

A A list of symbols used

Symbols

Norms

$\|\cdot\|_p$

ℓ_p -norm for sequences.

$\|\cdot\|_F$

Frobenius norm for matrices.

$\|\cdot\|$

Spectral norm for matrices.

$d_{TV}(\cdot, \cdot)$

Total variation distance between two distributions.

Sets

\mathbb{S}^{n-1}

The n -dimensional unit sphere

Δ^{n-1}

Probability simplex of dimension $n - 1$.

$\mathbb{A}^{n \times (n-1)}$

Set of n -dimensional left-stochastic matrices.

$\mathbb{N}_0 \triangleq \{0, 1, 2, \dots\}$

Set of nonnegative integers.

$\mathbb{N}_+ \triangleq \{1, 2, 3, \dots\}$

Set of strictly positive integers.

$\mathbb{R} \triangleq (-\infty, \infty)$

Set of reals.

$[n] = \{1, 2, \dots, n\}$

Set of integers 1 through n .

$\text{Perm}(n)$

All permutations of length n .

Asymptotics

$f(n) = \omega(g(n))$

Small-omega notation.

$f(n) = \Omega(g(n))$

$\liminf_{n \rightarrow \infty} f(n)/g(n) = \infty$

Big-omega notation.

$f(n) = O(g(n))$

$\liminf_{n \rightarrow \infty} f(n)/g(n) > 0$

Big-o notation.

$f(n) \sim g(n)$

$\limsup_{n \rightarrow \infty} f(n)/g(n) < \infty$

Asymptotic equivalence.

$f(n) = o(g(n))$

$\lim_{n \rightarrow \infty} f(n)/g(n) = 1$

Little-o notation.

$X_n = \Omega_{\mathbb{P}}(a_n)$

$\limsup_{n \rightarrow \infty} f(n)/g(n) = 0$

Omega-p notation.

$X_n = O_{\mathbb{P}}(a_n)$

$\forall \varepsilon \exists \delta_{\varepsilon, N_{\varepsilon}} : \mathbb{P}[|X_n/a_n| \leq \delta_{\varepsilon}] \leq \varepsilon \forall n > N_{\varepsilon}$

Stochastic boundedness.

$X_n \asymp_{\mathbb{P}}(a_n)$

$\forall \varepsilon \exists \delta_{\varepsilon, N_{\varepsilon}} : \mathbb{P}\left[\left|\frac{X_n}{a_n}\right| \geq \delta_{\varepsilon}\right] \leq \varepsilon \forall n > N_{\varepsilon}$

Equivalence-p notation.

$X_n = o_{\mathbb{P}}(a_n)$

$\forall \varepsilon \exists \delta_{\varepsilon}^-, \delta_{\varepsilon}^+, N_{\varepsilon} : \mathbb{P}[\delta_{\varepsilon}^- \leq |X_n/a_n| \leq \delta_{\varepsilon}^+] \geq 1 - \varepsilon \forall n > N_{\varepsilon}$

Convergence in probability.

“an absolute constant”

$\mathbb{P}\left[\left|\frac{X_n}{a_n}\right| \geq \delta\right] \rightarrow 0 \forall \delta > 0 \Leftrightarrow \forall \varepsilon, \delta \exists N_{\varepsilon, \delta} : \mathbb{P}\left[\left|\frac{X_n}{a_n}\right| \geq \delta\right] \leq \varepsilon \forall n > N_{\varepsilon, \delta}$

A constant independent of n .

Block Markov Chain

$n \in \mathbb{N}_+$

Number of states.

$K \in \mathbb{N}_+$

Number of clusters.

$\mathcal{V}_1, \dots, \mathcal{V}_K \subseteq [n] \triangleq \mathcal{V}$

Clusters, and set of all states.

$\alpha_1, \dots, \alpha_K \in (0, 1)$

$\mathcal{V} = \cup_{k=1}^K \mathcal{V}_k$, and $\mathcal{V}_a \cap \mathcal{V}_b = \emptyset \forall a \neq b$

Relative sizes of the clusters.

$\alpha_{\min} > 0, \alpha_{\max}$

$\alpha_k \triangleq \lim_{n \rightarrow \infty} |\mathcal{V}_k|/n$.

Minimum and maximum relative sizes of the clusters.

$\sigma : [n] \rightarrow [K]$

Cluster assignment.

$\{X_t\}_{t \geq 0}$

Markov chain.

$P \in \mathbb{A}^{K \times (K-1)}, P, Q \in \mathbb{A}^{n \times (n-1)}$

Transition matrices.

$\eta > 1$

$P_{x,y} \triangleq \frac{P_{\sigma(x), \sigma(y)}}{|\mathcal{V}_{\sigma(y)}| - \mathbb{1}[\sigma(x) = \sigma(y)]} \mathbb{1}[x \neq y] \forall x, y \in \mathcal{V}$

Assumed separability constant.

$\pi \in \Delta^{K-1}, \Pi \in \Delta^{n-1}$

$\exists 1 < \eta : \max_{a,b,c} \{p_{b,a}/p_{c,a}, p_{a,b}/p_{a,c}\} \leq \eta$

(Limiting) Equilibrium distribution.

$\Pi_x \triangleq \lim_{t \rightarrow \infty} \mathbb{P}[X_t = x] \forall x \in \mathcal{V}$

$\pi_k \triangleq \lim_{n \rightarrow \infty} \sum_{x \in \mathcal{V}_k} \Pi_x \forall k \in [K]$

Time-reversed transition matrix.

$P^* \in \mathbb{A}^{n \times (n-1)}$

$P_{x,y}^* \triangleq \frac{P_{x,y} \Pi_y}{\Pi_x}$

Dobrushin's ergodic coefficient.

$\delta(P) \in (0, \infty)$

$\delta(P) \triangleq \frac{1}{2} \sup_{x,y \in \mathcal{V}} \sum_{z \in \mathcal{V}} |P_{x,z} - P_{y,z}|$

$$0 \leq t_{\text{mix}}(\varepsilon) \leq -c_{\text{mix}} \ln \varepsilon$$

γ_{ps}

Information bound

$$\mathbb{P}_P[X_0 = x_0, \dots, X_T = x_T]$$

V^*

Φ

Ψ

\mathcal{Q}

$\mathcal{Q}(k, l)$

$$L \triangleq \ln \frac{\mathbb{P}_Q[X_0, X_1, \dots, X_T]}{\mathbb{P}_P[X_0, X_1, \dots, X_T]}$$

$$I_c(q||p)$$

$I_{a,b}(q||p)$

$I_{a,b}(\alpha, p)$

$$0 \leq J(\alpha, p) \leq I(\alpha, p) < \infty$$

Generic algorithm quantities

$T \in \mathbb{N}_+$

X_0, X_1, \dots, X_T

$\hat{N} \in \mathbb{N}_0^{n \times n}$

$\hat{\mathcal{V}}_1, \dots, \hat{\mathcal{V}}_K$

$\gamma^{\text{opt}} \in \min_{\gamma \in \text{Perm}(K)} \left| \bigcup_{k=1}^K \hat{\mathcal{V}}_{\gamma(k)} \setminus \mathcal{V}_k \right|$

$\mathcal{E} \subset \mathcal{V}$

Regime terminologies

“asymptotically accurate detection”

“asymptotically exact detection”

“dense regime”

“critical regime”

“sparse regime”

Spectral clustering algorithm

$\Gamma \subseteq \mathcal{V}$

\hat{N}_Γ

$U \Sigma V^T$

$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n \geq 0$

$\hat{R} \triangleq \sum_{k=1}^K \sigma_k U_{\cdot, k} V_{\cdot, k}^T$

Mixing time.

$t_{\text{mix}}(\varepsilon) \triangleq \min\{t \geq 0 : \sup_{x \in \mathcal{V}} d_{\text{TV}}(P_{x, \cdot}^t, \Pi) \leq \varepsilon\}$

We prove the second inequality, with c_{mix} being an absolute constant.

Pseudo spectral gap.

$\gamma_{\text{ps}} \triangleq \max_{i \geq 1} \frac{1 - \lambda((P^*)^i P^i)}{i}$

Probability of a sample path.

$\mathbb{P}_P[X_0 = x_0, \dots, X_T = x_T] \triangleq \prod_{t=1}^T P_{x_{t-1}, x_t}$

Vertex chosen uniformly at random from two clusters.

$V^* \stackrel{d}{=} \text{Unif}(\mathcal{V}_a \cup \mathcal{V}_b), a, b \in [K], a \neq b$

Probability measure of the true model.

I.e., under P and cluster assignments $\mathcal{V}_1, \dots, \mathcal{V}_K$.

Probability measure of the modified model.

I.e., under Q which is a perturbation of P constructed after the random vertex V^* was moved into its own cluster.

Set of all possible change-of-measure parameters.

$\mathcal{Q} \triangleq \{(q_{k,0}, q_{0,k})_{k=0, \dots, K} \in (0, \infty) \mid q_{0,0} = 0, \sum_{l=1}^K q_{0,l} = 1\}$

Sets of change-of-measure parameters leading to confusion between assigning to either cluster k, l .

$\mathcal{Q}(k, l) \triangleq \{q \in \mathcal{Q} \mid I_k(q||p) = I_l(q||p)\} \neq \emptyset, k \neq l$

Log-likelihood ratio.

Leading order coefficient in an asymptotic expansion of the log-likelihood ratio.

$I_c(q||p) \triangleq \lim_{n \rightarrow \infty} \frac{n}{T} \mathbb{E}_Q[L \mid \sigma(V^*) = c]$

Deconditioned leading order coefficient in an asymptotic expanding of the log-likelihood ratio.

$I_{a,b}(q||p) \triangleq \lim_{n \rightarrow \infty} \frac{n}{T} \mathbb{E}_\Psi[L]$

Separation between cluster a and b .

$I_{a,b}(\alpha, p) \triangleq \sum_{k=1}^K \frac{1}{\alpha_a} \left(\pi_a p_{a,k} \ln \frac{p_{a,k}}{p_{b,k}} + \pi_k p_{k,a} \ln \frac{p_{k,a}}{p_{k,b}} \right) + \left(\frac{\pi_b}{\alpha_b} - \frac{\pi_a}{\alpha_a} \right)$

Information quantities.

$J(\alpha, p) \triangleq \min_{k \neq l} \min_{q \in \mathcal{Q}(k,l)} I_{k,l}(q||p)$

$I(\alpha, p) \triangleq \min_{a \neq b} I_{a,b}(\alpha, p)$

Observation length.

Sample path of our Markov chain.

Observation matrix.

Each entry contains the number of times the Markov chain jumped from x to y , i.e.,

$\hat{N}_{x,y} \triangleq \sum_{t=0}^{T-1} \mathbb{1}[X_t = x, X_{t+1} = y] \forall x, y \in \mathcal{V}$

Approximated cluster assignments.

Permutation that minimizes the overlap between the true clusters and a cluster assignment.

Set of misclassified vertices.

$\mathcal{E} \triangleq \bigcup_{k=1}^K \hat{\mathcal{V}}_{\gamma^{\text{opt}}(k)} \setminus \mathcal{V}_k$

$\mathbb{E}_P[|\mathcal{E}|] = o(n)$

$\mathbb{E}_P[|\mathcal{E}|] = o(1)$

$T = \omega(n \ln n)$

$T \sim cn \ln n$ for some absolute constant $c > 0$

$\omega(n) = T = o(n \ln n)$

Truncated vertices.

This set is obtained from \mathcal{V} by removing the $\lfloor n \exp(-(T/n) \ln(T/n)) \rfloor$ states with the highest numbers of visits in the observed sample path of length T .

Truncated observation matrix.

This matrix is obtained from \hat{N} by setting all entries on the rows and columns corresponding to setates not in Γ to zero.

Singular value decomposition of \hat{N}_Γ .

Singular values of \hat{N}_Γ .

Best rank- K approximation of \hat{N}_Γ .

\mathcal{N}_x

Neighborhood.

$$\mathcal{N}_x \triangleq \left\{ y \in \mathcal{V} \mid \sqrt{\|\hat{R}_{x,\cdot} - \hat{R}_{y,\cdot}\|_2^2 + \|\hat{R}_{\cdot,x} - \hat{R}_{\cdot,y}\|_2^2} \leq \frac{1}{n} \right\} \\ \left(\frac{T}{n} \right)^{3/2} \left(\ln \frac{T}{n} \right)^{4/3}$$

$z_1^*, \dots, z_K^* \in \mathcal{V}$

Iteratively constructed cluster centers.

Cluster improvement algorithm

$\hat{p}, \hat{\pi}, \hat{\alpha}$

Approximated BMC parameters.

$u_x^{[t]}(c)$

Approximated difference between two log-likelihood functions.

$$u_x^{[t]}(c) \triangleq \left\{ \sum_{k=1}^K (\hat{N}_{x, \mathcal{V}_k^{[t]}} \ln \hat{p}_{c,k} + \hat{N}_{\mathcal{V}_k^{[t]}, x} \ln \frac{\hat{p}_{k,c}}{\hat{\alpha}_c}) - \frac{T}{n} \cdot \frac{\hat{\pi}_c}{\hat{\alpha}_c} \right\}$$

$\mathcal{H} \subseteq \mathcal{V}$

Set of well-behaved vertices.