mathbb{E}_\rho[L(h)] \leq \frac{\mathbb{E}_\rho[\hat{L}(h, S)]}{1 - \frac{\lambda}{2}} + \frac{\text{KL}(\rho || \pi) + \ln \frac{2\sqrt{n}}{\delta}}{n\lambda \left(1 - \frac{\lambda}{2} \right)}$$