


Lec 2 : KL Divergence

Yaejin Han



Defn. (KL Divergence)

For two probability distributions P, Q over the same space, the Kullback-Leibler divergence (or the relative entropy) of P w.r.t. Q is

$$D_{KL}(P \parallel Q) = \begin{cases} \mathbb{E}_{X \sim P} \left[\log \frac{dP}{dQ}(X) \right] & \text{if } P \ll Q \\ +\infty & \text{o.w.} \end{cases}$$

Remark: 1. The above defn. covers both discrete and continuous cases, i.e.

$$D_{KL}(P \parallel Q) = \sum_{x \in X} p(x) \log \frac{p(x)}{q(x)} \quad \text{if } p, q \text{ are pmfs}$$

and

$$D_{KL}(P \parallel Q) = \int p(x) \log \frac{p(x)}{q(x)} d\mu(x) \quad \text{if } p, q \text{ are pdfs w.r.t. } \mu.$$

2. This is a divergence rather than a distance, i.e. $D_{KL}(P \parallel Q) \neq D_{KL}(Q \parallel P)$.

For this reason, we write $D_{KL}(P \parallel Q)$ instead of $D_{KL}(P, Q)$.

3. IT origin: $D_{KL}(P \parallel Q)$ is the "redundancy" of using Q for source coding while the true distribution is P :

$$D_{KL}(P \parallel Q) = \underbrace{\sum_x p(x) \log \frac{1}{q(x)}}_{\text{expected codelength of using } Q} - \underbrace{H(P)}_{\text{optimal expected codelength for source } P}$$

Basic properties

Property I: $D_{KL}(P \parallel Q) \geq 0$, with equality iff $P = Q$.

Pf. $D_{KL}(P \parallel Q) = \mathbb{E}_P \left[\log \frac{dP}{dQ} \right] = \mathbb{E}_P \left[-\log \frac{dQ}{dP} \right] \geq -\log \mathbb{E}_P \left[\frac{dQ}{dP} \right] = 0. \quad \square$

Note: this gives the usual proof of

$$I(X; Y) = \mathbb{E}_{P_{XY}} \left[\log \frac{P_{XY}(X, Y)}{P_X(X)P_Y(Y)} \right] = D_{KL}(P_{XY} \parallel P_X P_Y) \geq 0.$$

Also, equality holds iff $P_{XY} = P_X P_Y$, i.e. X and Y are independent.

Property II : $(P, Q) \mapsto D_{KL}(P \parallel Q)$ is joint convex.

Pf. Follow from the joint convexity of $(x, y) \mapsto x \log \frac{x}{y}$ over \mathbb{R}_+^2 ,
whose Hessian is $\begin{bmatrix} 1/x & -1/y \\ -1/y & x/y^2 \end{bmatrix} \succeq 0$. \square

Property III (Chain rule): $D_{KL}(P_{X^n} \parallel Q_{X^n}) = \sum_{i=1}^n \mathbb{E}_{P_{X^{i-1}}} [D_{KL}(P_{X_i | X^{i-1}} \parallel Q_{X_i | X^{i-1}})]$.

Pf.

$$\begin{aligned} D_{KL}(P_{X^n} \parallel Q_{X^n}) &= \mathbb{E}_{P_{X^n}} \left[\log \frac{P_{X^n}}{Q_{X^n}} \right] \\ &= \mathbb{E}_{P_{X^n}} \left[\sum_{i=1}^n \log \frac{P_{X_i | X^{i-1}}}{Q_{X_i | X^{i-1}}} \right] \\ &= \sum_{i=1}^n \mathbb{E}_{P_{X^{i-1}}} \underbrace{\mathbb{E}_{P_{X_i | X^{i-1}}} \left[\log \frac{P_{X_i | X^{i-1}}}{Q_{X_i | X^{i-1}}} \right]}_{= D_{KL}(P_{X_i | X^{i-1}} \parallel Q_{X_i | X^{i-1}})} \end{aligned} \quad \square$$

Data processing inequality (DPI) : an important property of KL divergence

If

$$\begin{array}{ccccc} P_X & & & & P_Y \\ & \searrow & P_{Y|X} & \swarrow & \\ Q_X & & & & Q_Y \end{array}$$

then

$$D_{KL}(P_X \parallel Q_X) \geq D_{KL}(P_Y \parallel Q_Y)$$

(i.e. distributions become "closer" after processing)

Pf. (Method 1: convexity) For conditional kernel $K(x|y) = \frac{P_{Y|X}(y|x)}{Z_y}$,

$$P_Y(y) = \sum_x P_X(x) P_{Y|X}(y|x) = Z_y \cdot \sum_x P_X(x) K(x|y) = Z_y \mathbb{E}_{X \sim K(\cdot|y)} [P_X(x)]$$

(Q_Y) (Q_X) (Q_X) (Q_X)

$$\Rightarrow D_{KL}(P_Y \parallel Q_Y) = \sum_y P_Y(y) \log \frac{P_Y(y)}{Q_Y(y)}$$

$$= \sum_y Z_y \cdot \mathbb{E}_{X \sim K(\cdot|y)} [P_X(x)] \log \frac{\mathbb{E}_{X \sim K(\cdot|y)} [P_X(x)]}{\mathbb{E}_{X \sim K(\cdot|y)} [Q_X(x)]}$$

convexity of
 $(x, y) \mapsto x \log \frac{x}{y}$

$$\leq \sum_y Z_y \cdot \mathbb{E}_{X \sim K(\cdot|y)} \left[P_X(x) \log \frac{P_X(x)}{Q_X(x)} \right]$$

$$= \sum_y \sum_x P_{Y|X}(y|x) P_X(x) \log \frac{P_X(x)}{Q_X(x)} = D_{KL}(P_X \parallel Q_X).$$

(Method 2: chain rule) Let $P_{XY} = P_X P_{Y|X}$, $Q_{XY} = Q_X P_{Y|X}$.

$$\begin{aligned}
 D_{KL}(P_X \parallel Q_X) &= D_{KL}(P_X \parallel Q_X) + \underbrace{\mathbb{E}_{P_X}[D_{KL}(P_{Y|X} \parallel Q_{Y|X})]}_{=0} \\
 &= D_{KL}(P_{XY} \parallel Q_{XY}) \\
 &= D_{KL}(P_Y \parallel Q_Y) + \underbrace{\mathbb{E}_{P_Y}[D_{KL}(P_{X|Y} \parallel Q_{X|Y})]}_{\geq 0} \quad \square
 \end{aligned}$$

Applications of DPI

① DPI of mutual information: if $X - Y - Z$, then

$$I(X; Y) \geq I(X; Z)$$

Pf.

$$\begin{array}{ccc}
 P_{XY} & \xrightarrow{Id \otimes P_{Z|Y}} & P_{XZ} \text{ (by Markov)} \\
 P_X P_Y & & P_X P_Z
 \end{array}$$

$$\Rightarrow I(X; Y) = D_{KL}(P_{XY} \parallel P_X P_Y) \geq D_{KL}(P_{XZ} \parallel P_X P_Z) = I(X; Z) \quad \square$$

② Fano's inequality: if $X \sim \text{Unif}([M])$, then

$$P(X \neq Y) \geq 1 - \frac{I(X; Y) + \log 2}{\log M}.$$

Pf.

$$\begin{array}{ccc}
 P_{XY} & \xrightarrow{(x, y) \mapsto \mathbb{1}(x=y)} & \text{Bern}(P(X=Y)) \\
 P_X P_Y & & \text{Bern}(\frac{1}{M}) \text{ (since } X \sim \text{Unif}([M]) \text{)}
 \end{array}$$

$$\Rightarrow I(X; Y) = D_{KL}(P_{XY} \parallel P_X P_Y)$$

$$\geq D_{KL}(\text{Bern}(P(X=Y)) \parallel \text{Bern}(\frac{1}{M}))$$

$$= (1 - P(X \neq Y)) \log \frac{1 - P(X \neq Y)}{1/M} + P(X \neq Y) \log \frac{P(X \neq Y)}{1 - \frac{1}{M}}$$

$$\geq (1 - P(X \neq Y)) \log M - \log 2 \quad \square$$

③ Contiguity: \forall event A : $P(A) \log \frac{e^{P(A)}}{Q(A)} \leq D_{KL}(P \parallel Q)$

(so if $D_{KL}(P \parallel Q) = 0$, then $Q(A) \rightarrow 0 \Rightarrow P(A) \rightarrow 0$)

Pf.

$$\begin{array}{ccc} P & \searrow & \text{Bern}(P(A)) \\ Q & \nearrow & \text{Bern}(Q(A)) \end{array} \quad x \mapsto 1(x \in A)$$

$$\Rightarrow D_{KL}(P \parallel Q) \geq D_{KL}(\text{Bern}(P(A)) \parallel \text{Bern}(Q(A))) \geq P(A) \log \frac{e^{P(A)}}{Q(A)} \quad \square$$

Dual representation of KL: move from distributions to functions

Donsker-Varadhan. $D_{KL}(P \parallel Q) = \sup_f \mathbb{E}_P f - \log \mathbb{E}_Q[e^f]$.
where the sup is taken over all functions f with $\mathbb{E}_Q[e^f] < \infty$.

Pf (\leq) Take $f = \log \frac{dP}{dQ}$.

(\geq) By replacing f by $f - c$, WLOG can assume $\mathbb{E}_Q[e^f] = 1$.

In this case,

$$\tilde{Q}(dx) = e^{f(x)} Q(dx) \text{ is also a distribution.}$$

So

$$\begin{aligned} D_{KL}(P \parallel Q) - \mathbb{E}_P f &= \mathbb{E}_P \left[\log \frac{dP}{dQ} \right] = \mathbb{E}_P \left[\log \frac{dP}{d\tilde{Q}} \right] \\ &= D_{KL}(P \parallel \tilde{Q}) \geq 0 \end{aligned} \quad \square$$

Gibbs variational principle.

$$\log \mathbb{E}_Q[e^f] = \sup_P \mathbb{E}_P f - D_{KL}(P \parallel Q)$$

Pf. (\leq) Take $P(dx) = \frac{e^{f(x)} Q(dx)}{\mathbb{E}_Q[e^f]}$.

(\geq) By Donsker-Varadhan. □

Both results have numerous applications in practice.

Application 1: transportation inequalities

Example 1.1. Restricting Donsker-Vardhan to $f = \lambda g$ with $\|g\|_\infty \leq 1$:

$$\begin{aligned} D_{KL}(P \| Q) &\geq \sup_{\substack{\lambda \in \mathbb{R} \\ \|g\|_\infty \leq 1}} \lambda \mathbb{E}_P[g] - \underbrace{\log \mathbb{E}_Q[e^{\lambda g}]}_{\leq \lambda \mathbb{E}_Q[g] + \frac{\lambda^2}{8} \text{ by Hoeffding's ineq.}} \\ &\geq \sup_{\substack{\lambda \in \mathbb{R} \\ \|g\|_\infty \leq 1}} \lambda (\mathbb{E}_P[g] - \mathbb{E}_Q[g]) - \frac{\lambda^2}{8} \\ &= 2 \left(\sup_{\|g\|_\infty \leq 1} \mathbb{E}_P[g] - \mathbb{E}_Q[g] \right)^2 \\ &= 2 \cdot TV(P, Q)^2, \end{aligned}$$

which is Pinsker's inequality (see next lecture, also for an alternative proof)

Example 1.2 (Bobbov & Götze): The following are equivalent:

① $\mathbb{E}_Q[e^{\lambda(f - \mathbb{E}_Q f)}] \leq \exp(\frac{\lambda^2}{2} C)$ for all 1-Lip functions f and $\lambda \in \mathbb{R}$:

② $W_1(P, Q) \leq \sqrt{2C \cdot D_{KL}(P \| Q)}$ holds for all P . Lipschitz:

Wasserstein-1 distance:

$$\inf_{\pi \in \Pi(P, Q)} \mathbb{E}_{(X, Y) \sim \pi} [d(X, Y)]$$

$$= \sup_{f: 1\text{-Lip}} \mathbb{E}_P f - \mathbb{E}_Q f.$$

$$|f(x) - f(y)| \leq d(x, y)$$

for a given metric d .

Pf. (① \Rightarrow ②) $KL(P \| Q) \geq \sup_{\substack{\lambda \in \mathbb{R} \\ f: 1\text{-Lip}}} \lambda \cdot \mathbb{E}_P f - \log \mathbb{E}_Q[e^{\lambda f}]$

$$\geq \sup_{\substack{\lambda \in \mathbb{R} \\ f: 1\text{-Lip}}} \lambda (\mathbb{E}_P f - \mathbb{E}_Q f) - \frac{\lambda^2 C}{2}$$

$$= \frac{1}{2C} \left(\sup_{f: 1\text{-Lip}} \mathbb{E}_P f - \mathbb{E}_Q f \right)^2 = \frac{W_1(P, Q)^2}{2C}.$$

$$\begin{aligned} (② \Rightarrow ①) \log \mathbb{E}_Q[e^{\lambda(f - \mathbb{E}_Q f)}] &= \sup_P \mathbb{E}_P[\lambda(f - \mathbb{E}_Q f)] - D_{KL}(P \| Q) \\ &\leq \sup_P \lambda (\mathbb{E}_P f - \mathbb{E}_Q f) - \frac{(\mathbb{E}_P f - \mathbb{E}_Q f)^2}{2C} \\ &\leq \frac{\lambda^2}{2} C \quad \square \end{aligned}$$

Application 2: variational inference.

Setting: a family of distributions $p_\theta(x^*, y^*)$ where both $p_\theta(x^*)$ and $p_\theta(y^* | x^*)$ are tractable

Problem: estimate θ given only y^* (x^* not observable: missing data/latent variable)

Difficulty: $p_\theta(y^*) = \int p_\theta(x^*) p_\theta(y^* | x^*) dx^*$ often not log-concave or tractable

Evidence Lower Bound (ELBO)

$$\log p_\theta(y^*) = \sup_q \mathbb{E}_{x^* \sim q} \left[\log \frac{p_\theta(x^*, y^*)}{q(x^*)} \right]$$

Pf. Gibbs variational principle

$$\Rightarrow \log p_\theta(y^*) = \log \mathbb{E}_{p_\theta(x^*)} e^{\log p_\theta(y^* | x^*)}$$

$$= \sup_q \mathbb{E}_{q(x^*)} [\log p_\theta(y^* | x^*)] - D_{KL}(q \| p)$$

$$= \text{ELBO}$$



Example 2.1 (Ising model). $p(y^*) = \frac{1}{Z} \exp\left(\sum_{i,j} A_{ij} y_i y_j + \sum_i b_i y_i\right)$, $y^* \in \{\pm 1\}^n$

Variational inference of $\log Z$:

$$\log Z = \log \left(2^n \mathbb{E}_{y^* \sim \text{Unif}(\{\pm 1\}^n)} \exp\left(\sum_{i,j} A_{ij} y_i y_j + \sum_i b_i y_i\right) \right)$$

$$= n \log 2 + \sup_p \left(\mathbb{E}_p \left[\sum_{i,j} A_{ij} y_i y_j + \sum_i b_i y_i \right] - D_{KL}(p \| \text{Unif}(\{\pm 1\}^n)) \right)$$

$$= \sup_p \mathbb{E}_p \left[\sum_{i,j} A_{ij} y_i y_j + \sum_i b_i y_i \right] + H(p).$$

Relaxing to $p = \prod_i \text{Bern}(p_i)$ and optimizing over (p_1, \dots, p_n) yield a tractable lower bound.

Example 2.2 (EM algorithm): aim to find the MLE

$$\operatorname{argmax}_{\theta} \log p_{\theta}(y^n) = \operatorname{argmax}_{\theta} \sup_q \mathbb{E}_{x^n \sim q} \left[\log \frac{p_{\theta}(x^n, y^n)}{q(x^n)} \right].$$

Successive maximization:

- E step: fix $\theta = \theta^{(t)}$, the maximizer is $q^{(t)}(x^n) = p_{\theta^{(t)}}(x^n | y^n)$
- M step: fix $q = q^{(t)}$, the maximizer is $\theta^{(t+1)} = \operatorname{argmax}_{\theta} \mathbb{E}_{x^n \sim q^{(t)}} [\log p_{\theta}(x^n, y^n)]$

factorizable in the missing data case

$$\theta^{(t+1)} = \operatorname{argmax}_{\theta} \mathbb{E}_{x^n \sim q^{(t)}} [\log p_{\theta}(x^n, y^n)]$$

no integral; tractable

(For example, in exponential families $p_{\theta}(x, y) \propto \exp(\langle \theta, T(x, y) \rangle - A(\theta))$,

E-step corresponds to the computation of $\mu_i \triangleq \mathbb{E}_{x^n \sim p_{\theta^{(t)}}(\cdot | y^n)} [T(x_i, y_i)]$,

and M-step corresponds to the usual MLE computation $\nabla A(\theta^{(t+1)}) = \frac{1}{n} \sum_{i=1}^n \mu_i$.)

Example 2.3 (VAE): given images y_1, \dots, y_n , aim to find a generative model:

$$x_i \sim N(0, I), \quad y_i \sim N(\mu_{\theta}(x_i), \sigma_{\theta}^2(x_i)I)$$

parametrized by neural nets

Using ELBO:
$$\max_{\theta} \log p_{\theta}(y^n) \geq \max_{\theta} \max_{q} \mathbb{E}_{x^n \sim q} \left[\log \frac{p_{\theta}(x^n, y^n)}{q(x^n)} \right]$$

another Gaussian
 $q_{\theta}(y_i) = N(\mu_{\theta}(y_i), \sigma_{\theta}^2(y_i)I)$

Idea of VAE: ① Replace $\mathbb{E}_{x^n \sim q}$ by empirical mean of simulated samples

$$x_{ij} \sim N(\mu_{\theta}(y_j), \sigma_{\theta}^2(y_j)I), \quad j=1, 2, \dots, M;$$

② Compute ∇_{θ} by the explicit expression of $\log p_{\theta}(x^n, y^n)$;

③ Compute ∇_{ϕ} by the reparametrization trick:

$$\nabla_{\phi} \mathbb{E}_{x \sim N(\mu_{\phi}, \sigma_{\phi}^2 I)} [f(x)] = \nabla_{\phi} \mathbb{E}_{\varepsilon \sim N(0, I)} [f(\mu_{\phi} + \sigma_{\phi} \varepsilon)]$$

$$= \mathbb{E}_{\varepsilon \sim N(0, I)} [\nabla_{\phi} f(\mu_{\phi} + \sigma_{\phi} \varepsilon)]$$

$$\approx \frac{1}{M} \sum_{i=1}^M \nabla_{\phi} f(\mu_{\phi} + \sigma_{\phi} \varepsilon_i).$$

Application 3: adaptive data analysis

Problem: data $X^n \stackrel{\text{i.i.d.}}{\sim} P$, a class of functions $\{\phi_t: X \rightarrow \mathbb{R}\}$

For each given ϕ_t , we have

$$P_n \phi_t := \frac{1}{n} \sum_{i=1}^n \phi_t(X_i) \approx \mathbb{E}_P[\phi_t(X_i)] =: P \phi_t$$

What happens to $P_n \phi_T$ if the index T depends on the data X^n ?

Example 3.1 (Russo & Zou '16) If each ϕ_t is σ^2 -sub-Gaussian under P ,

then
$$\left| \mathbb{E}[P_n \phi_T] - \mathbb{E}[P \phi_T] \right| \leq \sqrt{\frac{2\sigma^2}{n} I(X^n; T)}.$$

Remark: ① If $I(T; X^n) = 0$, i.e. T is independent of X^n ,

then $P_n \phi_T$ is unbiased for $P \phi_T$.

② If $T \in \{1, \dots, m\}$, then $I(X^n; T) \leq H(T) \leq \log m$,

and the upper bound $\sqrt{\frac{2\sigma^2 \log m}{n}}$ can be shown via union bound

Pf. Define two distributions: $P_{X^n, T}$: the joint distribution in the problem

$Q_{X^n, T} = P_{X^n} P_T$: an auxiliary distribution where

X^n and T are independent

Then
$$\mathbb{E}[P_n \phi_T] = \mathbb{E}_{P_{X^n, T}} \left[\frac{1}{n} \sum_i \phi_T(X_i) \right]$$

$$\mathbb{E}[P \phi_T] = \mathbb{E}_{Q_{X^n, T}} \left[\frac{1}{n} \sum_i \phi_T(X_i) \right]$$

Donsker-Varadhan $\Rightarrow I(X^n; T) = D_{KL}(P_{X^n, T} \| Q_{X^n, T})$

$$\begin{aligned} &\geq \sup_{\lambda \in \mathbb{R}} \mathbb{E}_{P_{X^n, T}} \left[\frac{\lambda}{n} \sum_i \phi_T(X_i) \right] \\ &\quad - \underbrace{\log \mathbb{E}_{Q_{X^n, T}} \left[e^{\frac{\lambda}{n} \sum_i \phi_T(X_i)} \right]}_{\leq \lambda \cdot \mathbb{E}[P \phi_T] + \frac{\lambda^2 \sigma^2}{2n}} \\ &\quad \text{by sub-Gaussian assumption} \end{aligned}$$

$$\geq \sup_{\lambda \in \mathbb{R}} \lambda (\mathbb{E}[P_n \phi_T] - \mathbb{E}[P \phi_T]) - \frac{\lambda^2 \sigma^2}{2n}$$

$$= (\mathbb{E}[P_n \phi_T] - \mathbb{E}[P \phi_T])^2 \cdot \frac{n}{2\sigma^2} \quad \square$$

Application 4: PAC-Bayes.

PAC-Bayes inequality. Let $X \sim P$, and consider a class of functions $\{f_\theta: X \rightarrow \mathbb{R}\}$.

Fix any prior distribution π of θ . Then w.p. $\geq 1 - \delta$ (over the randomness in X),
for all distributions p over θ ,

$$\mathbb{E}_{\theta \sim p} [f_\theta(X) - \psi(\theta)] \leq D_{KL}(p \parallel \pi) + \log \frac{1}{\delta},$$

where $\psi(\theta) := \log \mathbb{E}_{X \sim P} e^{f_\theta(X)}$.

Remark: ① The exception set depends on π , but not on p .

② This inequality holds for all p , which generalizes the union bound where p is usually taken to be a point mass $p = \delta_{\theta_0}$.

③ By taking $p = P_{\theta|X}$ to be a data-dependent distribution, we'll have

$$\begin{aligned} \mathbb{E}_{(\theta, X) \sim P_{\theta \times X}} [f_\theta(X) - \psi(\theta)] &\leq \inf_{\pi} \mathbb{E}_{P_X} [D_{KL}(P_{\theta|X} \parallel \pi)] \\ &= I(\theta; X) \quad (\text{exercise!}) \end{aligned}$$

Pf. By Markov's inequality, suffices to prove

$$\mathbb{E}_{X \sim P} \left[\sup_p \exp(\mathbb{E}_{\theta \sim p} [f_\theta(X) - \psi(\theta)] - D_{KL}(p \parallel \pi)) \right] \leq 1.$$

By Gibbs variational principle, the LHS is

$$\begin{aligned} &\mathbb{E}_{X \sim P} \left[\exp(\log \mathbb{E}_{\theta \sim \pi} e^{f_\theta(X) - \psi(\theta)}) \right] \\ &= \mathbb{E}_{X \sim P} \mathbb{E}_{\theta \sim \pi} [e^{f_\theta(X) - \psi(\theta)}] \\ &= \mathbb{E}_{\theta \sim \pi} \left[\underbrace{\mathbb{E}_{X \sim P} e^{f_\theta(X) - \psi(\theta)}}_{=1} \right] = 1 \end{aligned}$$

□

Why call it PAC-Bayes? Come from the following application in statistical learning theory:

Example 4.1. Let $f: X \rightarrow [0, 1]$, $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} P$

$$P_n f := \frac{1}{n} \sum_{i=1}^n f(X_i), \quad P f := \mathbb{E}_{X \sim P} [f(X)].$$

For fixed f , sub-Gaussian concentration (Hoeffding's inequality) gives

$$(P_n f - P f)^2 \leq \frac{1}{2n} \log \frac{2}{\delta} \quad \text{w.p. } \geq 1 - \delta.$$

By PAC-Bayes, fix a prior π , then w.p. $\geq 1 - \delta$, for any p :

$$\begin{aligned} \mathbb{E}_{f \sim p} \left[\lambda (P_n f - P f)^2 - \log \mathbb{E}_X e^{\underbrace{\lambda (P_n f - P f)^2}_{\frac{1}{8n} \text{-subGaussian}}} \right] &\leq D_{KL}(p \| \pi) + \log \frac{1}{\delta} \\ &\leq \frac{1}{2} \log \frac{1}{1 - \frac{\lambda}{4n}} \quad \text{if } \lambda < 4n. \end{aligned}$$

Choosing $\lambda = 2n$ gives
$$\mathbb{E}_{f \sim p} [(P_n f - P f)^2] \leq \frac{D_{KL}(p \| \pi) + \log \frac{2}{\delta}}{2n} \quad \forall p.$$

PAC-Bayes also has surprising applications to concentration inequalities, by choosing p and π appropriately.

Example 4.2. If $X \sim N(0, \Sigma)$, then w.p. $\geq 1 - \delta$,

$$\|X\|_2 \leq \sqrt{\text{Tr } \Sigma} + \sqrt{2\|\Sigma\|_{\text{op}} \log \frac{1}{\delta}}.$$

Remark: Try union bound to $\|X\|_2 = \sup_{\|v\|_2 \leq 1} \langle X, v \rangle$ yourself!

It's very difficult to make covering/chaining arguments give such sharp bound, because of a general shaped Σ . One need to invoke Talagrand's generic chaining to this example, but it's very difficult to carry out.

Pf. $\|X\|_2 = \sup_{\|v\|_2=1} \langle v, X \rangle.$

To apply PAC-Bayes, we construct a prior p_v such that $\mathbb{E}_{\theta \sim p_v}[\langle \theta, X \rangle] = \langle v, X \rangle.$
 A natural choice is $p_v = N(v, \sigma^2 I).$ Then for $\pi = N(0, \sigma^2 I):$ w.p. $\geq 1-\delta.$

$$\sup_{\|v\|_2 \leq 1} \mathbb{E}_{p_v} \left[\lambda \langle v, X \rangle - \underbrace{\log \mathbb{E}_\pi e^{\lambda \langle \theta, X \rangle}}_{= \frac{\lambda^2}{2} \theta^T \Sigma \theta} \right] - \underbrace{D_{KL}(p_v \| \pi)}_{= \frac{\|v\|_2^2}{2\sigma^2} = \frac{1}{2\sigma^2}} \leq \log \frac{1}{\delta}.$$

$$\Rightarrow \sup_{\|v\|_2 \leq 1} \lambda \langle v, X \rangle - \frac{\lambda^2}{2} (v^T \Sigma v + \sigma^2 \text{Tr}(\Sigma)) - \frac{1}{2\sigma^2} \leq \log \frac{1}{\delta}.$$

$$\Rightarrow \langle v, X \rangle \leq \underbrace{\frac{\lambda}{2} (v^T \Sigma v + \sigma^2 \text{Tr}(\Sigma))}_{\leq \|\Sigma\|_{op}} + \frac{1}{\lambda} \left(\frac{1}{2\sigma^2} + \log \frac{1}{\delta} \right), \quad \forall v.$$

$$\left. \begin{array}{l} \text{Optimize over } \sigma^2: \quad \sigma^2 = \frac{1}{\lambda \cdot \sqrt{\text{Tr}(\Sigma)}} \\ \text{Optimize over } \lambda: \quad \lambda = \sqrt{\frac{2 \log(1/\delta)}{\|\Sigma\|_{op}}} \end{array} \right\} \Rightarrow \|X\|_2 \leq \sqrt{\text{Tr}(\Sigma)} + \sqrt{2\|\Sigma\|_{op} \log \frac{1}{\delta}} \quad \text{w.p. } \geq 1-\delta. \quad \square$$

Example 4.3. Let X_1, \dots, X_n be i.i.d with $\mathbb{E}[X_i] = 0, \mathbb{E}[X_i X_i^T] = \Sigma,$ and that $v^T X_i$ is $v^T \Sigma v$ -sub Gaussian for any $v \in \mathbb{R}^n.$ Let $\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n X_i X_i^T$ be the sample covariance. Then w.p. $\geq 1-\delta,$

$$\|\hat{\Sigma} - \Sigma\|_{op} \leq C \|\Sigma\|_{op} \left(\sqrt{\frac{r(\Sigma) + \log \frac{1}{\delta}}{n}} + \frac{r(\Sigma) + \log \frac{1}{\delta}}{n} \right)$$

where $r(\Sigma) = \frac{\text{Tr}(\Sigma)}{\|\Sigma\|_{op}}$ is called the effective rank.

Remark: This is the result of [Koltchinskii & Lounici '17], where the key challenge is to arrive at the tight factor $r(\Sigma).$ Our proof is taken from [Zhivotovskiy '21] via PAC-Bayes.

Pf. (Throughout the proof, C denotes a large universal constant which may change from line to line).

$$\|\hat{\Sigma} - \Sigma\|_{\text{op}} = \sup_{\|u\|_2, \|v\|_2=1} u^T (\hat{\Sigma} - \Sigma) v.$$

Consider $(\theta, \theta') \sim p_{u,v} := f_u \otimes f_v$, where f_u is the density of

$$N(u, \sigma^2 I) \text{ conditioned on } (x-u)^T \Sigma (x-u) \leq r^2.$$

Clearly $\mathbb{E}_{(\theta, \theta') \sim p_{u,v}} [\theta^T (\hat{\Sigma} - \Sigma) \theta'] = u^T (\hat{\Sigma} - \Sigma) v$, and

$$p := \mathbb{P}(Z^T \Sigma^{-1} Z \leq r^2) \geq 1 - \frac{\mathbb{E}[Z^T \Sigma^{-1} Z]}{r^2} = 1 - \frac{\sigma^2 \text{Tr}(\Sigma)}{r^2} \text{ for } Z \sim N(0, \sigma^2 I).$$

Let $\pi = N(0, \sigma^2 I) \otimes N(0, \sigma^2 I)$. One can compute

$$D_{\text{KL}}(f_u \| N(0, \sigma^2 I)) = \frac{1}{2\sigma^2} + \log\left(\frac{1}{p}\right).$$

$$\text{so that } D_{\text{KL}}(p_{u,v} \| \pi) = \frac{1}{\sigma^2} + 2\log\left(\frac{1}{p}\right).$$

Now by PAC-Bayes, w.p. $\geq 1 - \delta$,

$$\begin{aligned} \sup_{\|u\|_2, \|v\|_2=1} \mathbb{E}_{(\theta, \theta') \sim p_{u,v}} \left[\underbrace{\lambda \theta^T (\hat{\Sigma} - \Sigma) \theta' - \log \mathbb{E} e^{\lambda \theta^T (\hat{\Sigma} - \Sigma) \theta'}}_{\leq \frac{C\lambda^2}{n} (\theta^T \Sigma \theta + \theta'^T \Sigma \theta')} \right] - D_{\text{KL}}(p_{u,v} \| \pi) &\leq \log \frac{1}{\delta} \\ &\leq \frac{C\lambda^2}{n} (\theta^T \Sigma \theta + \theta'^T \Sigma \theta')^2 \text{ for } \lambda \leq \frac{n}{C(\theta^T \Sigma \theta + \theta'^T \Sigma \theta')}. \end{aligned}$$

Since $\theta^T \Sigma \theta \leq (\sqrt{u^T \Sigma u} + \sqrt{(\theta-u)^T \Sigma (\theta-u)})^2 \leq (\sqrt{\|\Sigma\|_{\text{op}}} + r)^2$, we get

$$\|\hat{\Sigma} - \Sigma\|_{\text{op}} \leq \frac{C\lambda}{n} (\sqrt{\|\Sigma\|_{\text{op}}} + r)^4 + \frac{1}{\lambda} \left(\frac{1}{\sigma^2} + 2\log\left(\frac{1}{p}\right) + \log \frac{1}{\delta} \right) \text{ if } \lambda \leq \frac{n}{C(\sqrt{\|\Sigma\|_{\text{op}}} + r)^2}.$$

Choose $r^2 = 2\|\Sigma\|_{\text{op}}$, $\sigma^2 = \frac{\|\Sigma\|_{\text{op}}}{\text{Tr}(\Sigma)} = \frac{1}{r(\Sigma)}$, then $p \geq \frac{1}{2}$, and

$$\|\hat{\Sigma} - \Sigma\|_{\text{op}} \leq C \left(\frac{\lambda}{n} \|\Sigma\|_{\text{op}}^2 + \frac{1}{\lambda} (r(\Sigma) + \log \frac{1}{\delta}) \right), \text{ if } \lambda \leq \frac{n}{C\|\Sigma\|_{\text{op}}}.$$

Finally, choosing

$$\lambda \asymp \begin{cases} \frac{n}{\|\Sigma\|_{\text{op}}} \sqrt{\frac{r(\Sigma) + \log\left(\frac{1}{\delta}\right)}{n}} & \text{if } \frac{r(\Sigma) + \log\frac{1}{\delta}}{n} \leq 1 \\ \frac{n}{\|\Sigma\|_{\text{op}}} & \text{if } \frac{r(\Sigma) + \log\frac{1}{\delta}}{n} > 1 \end{cases}$$

leads to the claimed result.

