

1. Metric Entropy

Covering number: (T, d) metric space, $K \subset T$

$$\mathcal{N}(K, d, \varepsilon) = \min\{|\mathcal{N}| : \mathcal{N} \subset K \text{ is an } \varepsilon\text{-net of } K\}$$

Metric entropy: of K at scale ε ,

$$H(K, d, \varepsilon) = \log \mathcal{N}(K, d, \varepsilon)$$

Packing number: (T, d) metric space, $K \subset T$

$$\mathcal{P}(K, d, \varepsilon) = \max\{|\mathcal{M}| : \mathcal{M} \subset K \text{ is } \varepsilon\text{-separated}\}$$

Covering vs Packing: Any $K \subset T$ and $\varepsilon > 0$,

$$\mathcal{P}(K, d, 2\varepsilon) \leq \mathcal{N}(K, d, \varepsilon) \leq \mathcal{P}(K, d, \varepsilon)$$

ε -neighborhood of A :

$$A + \varepsilon B_2^n = \{x : d(x, A) \leq \varepsilon\}$$

Covering and Packing vs Volumes:

$$\frac{\text{vol}(K)}{\text{vol}(\varepsilon B_2^n)} \leq \mathcal{N}(K, \varepsilon) \leq \mathcal{P}(K, \varepsilon) \leq \frac{\text{vol}(K + (\varepsilon/2)B_2^n)}{\text{vol}((\varepsilon/2)B_2^n)}$$

Covering numbers of Euclidean ball:

$$\left(\frac{1}{\varepsilon}\right)^n \leq \mathcal{N}(B_2^n, \varepsilon) \leq \left(1 + \frac{2}{\varepsilon}\right)^n$$

Same upper bound holds for **sphere** S^{n-1} , and for $\varepsilon \in (0, 1]$,

$$n \log(1/\varepsilon) \leq \log \mathcal{N}(B_2^n, \varepsilon) \leq n \log(3/\varepsilon)$$

Convex Hull: $T \subset \mathbb{R}^n$

$$\text{conv}(T) = \left\{ \sum_{i=1}^m \lambda_i t_i : m \in \mathbb{N}, t_i \in T, \lambda_i \geq 0, \sum_{i=1}^m \lambda_i = 1 \right\}$$

Entropy of convex hulls: $P = \text{conv}(v_1, \dots, v_n) \subset B_2^n$

$$H(P, \varepsilon) = \log \mathcal{N}(P, \varepsilon) \leq C \frac{\log N}{\varepsilon^2}$$

2. Covariance Estimation

Operator norm net reduction: $A \in \mathbb{R}^{m \times n}$, $\mathcal{N} \subset S^{n-1}$

$$\|A\|_{\text{op}} \leq \frac{1}{1-\varepsilon} \sup_{x \in \mathcal{N}} \|Ax\|_2, \quad \varepsilon \in (0, 1)$$

Symmetric quadratic-form net reduction:

$$\|A\|_{\text{op}} \leq \frac{1}{1-2\varepsilon} \sup_{x \in \mathcal{N}} |x^\top A x|, \quad \varepsilon \in (0, 1/2)$$

Isotropic sub-Gaussian data: $X \in \mathbb{R}^d$ is isotropic if

$$\mathbb{E}X = 0 \text{ and } \mathbb{E}[XX^\top] = I_d \implies \mathbb{E}\langle X, u \rangle^2 = 1, \quad \forall u \in S^{d-1}$$

Sub-Gaussian norm for a vector: $X \in \mathbb{R}^d$,

$$\|X\|_{\psi_2} = \sup_{u \in S^{d-1}} \|\langle X, u \rangle\|_{\psi_2}$$

Useful result: X_1, \dots, X_d are independent, mean-zero, sub-Gaussian r.v., and let $X = (X_1, \dots, X_d) \in \mathbb{R}^d$,

$$\max_{1 \leq i \leq d} \|X_i\|_{\psi_2} \leq \|X\|_{\psi_2} \leq C \max_{1 \leq i \leq d} \|X_i\|_{\psi_2}.$$

ψ_2 -norm vs spectral norm: $X \sim \mathcal{N}(0, \Sigma) \in \mathbb{R}^d$, where $\Sigma \succeq 0$,

$$\|X\|_{\psi_2} \leq C \sqrt{\|\Sigma\|_{\text{op}}}.$$

Sample Covariance for Isotropic Sub-Gaussian: $s \geq 0$,

$$\mathbb{P} \left\{ \left\| \widehat{\Sigma} - I_d \right\|_{\text{op}} \geq CK^2 \left(\sqrt{\frac{d+s}{N}} + \frac{d+s}{N} \right) \right\} \leq 2e^{-cs}$$

Expectation bound:

$$\mathbb{E} \left\| \widehat{\Sigma} - I_d \right\|_{\text{op}} \lesssim K^2 \left(\sqrt{\frac{d}{N}} + \frac{d}{N} \right)$$

Equivalently, when $N \gg d$, the leading behavior is

$$\left\| \widehat{\Sigma} - I_d \right\|_{\text{op}} \approx \sqrt{d/N}$$

General covariance by whitening: If $\Sigma \succ 0$, define

$$Y = \Sigma^{-1/2} X.$$

Then Y is isotropic with $\widehat{\Sigma}_Y = \Sigma^{-1/2} \widehat{\Sigma} \Sigma^{-1/2}$. This gives

$$\widehat{\Sigma} - \Sigma = \Sigma^{1/2} (\widehat{\Sigma}_Y - I_d) \Sigma^{1/2},$$

and

$$\|\widehat{\Sigma} - \Sigma\|_{\text{op}} \leq \|\Sigma\| \|\widehat{\Sigma}_Y - I_d\|_{\text{op}}$$

Hence, given $N \geq d$,

$$\mathbb{P} \left\{ \left\| \widehat{\Sigma} - \Sigma \right\|_{\text{op}} \geq CK^2 \|\Sigma\|_{\text{op}} \sqrt{\frac{d}{N}} \right\} \leq 2e^{-d}$$

Weyl: For symmetric $d \times d$ matrices A, B :

$$\max_{1 \leq i \leq d} |\lambda_i(A) - \lambda_i(B)| \leq \|A - B\|_{\text{op}}$$

Davis-Kahan: Let A, B be symmetric $d \times d$ matrices. Let P_A be the orthogonal projector onto $\text{span}\{v_1(A), \dots, v_k(A)\}$. Assume there is an eigengap at k for A and that $\|A - B\| \leq \delta/2$. Let P_B be the projector onto $\text{span}\{v_1(B), \dots, v_k(B)\}$. Then

$$\|P_A - P_B\|_{\text{op}} \leq \frac{2\|A - B\|_{\text{op}}}{\delta}$$

Projection matrices: $u, v \in S^{d-1}$ and define $P_u = uu^\top$, then there exists a sign $s \in \{-1, 1\}$ such that

$$\frac{1}{2} \|u - sv\|_2 \leq \|P_u - P_v\|_{\text{op}} \leq 2 \|u - sv\|_2.$$

Perturbation of top eigenvectors: $A, B \in \mathbb{R}^{d \times d}$ symmetric matrices with $\delta := \lambda_1(A) - \lambda_2(A) > 0$. Then there exists a sign $s \in \{-1, 1\}$ such that

$$\|v_1(A) - sv_1(B)\|_2 \leq C \frac{\|A - B\|_{\text{op}}}{\delta}.$$

3. Random Processes

Random process: Let T be a nonempty set. A *random process* indexed by T is a family of real-valued random variables

$$(X_t)_{t \in T}.$$

Canonical pseudo-metric: Let $(X_t)_{t \in T}$ be a square-integrable random process. The *Canonical pseudo-metric* is:

$$d_X(s, t) := \|X_t - X_s\|_{L_2} = (\mathbb{E}|X_t - X_s|^2)^{1/2}, \quad s, t \in T.$$

Gaussian process: A random process $(X_t)_{t \in T}$ is a *Gaussian process* if every finite subcollection has a multivariate Gaussian distribution, i.e.,

$$(X_{t_1}, \dots, X_{t_m}) \sim \text{MVN}(\cdot, \cdot)$$

Canonical Gaussian process on \mathbb{R}^n : Let $g \sim \mathcal{N}(0, I_n)$, and let $T \subset \mathbb{R}^n$. Define

$$X_t := \langle g, t \rangle, \quad t \in T.$$

Then $(X_t)_{t \in T}$ is a centered Gaussian process, and

$$X_t - X_s = \langle g, t - s \rangle \sim \mathcal{N}(0, \|t - s\|_2^2).$$

Therefore, the canonical metric is

$$d_X(s, t) = \|t - s\|_2.$$

Finite maximal inequality: Let Z_1, \dots, Z_M be centered random variables s.t. $\|Z_j\|_{\psi_2} \leq L$ for all j . Then

$$\mathbb{E} \max_{1 \leq j \leq M} Z_j \leq CL \sqrt{\log M}$$

4. Gaussian Width

Definition: For a bounded set $T \subset \mathbb{R}^d$

$$w(T) := \mathbb{E} \sup_{t \in T} \langle g, t \rangle, \quad g \sim \mathcal{N}(0, I_d).$$

Euclidean ball:

$$w(B_2^d) = \mathbb{E} \|g\|_2 \asymp \sqrt{d}.$$

ℓ_1 -ball:

$$w(B_1^d) = \mathbb{E} \|g\|_\infty \asymp \sqrt{\log d}.$$

Gaussian tail upper bound:

$$G \sim \mathcal{N}(0, 1) \implies \mathbb{P}\{|G| > t\} \geq c_0 \frac{e^{-t^2/2}}{t}, \quad t \geq 1$$

Finite sets: If $T = \{t_1, \dots, t_M\}$, then

$$w(T) \leq \left(\max_{1 \leq j \leq M} \|t_j\|_2 \right) \sqrt{2 \log M}$$

Sparse unit vectors: If $T_k := \{x \in S^{d-1} : \|x\|_0 \leq k\}$, then

$$w(T_k) \lesssim \sqrt{k \log(ed/k)}$$

5. Chaining

Dudley's discrete inequality:

$$\mathbb{E} \sup_{t \in T} X_t \leq CK \sum_{k \in \mathbb{Z}} 2^{-k} \sqrt{\log \mathcal{N}(T, d, 2^{-k})}$$

Dudley's entropy integral inequality: Let $(X_t)_{t \in T}$ be a random process on metric space (T, d) s.t. $\|X_t - X_s\|_{\psi_2} \leq Kd(t, s)$ for all $s, t \in T$. Then for all fixed $t_0 \in T$,

$$\mathbb{E} \sup_{t \in T} (X_t - X_{t_0}) \leq CK \int_0^{\text{diam}(T, d)} \sqrt{\log \mathcal{N}(T, d, \varepsilon)} d\varepsilon$$

Supremum of Increments:

$$\mathbb{E} \sup_{s, t \in T} |X_t - X_s| \leq 2CK \int_0^{\text{diam}(T, d)} \sqrt{\log \mathcal{N}(T, d, \varepsilon)} d\varepsilon$$

High-Probability Dudley Bound: Let $D_T = \text{diam}(T, d)$, then with probability at least $1 - 2e^{-cu^2}$

$$\sup_{s, t \in T} |X_t - X_s| \leq CK \left[\int_0^{D_T} \sqrt{\log \mathcal{N}(T, d, \varepsilon)} d\varepsilon + u \cdot D_T \right]$$

Sub-exponential increments: Suppose $(X_t)_{t \in T}$ satisfies

$$\|X_t - X_s\|_{\psi_1} \leq Kd(t, s) \quad \text{for all } s, t \in T.$$

Then for every $t_0 \in T$,

$$\mathbb{E} \sup_{t \in T} (X_t - X_{t_0}) \leq CK \int_0^{\text{diam}(T, d)} \log \mathcal{N}(T, d, \varepsilon) d\varepsilon$$

Entropy integral bound for Gaussian width:

$$w(T) \leq C \int_0^{\text{diam}(T, \|\cdot\|_2)} \sqrt{\log \mathcal{N}(T, \|\cdot\|_2, \varepsilon)} d\varepsilon$$

6. Empirical Processes

Definition: Given a class \mathcal{F} of real-valued functions on Ω , define

$$Z_f := \mu_n(f) - \mu(f) = \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E}f(X), \quad f \in \mathcal{F}.$$

The random family $(Z_f)_{f \in \mathcal{F}}$ is called the *empirical process* indexed by \mathcal{F} , and $\sup_{f \in \mathcal{F}} |Z_f|$ the uniform empirical error.

Uniform LLN for Lipschitz Functions: Let μ be any probability measure on $[0, 1]$, and let X_1, \dots, X_n be iid with law μ . For $L > 0$, define the anchored Lipschitz class

$$\mathcal{F}_L = \{f : [0, 1] \rightarrow \mathbb{R} : f(0) = 0, \|f\|_{\text{Lip}} \leq L\}.$$

Then

$$\mathbb{E} \sup_{f \in \mathcal{F}_L} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E}f(X) \right| \leq \frac{CL}{\sqrt{n}}.$$

Covering Lipschitz functions: For $L > 0$, define

$$\mathcal{F}_{L, d} := \left\{ f : [0, 1]^d \rightarrow [0, 1] : \|f\|_{\text{Lip}} \leq L \right\}.$$

Then

$$\log \mathcal{N}(\mathcal{F}_{L, 1}, \|\cdot\|_{\infty}, \varepsilon) \leq C \frac{L}{\varepsilon}, \quad 0 < \varepsilon \leq 1,$$

$$\log \mathcal{N}(\mathcal{F}_{L, d}, \|\cdot\|_{\infty}, \varepsilon) \leq C_d \left(\frac{L}{\varepsilon} \right)^d, \quad 0 < \varepsilon \leq 1,$$

where C_d may depend on the ambient dimension d .

Giné-Zywn Symmetrization: Let \mathcal{F} be a class of measurable real-valued functions on Ω , and let $\varepsilon_1, \dots, \varepsilon_n$ be iid Rademacher random variables, independent of everything else. Then

$$\mathbb{E} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n (f(X_i) - \mathbb{E}f(X)) \right| \leq 2 \mathbb{E} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(X_i) \right|.$$

Empirical and expected Rademacher complexity:

$$\widehat{\mathfrak{R}}_n(\mathcal{F}) = \mathbb{E}_{\varepsilon} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(X_i) \right| \implies \mathfrak{R}_n(\mathcal{F}) = \mathbb{E}_{X_{1:n}} \widehat{\mathfrak{R}}_n(\mathcal{F})$$

Empirical L_2 pseudometric:

$$\|f - g\|_{L_2(\mu_n)} := \left(\frac{1}{n} \sum_{i=1}^n (f(X_i) - g(X_i))^2 \right)^{1/2}.$$

Conditional chaining bound for Rademacher averages:

$$\widehat{\mathfrak{R}}_n(\mathcal{F}) \leq \frac{C}{\sqrt{n}} \int_0^{\text{diam}(\mathcal{F}, L_2(\mu_n))} \sqrt{\log \mathcal{N}(\mathcal{F}, L_2(\mu_n), \varepsilon)} d\varepsilon$$

Finite-class Massart-type bound: \mathcal{F} is finite, and $\forall f \in \mathcal{F}$,

$$\|f\|_n^2 = \frac{1}{n} \sum_{i=1}^n f(x_i)^2 \leq r^2.$$

Then

$$\widehat{\mathfrak{R}}_n(\mathcal{F}; x_1, \dots, x_n) \leq r \sqrt{\frac{2 \log |\mathcal{F}|}{n}}.$$

7. Boolean Classes/VC Bounds

Definition: Let \mathcal{H} be a class of Boolean functions on Ω . Given sample points $x_1, \dots, x_n \in \Omega$, define the set of **traces** of \mathcal{H} on the sample by

$$\mathcal{H}|_{x_{1:n}} := \{(h(x_1), \dots, h(x_n)) : h \in \mathcal{H}\} \subseteq \{0, 1\}^n.$$

$\mathcal{H}|_{x_{1:n}}$ records all labelings that the class can realize on sample.

Growth function: For a class \mathcal{H} of Boolean functions,

$$\Pi_{\mathcal{H}}(n) := \sup_{x_1, \dots, x_n \in \Omega} |\mathcal{H}|_{x_{1:n}}|.$$

Pajor Lemma: \mathcal{H} is finite class of Boolean functions on Ω ,

$$|\mathcal{H}| \leq \#\{\Lambda \subseteq \Omega : \Lambda \text{ is shattered by } \mathcal{H}\}$$

Sauer-Shelah Lemma: Let $\text{VC}(\mathcal{H}) = d < \infty$,

$$\Pi_{\mathcal{H}}(n) \leq \sum_{k=0}^d \binom{n}{k} \leq \left(\frac{en}{d} \right)^d \quad (n \geq d)$$

Preliminary VC bound for uniform empirical error:

Let \mathcal{H} be a class of Boolean functions on Ω , and $\text{VC}(\mathcal{H}) = d$. Let X_1, \dots, X_n be i.i.d. with law μ , and assume $n \geq d$. Then

$$\mathbb{E} \sup_{h \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n h(X_i) - \mathbb{E}h(X) \right| \leq C \sqrt{\frac{d \log(en/d)}{n}},$$

Good to remember: Cond. on the sample $X_{1:n}$, the process

$$R_h := \frac{1}{n} \sum_{i=1}^n \varepsilon_i h(X_i), \quad h \in \mathcal{H},$$

has sub-Gaussian increments with respect to

$$\frac{1}{\sqrt{n}} \|h - g\|_{L_{\infty}(\mu_n)}.$$

8. Statistical Learning

Population and empirical risk:

$$R(h) := \mathbb{E}[\ell(h(X), Y)], \quad R_n(h) := \frac{1}{n} \sum_{i=1}^n \ell(h(X_i), Y_i)$$

Population and empirical risk minimizer:

$$h^* \in \arg \min_{h \in \mathcal{H}} R(h), \quad \widehat{h}_n \in \arg \min_{h \in \mathcal{H}} R_n(h)$$

Excess Risk bound: For any hypothesis class \mathcal{H} ,

$$R(\widehat{h}_n) - R(h^*) \leq 2 \sup_{h \in \mathcal{H}} |R_n(h) - R(h)|$$

Finite-class generalization bound: For every $\delta \in (0, 1)$, with probability at least $1 - \delta$,

$$R(\widehat{h}_n) - R(h^*) \leq 2 \sqrt{\frac{\log(2|\mathcal{H}|/\delta)}{2n}}$$

and using the tail integral identity gives

$$\mathbb{E}[R(\widehat{h}_n) - R(h^*)] \lesssim \sqrt{\frac{\log |\mathcal{H}|}{n}}.$$

Loss class has the same trace complexity: The loss class:

$$\mathcal{L} := \{(x, y) \mapsto \mathbf{1}_{\{h(x) \neq y\}} : h \in \mathcal{H}\}.$$

Then for every sample $(x_1, y_1), \dots, (x_n, y_n)$,

$$\left| \mathcal{L}|_{(x_i, y_i)_{i=1}^n} \right| = |\mathcal{H}|_{x_{1:n}}|.$$

Consequently,

$$\Pi_{\mathcal{L}}(n) \leq \Pi_{\mathcal{H}}(n) \quad \text{and} \quad \text{VC}(\mathcal{L}) \leq \text{VC}(\mathcal{H}).$$

VC Generalization Bound: Assume $n \geq d$,

$$\mathbb{E}[R(\widehat{h}_n) - R(h^*)] \lesssim \sqrt{\frac{\text{VC}(\mathcal{L}) \log(en/\text{VC}(\mathcal{L}))}{n}}.$$

One-dimensional Lipschitz regression: with squared loss

$$\mathbb{E}[R(\widehat{h}_n) - R(h^*)] \lesssim \sqrt{L/n}$$

Higher-dimensional Lipschitz regression: for all $\varepsilon \in (0, 1)$

$$\mathbb{E}[R(\widehat{h}_n) - R(h^*)] \lesssim \varepsilon + \sqrt{\frac{\log \mathcal{N}(\mathcal{F}_{L, d}, \|\cdot\|_{\infty}, \varepsilon)}{n}}.$$