# 1 Basics

- General p-norm: $\|x\|_p = (\sum_{i=1}^n |x_i|^p)^{1/p}$
- Hoelder: $\|uv\|_1 \le \|u\|_p \|v\|_q$, $\quad \|u+v\|_p \le \|u\|_p + \|v\|_p$
- Triangle: $\|u+v\|_p \le \|u\|_p + \|v\|_p$
- Cauchy-Schwarz: $|\langle u,v\rangle|^2 \le \|u\|^2 \|v\|^2$
- Cauchy-Schwarz: $\mathbb{E}[X,Y]^2 \le \mathbb{E}[X^2]\mathbb{E}[Y^2]$
- Taylor: $f(x) \approx f(a) + \frac{f'(a)}{1!}(x-a) + \frac{f''(a)}{2!}(x-a)^2 + \dots$
  - $f(\boldsymbol{x}) \approx f(\boldsymbol{a}) + \frac{\partial f(\boldsymbol{x})}{\partial \boldsymbol{x}}\big|_a - \frac{1}{2}(\boldsymbol{x}-\boldsymbol{a})^\top \left(\frac{\partial^2 f(\boldsymbol{x})}{\partial \boldsymbol{x}\partial \boldsymbol{x}^\top}\right)\big|_a (\boldsymbol{x}-\boldsymbol{a})$
  - Power series of exp.: $\exp(x) := \sum_{k=0}^\infty \frac{x^k}{k!}$
- Entropy: $H(X) \equiv H(p_X) = \mathbb{E}_X[-\log \mathbb{P}(X=x)]$
  - $H(X \mid Y) = \sum_y \mathbb{P}(Y=y)H(X \mid Y=y) \le H(X)$
  - $H(X,Y) = H(X) + H(Y \mid X)$
  - $H(X \mid g(X)) \ge 0 \quad$ ○ $H(g(X) \mid X) = 0$
  - $H(cX) \begin{cases} = H(X) & \text{discrete} \\ = H(X) + \log|c| > H(X) & \text{continuous} \end{cases}$
- MI: $I(X;Y \mid Z) = H(X \mid Z) - H(X \mid Y, Z)$ (symmetric)
  - $I(X;Y) = D_{\text{KL}}(p(x,y) \| p(x)p(y)) \ge 0$
  - $I(X_1,\dots,X_n;Z) = \sum_{i=1}^n I(X_i;Z \mid X_1,\dots,X_{i-1})$
    Markov chain: $I(X_1;X_2,X_3,\dots) = I(X_1;X_2)$
  - $I(X,Y;Z) = I(X;Z) + I(Y;Z \mid X)$
- KL-divergence: $D_{\text{KL}}(p \| q) = \sum_x p(x) \log\left(\frac{p(x)}{q(x)}\right) \ge 0$
- $1 - z \le \exp(-z)$
- Jensen, $f(X)$ convex: $f(\mathbb{E}[X]) \le \mathbb{E}[f(X)]$

## 1.1 Calculus

- Partial: $\int uv' \, dx = uv - \int u'v \, dx$ $\quad$ • $\frac{\partial}{\partial x}\frac{g}{h} = \frac{g'h}{h^2} - \frac{gh'}{h^2}$
- $\frac{\partial}{\partial \boldsymbol{x}}(\|\boldsymbol{x}-\boldsymbol{b}\|_2) = \frac{\boldsymbol{x}-\boldsymbol{b}}{\|\boldsymbol{x}-\boldsymbol{b}\|_2}$ $\quad$ • $\frac{d}{dx}|x| = \frac{x}{|x|}$
- $\frac{\partial}{\partial X}\log|X| = X^{-\top}$ $\quad$ • $|X^{-1}| = |X|^{-1}$
- $\frac{\partial}{\partial \boldsymbol{x}}(\boldsymbol{b}^\top \boldsymbol{x}) = \frac{\partial}{\partial \boldsymbol{x}}(\boldsymbol{x}^\top \boldsymbol{b}) = \boldsymbol{b}$
- $\frac{\partial}{\partial \boldsymbol{x}}(\boldsymbol{b}^\top A\boldsymbol{x}) = A^\top \boldsymbol{b}$ $\quad$ • $\frac{\partial}{\partial X}(\boldsymbol{c}^\top X\boldsymbol{b}) = \boldsymbol{c}\boldsymbol{b}^\top$
- $\frac{\partial}{\partial X}(\boldsymbol{c}^\top X^\top \boldsymbol{b}) = \boldsymbol{b}\boldsymbol{c}^\top$ $\quad$ • $\frac{\partial}{\partial \boldsymbol{x}}(\boldsymbol{x}^\top \boldsymbol{x}) = 2\boldsymbol{x}$
- $\frac{\partial}{\partial \boldsymbol{x}}(\boldsymbol{x}^\top A\boldsymbol{x}) = (A^\top + A)\boldsymbol{x} \overset{A \text{ sym.}}{=} 2A\boldsymbol{x}$
- $\frac{\partial}{\partial X}Tr(X^\top A) = A$ $\quad$ • Trace trick: $\boldsymbol{x}^\top A\boldsymbol{x} = \dots$
  $\dots \overset{\text{inn. prod.}}{=} Tr(\boldsymbol{x}^\top A\boldsymbol{x}) \overset{\text{cycl. permut.}}{=} Tr(\boldsymbol{x}\boldsymbol{x}^\top A) = Tr(A\boldsymbol{x}\boldsymbol{x}^\top)$

# 2 Maximum Entropy Inference

Sample $c \sim p(\cdot \mid X)$ s.t. $H[p(\cdot \mid x)]$ is maximal,
$\mathbb{E}_{C|X}[R(C,X)] = \mu$ and $\sum_c p(c \mid X) = 1$.
$\implies$ **Gibbs dist.:** $p(c \mid X) = \frac{1}{Z(X)}\exp(-\beta R(c,X))$

**Free energy:** $F(X) := -\frac{1}{\beta}\log Z(X)$
$\iff$ $p(c \mid X) = \exp(-\beta[R(c,X) - F(X)])$
$\implies$ **entropy:** $H[c \mid X] = \beta \underbrace{\mathbb{E}_{C|X}[R(C,X)]}_{=\mu} - \beta F(X)$

**ME:** $\max H[c \mid X] \iff \max Z(X) \iff \min F(X)$

- Exp. generalisation costs: $\mathbb{E}_{X''}\mathbb{E}_X \overset{\mathbb{E}_C}{\mathbb{E}_{C|X'}}[R(c,X'')]$
- Min. out-of-sample descr. length per deg. of freedom
  $\min_{p(\cdot|\cdot)} \mathbb{E}_{X',X''}\mathbb{E}_{C|X'}\left[-\log \frac{p(c|X'')}{p(c)}\right]$ $\quad p(c) = \mathbb{E}_X[p(c \mid X)]$
  $\overset{\text{Jensen}}{\ge} \min_{p(\cdot|\cdot)} \mathbb{E}_{X',X''}\left[-\log \mathbb{E}_{C|X'}[p(c \mid X'')]\right] - H[c]$
  $= \max_{p(\cdot|\cdot)} \mathbb{E}_{X',X''}[e^{H[c]} \cdot \kappa(X',X'')]$

**PA:** $T^* = \arg\max_T \kappa(X',X'')$

- PA-kernel: $\kappa(X',X'') := \sum_c p(c \mid X')p(c \mid X'')$
- combined: $p(c \mid X',X'') \propto p(c \mid X')p(c \mid X'')$

# 3 Methods for intractable Gibbs distr.

## 3.1 Markov Chains

**Mixing time of MC:** $\|P^t(c,\cdot) - \pi(\cdot)\|_{TV} \le \epsilon$
$t_{mix} \propto \frac{1}{\lambda_1 - \lambda_2}$ where $1 = \lambda_1 > \lambda_2 \ge \dots$
*Well behaving* **Markov Chains** are
1. **irreducible:** can go from/to any state (n steps)
2. **aperiodic:** chain doesn't go back & forth forever
   ($\forall n > n(c,c')$ (no) path length n w/ non-zero prob.)
1. $\wedge$ 2. $\implies$ **unique stat. dist.** $p(c') = \sum_c \pi(c \mid c')p(c)$
1. $\wedge$ 2. $\wedge$ stat. $\implies \lim_{t\to\infty}\mathbb{P}[X_t = c] = p(c) \wedge$
$\lim_{t\to\infty}\frac{1}{t}\sum_{s=1}^t f(X_s) = \sum_c p(c)f(c)$

**DBE** $\pi(c' \mid c)p(c) = \pi(c \mid c')p(c') \implies$ **p stat.**

**MH:** $\lambda_2 = \max\left\{1 - \frac{q(y,x)}{p_x}, 1 - \frac{q(x,y)}{p_y}\right\} = 1 - \alpha - \beta$

## 3.2 Sampling and SA

**Metropolis-Hastings:** Assume $p(c) \propto f(c)$.
$\pi(c' \mid c) := \begin{cases} q(c' \mid c)A(c,c') & c \ne c' \\ 1 - \sum_{c'\ne c}q(c' \mid c)A(c,c') & \text{otw.} \end{cases}$
where $q(c' \mid c)$ : prob. to propose the move $c \to c'$,
and $A(c,c') := \min\left\{1, \frac{q(c|c')f(c')/\mathcal{Z}}{q(c'|c)f(c)/\mathcal{Z}}\right\}$ prob. accept move

**Metropolis Algorithm:** Assume $p(c) \propto f(c)$ and
$q(c' \mid c) = q(c \mid c')$, i.e. symmetric.
1. Define symmetric $\{q(\cdot \mid c)\}_{c\in\mathcal{C}}$ s.t. graph $G_q$ is connected and every vertex in $G_q$ has edge to itself.
2. $c_0^T \leftarrow \$$ $\quad$ while $T > \epsilon$ do:
   - for $t = 1,2,\dots N$, do:
     - $\tilde{c} \leftarrow q(\cdot \mid c_{t-1}^T)$ // sample
     - $b \leftarrow \text{Bern}\left(\min\left\{1, e^{-\frac{1}{T}[R(\tilde{c},X) - R(c_{t-1},X)]}\right\}\right)$
     - If $b = 1$ then $c_t^T \leftarrow \tilde{c}$ else $c_t^T \leftarrow c_{t-1}$.
   - $c_0 \leftarrow c_N^T$
   - $T \leftarrow \text{reduce}(T)$
   - $c_0^T \leftarrow c_0$

**Temperature:** high temperature $T \to$ closer to uniform i.e. worse ability to discriminate between good and bad models $\to$ more likely to accept moves i.e. exploration, not stuck in bad local minima
reduce temperature to find better local minima and get stuck there

## 3.3 Laplace's Method (Least angle clust.)

1. *Square the cost:* $e^{-\frac{1}{T}R(c,X)} = const \cdot e^{g(c)^\top g(c)}$
2. *Complete the square:* $\int e^{-\frac{1}{T}(y-g(c))^2}\, dy = (\pi T)^{d/2}$
   $\Rightarrow e^{g(c)^\top g(c)} = (\pi T)^{-d/2}\int \exp^{-y^\top y + 2y^\top g(c)}\, dy$
3. *Rewrite normalisation constant:*
   $Z = \sum_c e^{-\frac{1}{T}R(c,X)} = \dots = const\int e^{-\frac{1}{T}f(y)}\, dy$
4. *Apply Laplace's method:*
   If $f$ has unique min. $y_0$ and Hessian $H := \frac{\partial^2 f}{\partial y^2}\big|_{y_0}$
   $\int e^{-\frac{1}{T}f(y)}\, dy \overset{(T\to 0)}{\approx} e^{-\frac{1}{T}f(y_0)}\left|\frac{H}{2\pi T}\right|^{-1/2}$

## 3.4 Mean-field Approximation

**Idea:** Approximate $p_\beta$ (Gibbs) with a "simple", factorisable distribution $p = p_1 \cdots p_N$.

**Approach:** Minimise $D_{KL}(p \| p_\beta)$

$\iff$ Minimise **Gibbs free energy:**

$$G(p) = \tfrac{1}{\beta} D_{KL}(p \| p_\beta) + F(\beta) = \mathbb{E}_{c \sim p}[R(c)] - \tfrac{1}{\beta} H[p]$$

*Note:* $H[p] = \sum_{i=1}^{N} H[p_i]$ and $F(\beta) \le G(p)$

**Ising model:** $R(c \mid J) = -\tfrac{1}{2} \sum_{i,j} J_{ij} c_i c_j - \sum_i h_i c_i$

where $J_{ij}$: interaction between particles,
$h_i$: noisy image, $\sigma_i$: denoised image

**Problem:** $\frac{\partial G(p)}{\partial p_{i\ell}} = 0$ s.t. $\sum_{\ell'} p_{i\ell'} = 1 \; \forall i$

**Solution:** with the *mean field* $h_i = [\cdots h_{i\ell} \cdots]^\top$

$h_{i\ell} := \frac{\partial \mathbb{E}[R(c)]}{\partial p_{i\ell}} = \mathbb{E}_{c \sim p_{|i \to \ell}}[R(c)] \leftarrow$ object $i$ chooses class $\ell$

$p_{i\ell} = e^{-\beta h_{i\ell}} / Z_i$

**EM-like Algo:** Iteratively 1. Pick random $i$
2. $h_i^{new} \leftarrow p_j^{old}$ 3. $p_i^{new} \leftarrow h_i^{new}$ until converged.

### 3.4.1 Smooth $k$-means  <span style="color:gray">scr.20 (p. 39)</span>

$R(c \mid X) = \sum_i \|x_i - y_{c_i}\|^2 + \tfrac{\lambda}{2} \sum_i \sum_{j \in N(i)} \mathbb{I}_{\{c_i \ne c_j\}}$

where the second term measures #violations of these neighbourhood constraints.

$$\implies h_{i\ell} = \|x_i - y_\ell\|^2 + \lambda \sum_{j \in N(i)} p_{j\ell} + const_i$$

## 4 Deterministic Annealing  <span style="color:gray">($Z$ is tractable)</span>

**Lemma:** func's $\times$ domain $\to$ domain $\times$ co-dom.

$\mathcal{O}(K^N) \to \sum_c \prod_i \epsilon_{i,c(i)} = \prod_i \sum_k \epsilon_{ik} \leftarrow \mathcal{O}(NK)$

$p(c \mid \theta, X) = \prod_{i \le N} p_i(c(i) \mid \theta, X)$

where $p_i(k \mid \theta, X) \propto \exp(-\tfrac{1}{T} \|x_i - \theta_k\|^2)$

Max. entr. $\implies \frac{\partial \log Z}{\partial \theta_k} = \frac{\partial \sum_{i \le n} \log \sum_{\nu \le K} \exp(-\|x_i - \theta_\nu\|^2)}{\partial \theta_k} \triangleq$

$0 \implies \theta_k^* = \frac{\sum_i p_i(k \mid \theta^*, X) \cdot x_i}{\sum_i p_i(k \mid \theta^*, X)}$

---

```
do
```
**E-step:** $p_i(k \mid \theta^{old}, X) = \frac{\exp(-\frac{1}{T} \|x_i - \theta_k\|^2)}{\sum_{j \le K} \exp(-\frac{1}{T} \|x_i - \theta_j\|^2)}$

**M-step:** $\theta_k \leftarrow \frac{\sum_{i \le n} p_{ik} x_i}{\sum_{i \le n} p_{ik}}$

$\theta^{old} \leftarrow \theta$

until convergence of $\theta$

$\theta_k \leftarrow \theta_k + \epsilon$  <span style="color:gray">(noise s.t. centroids can separate)</span>

**Phase transitions:** For $T \to \infty$: $\theta_k^* = \overline{X} \; \forall k \le K$
Once $T = 2\lambda_{max}$, more centroids appear, where
$\lambda_{max} = $ max. eigenvalue of $\tfrac{1}{N} X^\top X$.  ($x_i$'s row-wise)

**DA vs MAP:**
1. MAP can get stuck in local maximum
2. MAP not robust against noisy data
3. DA guaranteed to obtain global optimum if annealing is slow enough and ergodicity is given
4. In DA $T>0$ gives entropic regularisation

## 5 Histogram Clustering

**Least Angle Clust. (LAC):** [Idea]
Similarity $S(x_i, x_j) = w_{ij} \cos(\phi_{ij}) = w_{ij} e_i \cdot e_j$ with
unit vectors $e_i := x_i / \|x_i\|$, e.g. choice $w_{ij} = \|x_i\| \cdot \|x_j\|$.

**Dyadic data:** $\mathcal{Z} = \{(x_{i(r)}, y_{j(r)}); 1 \le r \le \ell\}$

- prototype / "centroid": $q(y_j \mid \alpha)$
- empirical dist.: $\hat{p}(y_j \mid x_i) = \frac{\hat{p}(x_i, y_j)}{\hat{p}(x_i)}$ $\leftarrow$scr. (5.10) $\leftarrow$scr. (5.11)

Likelihood: $P(\mathcal{Z} \mid c, q) = \prod_{r \le \ell} p(x_{i(r)}, y_{j(r)} \mid c, q)$

$= {}^{\text{scr. (5.12)}}_{\cdots} = \prod_i \prod_j [q(y_j \mid c(i)) \cdot p(c(i)) \cdot p(x_i)]^{\ell \hat{p}(x_i, y_i)}$

*Assume* $p(\alpha) = 1/k$ and $\hat{p}(x_i) = 1/n$

$\Rightarrow$ **Cost:** $R^{hc}(c, q, \mathcal{Z}) = \tfrac{\ell}{n} \sum_{i \le n} D_{KL}[\hat{p}(\cdot \mid x_i) \| q(\cdot \mid c(i))]$

Solving the **Gibbs dist.** $p(c \mid q, \hat{p}) = \prod_{i \le n} P_{i, c(i)}$

via Lagrange yields $q^*(y_j \mid \alpha) = \frac{\sum_{i \le n} P_{i\alpha} \cdot \hat{p}(y_j \mid x_i)}{\sum_{i \le n} P_{i\alpha}}$ <span style="color:gray">Lemma 2 ch.3 p.36</span>

### 5.1 Information Bottleneck Method

Find efficient code $X \mapsto \hat{X}$ (codebook vector) and preserve relevant info. about context $Y$.

**Criterion:** $R^{IB}(q(\hat{x} \mid x)) = I(X; \hat{X}) - \beta I(\hat{X}; Y)$

---

**Markov chain:** $\hat{X} \xrightarrow{q(\hat{x} \mid x)} X \xrightarrow{p(y \mid x)} Y$
**Generation process:** w/ *distortion* $d(x, \hat{x}) = D_{KL}[\cdot]$

$\begin{cases} q_t(\hat{x} \mid x) & \propto q_t(\hat{x}) \cdot \exp(-\beta D_{KL}[p(y \mid x) \| p_t(y \mid \hat{x})]) \\ q_{t+1}(\hat{x}) & = \sum_x p(x) \cdot q_t(\hat{x} \mid x) \\ p_{t+1}(y \mid \hat{x}) & = \sum_x p(y \mid x) \cdot p(x) \cdot q_t(\hat{x} \mid x) / q_t(\hat{x}) \end{cases}$

### 5.2 Parametric Distributional Clustering

**Idea:** Use a mixture of Gaussian prototypes, i.e.

$$p(y_j \mid \nu) \equiv p(b \mid \nu) = \sum_{\alpha \le s} p(\alpha \mid \nu) \, G_\alpha(b).$$

$$x_i \xrightarrow{c(i) = \nu} \nu \xrightarrow{p(b \mid \nu)} \hat{p}(b \mid i)$$

*Note:* Feature values $y_j$ ("bins" $b$) only depend on cluster index $\nu$ and not explicitly on the site $x_i$!

**Notation:** $x_i \leftarrow i$, $y_j \leftarrow b$ (bins), $\nu \leftarrow$ clusters

**Likelihood:** (both equivalent if $p(i) = \tfrac{1}{n}$)

$P(X \mid c, \theta) = \prod_{i \le n} p(c(i)) \prod_{b \le m} [p(b \mid c(i))]^{\ell \hat{p}(i, b)}$,

$P(X, M \mid \theta) = \prod_{i \le n} \prod_{\nu \le k} \big[ p(\nu) \cdot \prod_{b \le m} p(b \mid \nu)^{n_{ib}} \big]^{M_{i\nu}}$

where $n_{ib}$: #occur. an observ. at site $i$ is inside $I_b$
$M_{i\nu} = p(\nu \mid i) \in \{0, 1\}$ clust. membersh. assign.

**Cost (IB):** $R^{PDC}(c, p_{\cdot \mid c}) = -\log P(X, M\theta) = \ldots$

$\ldots = -\sum_{i \le n} \big[ \log p_{c(i)} + \tfrac{\ell}{n} \sum_{b \le m} \hat{p}(b \mid i) \log p(b \mid c(i)) \big]$

**E-step:** $h_{i\nu} = -\log p_\nu - \sum_b \tfrac{\ell}{n} \hat{p}(b \mid i) \log p(b \mid \nu)$
$q_{i\nu} = \mathbb{E}[\mathbb{I}_{\{c(i) = \nu\}}] \propto \exp(-h_{i\nu} / T)$

**M-step:** $p_\nu = \tfrac{1}{n} \sum_{i \le n} q_{i\nu}$
No closed form sol. for $p(\alpha \mid \nu)$. Thus, iteratively optimize pairs s.t. $\sum_\alpha p(\alpha \mid \nu) = 1$.

## 6 Graph-based Clustering

**Non-metric relations:** might assume negative values or violate the triangular inequality.

**Setting:** objects $o_i, o_j \in \mathcal{O}$; relations with weights $\mathcal{D} := \{D_{ij}\}$ on the edges $(i, j)$.

- Cluster $\alpha$: $\mathcal{G}_\alpha \equiv \{o \in \mathcal{O} : c(o) = \alpha\}$
- Inter-cluster edges: $\mathcal{E}_{\alpha\beta} = \{(i, j) \in \mathcal{E} : o_i \in \mathcal{G}_\alpha \wedge o_j \in \mathcal{G}_\beta\}$
- cut$(A, B) = \sum_{i \in A, j \in B} W_{ij} \to$ weight matrix $W$

- $\text{assoc}(A, \mathcal{V}) = \sum_{i \in A, j \in \mathcal{V}} W_{ij} \quad \rightarrow$ total connection strength from nodes in $A$ to all nodes in the graph

## Correlation clustering:

Minimise the sum of *pairwise* intracluster distances.

$$R^{cc}(c; \mathcal{D}) = - \sum_{\nu \leq k} \sum_{(i,j) \in \mathcal{E}_{\nu\nu}} S_{ij} + \sum_{\nu \leq k} \sum_{\substack{\mu \leq k \\ \mu \neq \nu}} \sum_{(i,j) \in \mathcal{E}_{\nu\mu}} S_{ij}$$

$$= -2 \sum_{\nu \leq k} \sum_{(i,j) \in \mathcal{E}_{\nu\nu}} S_{ij} + \sum_{(i,j)} S_{ij}$$

$\hookrightarrow$ intra-cluster $\quad \hookrightarrow$ const

$$\overset{up\ to}{\underset{thresh.\ u}{\overset{*}{=}}} -\frac{1}{2} \sum_{\nu \leq k} \sum_{(i,j) \in \mathcal{E}_{\nu\nu}} (|S_{ij} - u| + S_{ij} - u)$$

$$+ \frac{1}{2} \sum_{\nu \leq k} \sum_{\substack{\mu \leq k \\ \mu \neq \nu}} \sum_{(i,j) \in \mathcal{E}_{\nu\mu}} (|S_{ij} + u| - S_{ij} - u)$$

$* :$ altern. def. where $\frac{1}{2}(|X| \pm X) = \max\{0, \pm X\}$

## Graph partitioning: $\quad D_{ij} \in \mathbb{R}$

$$R^{gp}(c; \mathcal{D}) = const - \sum_{\nu \leq k} \text{cut}(\mathcal{G}_\nu(\mathcal{D}), \mathcal{V} \setminus \mathcal{G}_\nu(\mathcal{D}))$$

$$= const + \sum_{\nu \leq k} \text{cut}(\mathcal{G}_\nu(\mathcal{S}), \mathcal{V} \setminus \mathcal{G}_\nu(\mathcal{S}))$$

**Bias in $R(c;\mathcal{D})$:** Cost should scale prop. to #objects, i.e. $R(c; \mathcal{D}) = \mathcal{O}(n)$. $\quad * :$ use $D_{ij} = D(1 - \delta_{ij})$

**Tipp:** $\frac{\text{cut}(\mathcal{G}_\alpha, \mathcal{V} \setminus \mathcal{G}_\alpha)}{\text{assoc}(\mathcal{G}_\alpha, \mathcal{V})} \overset{*}{=} \frac{n \cdot p_\alpha \cdot n(1 - p_\alpha) \cdot D}{n \cdot p_\alpha \cdot n \cdot D} = 1 - p_\alpha$

## 6.1 Pairwise Clustering

**Cost:** $R^{pc}(c; \mathcal{D}) = \sum_\alpha \sum_{(i,j) \in \mathcal{E}_{\alpha\alpha}} \frac{D_{ij}}{|\mathcal{G}_\alpha|} = \sum_\alpha \sum_{(i,j) \in \mathcal{E}_{\alpha\alpha}} |\mathcal{G}_\alpha| \frac{D_{ij}}{|\mathcal{E}_{\alpha\alpha}|}$

**Equivariance to $k$-means:** $\quad$ (if $D_{ij} = \|x_i - x_j\|^2$)

$$\sum_{i \leq n} \|x_i - y_{c(i)}\|^2 = \sum_{i \leq n} \sum_{j \leq n} \sum_{\alpha \leq k} \frac{\mathbb{I}_{\{c(i) = \alpha\}} \mathbb{I}_{\{c(j) = \alpha\}}}{|\mathcal{G}_\alpha|} D_{ij}$$

**Invariance properties:**

- Symmetrisation: $\quad R^{pc}(c; \mathcal{D}^s) \equiv R^{pc}(c; \mathcal{D})$
- Off-diagonal shift: $\quad R^{pc}(c; \tilde{\mathcal{D}}) = R^{pc}(c; \mathcal{D}) - \lambda_{\min} \cdot n$

**Theorem:** If $S^c$ is p.s.d., then $D$ derives from squared Eucl. space. $\implies$ Make $S$ **p.s.d.**: $\tilde{S} := S - \lambda_{\min} \mathbb{I}$

## Constant Shift Embedding:

1. **Symmetrise** $D \rightarrow D^s$: $\quad D_{ij}^s := \frac{1}{2}(D_{ij} + D_{ji})$
2. **Centralise** $D$, then $S$: $\quad X^c := Q X^s Q^\top$
   $Q = \mathbb{I} - \frac{1}{n} e_n e_n^\top \quad S^c = -\frac{1}{2} D^c$

---

$$X_{ij}^c = X_{ij} - \frac{1}{n} \sum_k X_{ik} - \frac{1}{n} \sum_k X_{kj} + \frac{1}{n^2} \sum_{k,\ell} X_{k\ell}$$

$\implies$ sum over column/rows = 0

3. **(Off-)Diagonal shift:** Find $\lambda_{\min}$ of $S^c$

$\tilde{S} := S^c - \lambda_{\min} \mathbb{I} \qquad \tilde{D} := D - \lambda_{\min}(\mathbf{1} - \mathbb{I})$

$\tilde{D}_{ij} = \tilde{S}_{ii} + \tilde{S}_{jj} - 2\tilde{S}_{ij} = \|x_i - x_j\|^2$

## Reconstruction:

1. EVD: $\tilde{S} = V \Lambda V^\top \quad$ via $\quad (\tilde{S} - \lambda \mathbb{I}) v \overset{!}{=} 0 \quad (|v| = 1)$
   where $\Lambda = \text{diag}(\lambda_1 \dots \lambda_n) \quad$ and $\quad V = [v_1 \dots v_n]$
2. Find $p$ s.t. $\lambda_1 \geq \dots \lambda_p > \lambda_{p+1} = \dots = \lambda_n = 0$
3. $\implies X_p = V_p(\Lambda_p)^{1/2} \quad$ (each row is a vector)
4. $\implies X_t = V_t(\Lambda_t)^{1/2} \quad$ (approx. & denoising)

## Cluster membership of new data:

*Note:* $S^{new}$ is def. by $D_{ij}^{new} = S_{ii}^{new} + \tilde{S}_{jj} - 2 S_{ij}^{new}$

1. $(S^{new})^c = -\frac{1}{2} \Big[ D^{new}(\mathbb{I}_n - \frac{1}{n} e_n e_n^\top)$

   $\qquad - \frac{1}{n} e_m e_n^\top + \tilde{D}(\mathbb{I}_n - \frac{1}{n} e_n e_n^\top) \Big]$

2. Project: $\quad X_p^{new} = (S^{new})^c V_p(\Lambda_p)^{-1/2}$
3. Assign: $\quad \hat{c}_i = \arg\min_c \|(x_p^{new})_i - y_{c(i)}\|$

## 7 Model Selection for Clustering

What is the appropriate #clusters $k$ for my data?

**General approach:** Measure quality (neg. log-likelihood) for different $k \quad \rightarrow$ **elbow**.

## 7.1 Complexity-based Model Selection

**Strategy:** add a complexity term to neg. log-likelihood

**Attention:** MDL/BIC rely on likelihood optimisation $\rightarrow$ not generally applicable

**Ocam's razor:** Choose the model that provides the shortest description of the data.

### 7.1.1 Min. Description Length (MDL)

Minimise **descr. length**: $-\log p(X \mid \theta) - \log p(\theta)$

Approx.: $\hat{k} \in \arg\min_k -\log p(X \mid \hat{\theta}) + \frac{k'}{2} \log n$

### 7.1.2 Bayesian Information Crit. (BIC)

Parametrise likelihood $p(X \mid M)$ by $\theta$:

---

$$p(X \mid M) = \int_{\Theta_M} \exp(\log p(X \mid M, \theta)) \cdot p(\theta \mid M) \, d\theta$$

Assume flat prior $p(\theta | M) \approx const$ and expand log-likelihood by ML estimator $\hat{\theta}$:

$\overline{\ell}(\theta) = \frac{\ell(\theta)}{n} = \frac{1}{n} \log p(X | M, \theta) \overset{i.i.d.}{=} \frac{1}{n} \sum_i \ell(\theta, X_i) \overset{Taylor}{\approx} \dots$

$\implies p(X \mid M) = const_2 \cdot \exp(\ell(\hat{\theta}) - \frac{k'}{2} \log n)$

where $k'$: dimension of (trainable) parameters

## 8 Model Validation

### 8.1 Stability-based Validation

**Stability:** Solutions on two data sets drawn from the same source should be similar.

### 8.2 Information-theoretic Validation

#### 8.2.1 Shannon's Channel Coding Thm.

- **Channel:** $(\mathcal{S}, \{p(\cdot \mid s)\}_{s \in \mathcal{S}})$, $\mathcal{S}$: alphabet
  - $\epsilon$-noisy binary channel: $p(\hat{s} \mid s) = \begin{cases} 1 - \epsilon & \text{if } \hat{s} = s \\ \epsilon & \text{if } \hat{s} \neq s \end{cases}$
- **Capacity:** $\text{cap} = \max_p I(S; \hat{S}) \rightsquigarrow p_S(s)$
- $(M, n)$-**code:** is a pair $(Enc, Dec) \qquad \leftarrow$ scr. p.87
  where $M$: #messages, $n$: code-length
  - **Rate:** $r = \frac{\log_2 M}{n} \Leftrightarrow M = \lfloor 2^{nr} \rfloor$
  - **Commu. err.:** $p_{err} := \max_{i \leq M} \mathbb{P}(Dec(\widehat{Enc(i)}) \neq i)$

Goal / **Best code:** $\lim_{n \to \infty} \frac{\log M}{n}$ s.t. $\lim_{n \to \infty} p_{err} \to 0$

**Asymptotic equiparition property (AEP):**

- $A_\epsilon^{(n)}$: Typical set of sequences $(s_1, \dots, s_n) \in \mathcal{S}^n$
  $\left| -\frac{1}{n} \log p_{S^n}(s^n) - H[S] \right| < \epsilon \qquad \leftarrow$ scr. p.89
- $\mathbb{P}\left((S^n, \hat{S}^n) \in A_\epsilon^{(n)}\right) \overset{n \to \infty}{\to} 1 \qquad \leftarrow$ scr. p.90
- $p_{err} \leq 2^{-n(\text{cap} - 3\epsilon - r)} \overset{n \to \infty}{\to} 0$ if $r < \text{cap}$

#### 8.2.2 Algorithm Validation

**Assumptions:**

- Exponential solution space, i.e. $\log|\mathcal{C}| = \mathcal{O}(n)$
- $\mathcal{A}$'s output is probabilistic, i.e. $p(\cdot \mid X')$

**Ideal variant:**

**Messages:** $\mathcal{M} = \{X_1', \dots, X_m'\}$

**Code:** $X_i' \xrightarrow{Enc_A} p(\cdot \mid X_i') \xrightarrow{C_A} p(\cdot \mid X_i'') \xrightarrow{Dec_A} \hat{X}$

**Empirical variant:**

**Messages:** $\mathcal{M} = \{\tau_1, \dots, \tau_m\}$ drawn u.a.r. from $\mathbb{T}$

- Require $\sum_\tau p(c \mid \tau \circ X') \approx \frac{|\mathbb{T}|}{|\mathcal{C}|} \pm \rho$ $\qquad \leftarrow$ scr. p.95

**Code:** $\tau_i \xrightarrow{Enc} p(\cdot \mid \tau_i \circ X') \xrightarrow{C_A} p(\cdot \mid \tau_i \circ X'') \xrightarrow{Dec} \hat{\tau}$

- $Enc_A$: encodes $\tau_i \in \mathcal{M}$ as $p(\cdot \mid \tau_i \circ X')$
- $Dec_A$: selects $\hat{\tau} = \arg\max_\tau \kappa(\tau_i \circ X'', \tau \circ X')$

 whereby $\kappa(X'', X') := \sum_c p(c \mid X'') p(c \mid X')$

**Asymptotic Equipartition Property (AEP):**

*AEP fulfilled* if $\log \kappa(X', X'') \overset{n \to \infty}{\to} \mathcal{E}$

 whereby $\mathcal{E} := \mathbb{E}_{X', X''}[\log \kappa(X', X'')]$

- $A_\epsilon^{(n)}$: set of $(\epsilon, n)$-typical pairs $X', X''$

 $|\log \kappa(X', X'') - \mathcal{E}| < \epsilon$

- $p_{\text{err}} \le P_{(n)}$ c.f. scr. (6.19) $\overset{n \to \infty}{\to} 0$ if $\frac{\log m}{\log |\mathcal{C}|} < I$

 where $I := \frac{1}{\log |\mathcal{C}|} \mathbb{E}_{X', X''}[\log(|\mathcal{C}| \kappa(X', X''))]$

## 8.3 Applications of PA

**PA:** *quantifies the amount of information that algorithms extract from phenomena.* $\to$ quantified by **capacity** (max. # distinguishable messages that can be communicated)

**Temperature:** $T^* = \arg\max_T \kappa(X', X'')$

**Cost functions:** Given $R_1(\cdot, \cdot), \dots, R_s(\cdot, \cdot)$
$$\max_{\ell \le s} \kappa_\ell(X', X'') = \max_{\ell \le s} \frac{1}{Z_{X'} Z_{X''}} \sum_c e^{-\frac{1}{T} R_\ell(c, X')} e^{-\frac{1}{T} R_\ell(c, X'')}$$

**Algorithms:** Many MST (min. spanning tree) algo's are **contractive** ($\to$ sequence of candidate sol's).

**Approximation Set Coding (ASC):**
$$p^{\text{ASC}}(c \mid X') = \begin{cases} 1/|G_\gamma(X')| & \text{if } c \in G_\gamma(X') \\ 0 & \text{otw.} \end{cases}$$

$$G_\gamma(X') := \left\{ c \in \mathcal{C} : R(c, X') - \min_{c \in \mathcal{C}} R(c, X') \le \gamma \right\}$$

1. Run $\mathcal{A}$ to compute $G_t^{\mathcal{A}}(X')$ and $G_t^{\mathcal{A}}(X'')$, for all $t$

2. $t^* = \arg\max_t \kappa(X', X'') = \arg\max_t \frac{|G_t^{\mathcal{A}}(X') \cap G_t^{\mathcal{A}}(X'')|}{|G_t^{\mathcal{A}}(X')| \cdot |G_t^{\mathcal{A}}(X'')|}$

3. $c^* \xleftarrow{\$ \text{ sample}} \text{Unif}\left( G_{t^*}^{\mathcal{A}}(X') \cap G_{t^*}^{\mathcal{A}}(X'') \right)$

## 9 Appendix

### 9.1 Tips and Tricks

**Complete the square:**
If $p(x) \propto \exp(-\frac{1}{2} x^\top A x + x^\top b)$,
then $p(x) = \mathcal{N}(x \mid A^{-1} b, A^{-1})$

**Constrained optimisation:**
*primal*: $\min_x f(x)$ s.t. $g_i(x) = 0$; $h_j(x) \le 0$

**Lagrangian:** with each $\alpha_j \ge 0$
$$\mathcal{L}(x, \lambda, \alpha) = f(x) + \sum_i \lambda_i g_i(x) + \sum_j \alpha_j h_j(x)$$
Solve: $\frac{\partial \mathcal{L}}{\partial x} = 0$; $g_i(x) = 0$; $\alpha_j \ge 0$; $h_j(x) \le 0$

If **Slater's cond.** holds, $\exists x : g_i(x) = 0, h_j(x) < 0$, then we can solve the *dual* instead:
$$\max_{\lambda, \alpha} \{ \min_x \mathcal{L}(x, \lambda, \alpha) \} \quad \text{s.t.} \quad \alpha_j \ge 0$$
Solve: $\frac{\partial \mathcal{L}}{\partial x} = 0$; $\frac{\partial \mathcal{L}}{\partial \lambda} = 0$; $\alpha_j h_j(x) = 0$; $\alpha_j \ge 0$

**Euler-Lagrange:** Find extrema of functional $\mathcal{F}[f] = \int G(x, f(x), f(x)) \, dx$, thus $\frac{\partial \mathcal{F}}{\partial f} \overset{!}{=} 0$.
If $G$ is twice diff'able, then
$$\frac{\partial \mathcal{F}}{\partial f} = \frac{\partial G}{\partial f(x)} - \frac{d}{dx}\left( \frac{\partial G}{\partial f'(x)} \right) \overset{(*)}{=} \frac{\partial G}{\partial f(x)}.$$
$(*)$: when $G$ does not depend on $f'$.

### 9.2 Approximations

**Laplace Approximation:** $\frac{df}{dx}\Big|_{x_0} = 0$
$$\implies \int_{\mathbb{R}} e^{Cf(x)} \, dx \approx \sqrt{2\pi} C \cdot |f''(x_0)| \cdot e^{Cf(x_0)}$$

**Hyperbolic Functions:**

- $\sinh(x) = \dfrac{e^x - e^{-x}}{2}, \quad \dfrac{d}{dx} \sinh(x) = \cosh(x)$
- $\cosh(x) = \dfrac{e^x + e^{-x}}{2}, \quad \dfrac{d}{dx} \cosh(x) = \sinh(x)$
- $\tanh(x) = \dfrac{\sinh(x)}{\cosh(x)} = \dfrac{e^x - e^{-x}}{e^x + e^{-x}}, \cosh^2(x) + \sinh^2(x) = 1$
- $\dfrac{d}{dx} \tanh(x) = 1 - \tanh^2(x) = \dfrac{1}{\cosh^2(x)} = \text{sech}^2(x)$