# Maximum Entropy Clustering   $O(N \cdot k^N) \to O(KN)$

Sum-product trick: $Z = (f_{11} z_1 \dots f_{n1}) + (f_{12} z_1 \dots f_{n1})$
$+ \dots + (f_{1K} f_{2K} \dots f_{nK}) = (f_{11} + f_{12} + \dots + f_{1K}) \dots (f_{n1} + f_{n2} + \dots + f_{nK})$

DB: $\exp(-R(c_2)/T)/Z = [\exp(-R(c_1)/T)/Z] \exp(R(c_2)-R(c_1)/T)$
$P(c_2) P(c_2 \to c_1) = P(c_1) P(c_1 \to c_2), R(c_1) \leq RC \leq c_2$

Gibbs: draw $\hat{\imath}$, $c(i) \sim P(c(i) | c(i), \dots c(i-1), c(i+1) \dots c(n))$
- Deterministic setting: $\nabla_\theta R = 0$
- Probabilistic setting: $E[\nabla_\theta R] = 0$

- kmeans: global minimum may be fragile (wrt noise), nearest neighbor rule is the culprit, very dependant on init. config., indep. makes altern. opt. tractable
- Fuzzy asgmt. algo. doesn't operate in the state of solutions, we get a summary statistic of all highly probable solutions (under Gibbs distr.)
- DA: $T \to 0$ is k-means, assign each point to closest $T \to \infty$ is uniform P, all centroids are equal to mean
- kmeans: hard assgmts are problematic, calculus doesn't represent doubt (uncertainty about assgmts
- Fuzzy: conditioning on prior nodes obts indep.
- $R^{kom}(c, \mu, X) = \sum_{i \in n} \sum_{v \in k} T_{c(i), v} D(x_i, \mu_v)$ - reduce data dims and display similarities among datapoints
- two sources: quantization error, codevector confusion
- $R^{cc}(c, \mu, X) = \frac{1}{n} \sum_{i \in n} D(x_i, \mu_{c(i)}) + \lambda \sum_{v \in k} P_v C(P_v)$ number
- $C^{ME}(P_v) = -\log P_v$ (entropy constr.) of clusters $\ln 4$
- $C^{lb}(P_v) = (P_v)^{-s}, s \geq 1$ (load balancing) set a priori $s = 1$ is k-means, pay price for heavily used indices
- $c^{opt}$: Huffman coding, $-\lim P_v \log P_v = 0$, optimal soln is uniform quant, codebook ignores data, we can invest a centroid for one datapoint
- least angle: groups cosine-similar vectors, gives too much weight to large clusters

$R^{la}(M, X) = -\frac{1}{n} \sum_{v \in k} \sum_{i \in n} \sum_{j \geq i} M_{iv} M_{jv} \cos \phi_{ij}$ ← on mm-
$\qquad = -\frac{1}{2n} \sum_{v \in k} \sum_{i \in n} \sum_{j \in n} M_{iv} M_{jv} e_i e_j + 1$ are the cost func
$\qquad = -\frac{1}{2n} \sum_{v \in k} (\sum_{i \in n} M_{iv} e_i)^2 + 1$ by cluster sizes

EM algo ⟹ init: $\theta^0$ init with random values
E-step: compute $p(\bar{z} | X, \theta^0)$
M-step: $Q \leftarrow \arg\max_\theta E_{p(\bar{z} | X, \theta^0)} [\log p(X, \bar{z} | \theta)]$
Repeat: (if $\theta^0$ and $\theta$ are close eno-gh finish) otherwise set $\theta^0 \leftarrow \theta$ and go to E-step

prereq: efficient to calc. $p(\bar{z} | X, \theta)$ for any $\theta$ efficient to calc. M-step for any $\theta_0$

---

Markov chains: irreducible if it is possible to get to any state from any state. Periodicity: a state $T$ has period $k$ if any return to state $T$ must occur in multiples of $k$ time steps. If $k=1$ state is aperiodic. Chain is aperiodic if every state is aperiodic. An irreducible chain needs one aperiod. state to imply all states are aperiodic.

Metro: always accepts improvements, accepts cost deteriorations with small prob. - prob. of being in a state is Boltzmann weight for that state - DB guarantees Gibbs distr. is the stationary distr., it is sufficient but not necessary cond.
- proposed distr. reflects search strategy, a flat prior doesn't exploit local structure, topology plays no role - areas of low cost doesn't mean solution look similar, then are only judged to be comparable
- simulated annealing: approx. global min in fixed time rather than finding precise local optimum
- MH: acceptance $A(x_t, x') = \min \{1, \frac{P(x') q(x' | x_t)}{P(x_t) q(x_t | x')}\}$
- Metropolis: MH with symmetric proposal: $q(x_t | x') = q(x' | x_t)$
- Simulated annealing with non-homogeneous Markov chain $P(x) \propto p^{1/T_t}(x)$, $T_0 = 1$, many cooling schedules
- exponential: $T(t) = T_0 \kappa^t$, linear: $T(t) = T_0 - \eta t$
- Gibbs: requires sampling from conditional posterior, reduces multidim sampling to sequence of 1D
- if the Markov chain is time-homogeneous, process is described by single, time-indep. matrix $P_{ij}$, $\pi$ is called stationary distribution s.t. $0 \leq \pi_i \leq 1, \sum_i \pi_i = 1$
$\pi_j = \sum_{i \in S} \pi_i p_{ij}, \to \pi P = \pi, P$ is transition matrix
$\pi = \frac{e}{\sum_i e_i}$ - $e$ is the left ergval of transition mat. $P$ with ergval 1
- Stationarity: $\sum_x p(x) P(x \to y) = p(y)$ (DB implies stat.)

For $X \in \mathbb{R}^{m \times n}$, $C_n X$ removes means from columns, $X C_n$ removes means from rows, scatter $S = X C_n (X C_n)^T = X C_n X$
$Q_n = I_n - \frac{1}{n} U(n), P^c = Q_n P Q_n, U(n) = 11^T$ (mat. of all 1's)

My implementation of CSE
1- Get diagonal (0 diag), 2- symmetrize $D$
3- Centralize $D$: $D^c \leftarrow Q D Q$, 4- similarities: $S^c \leftarrow -\frac{1}{2} D^c$
5- Shift sims: $\tilde{S}^c \leftarrow S^c - \lambda_n(S^c) I_n$, 6- $\tilde{D}_{ij} \leftarrow \tilde{S}_{ii}^c + \tilde{S}_{jj}^c - 2\tilde{S}_{ij}^c$
7- $\tilde{S}^c \leftarrow Q \tilde{S}^c Q$, 8- $\tilde{D}^c \leftarrow (-2) \tilde{S}^c$, 9- Assert $\tilde{D}^c = Q D Q$
- Shifted $D$ contains squared Euclidean distances in the high dimensional space

Gibbs free energy - goal is simultaneously max.ing entropy and minimizing cost, conv.g space is $X$, real-valued energy function is $E(x)$
- $G(p) = E_p[E] - \frac{1}{\beta} \cdot H(p) = \sum_x p(x) E(x) + \frac{1}{\beta} \sum_x p(x) \log p(x)$
- $P_B(x) = \exp[-\beta(E(x) - F(\beta))] = e^{-\beta E(x)} / [\sum_{x'} e^{-\beta E(x')}]$ Gibbs distr.
- $G(p) = \frac{1}{\beta} \cdot D_{KL}(p \| P_B) + F(\beta), G(p^* = P_B) = 0, G(p) \geq F(\beta)$
- $F(\beta) = -\frac{1}{\beta} \log Z(\beta), Z(\beta) = \sum_x e^{-\beta E(x)}$
- $F(\beta) = -\frac{1}{\beta} \log [\sum_x e^{-\beta E(x)}]$ can be derived from above using $\sum_x P_B(x) = 1$

Locally linear embedding - do local linear approx. at each point and smoothly interpolate them (as manifold locally looks linear)
- Find k-NN for each point, $Z$ contains neighbors of $X$, define local cover as $C = Z^T Z$, solve for $C w z \bar{1}$, assign normalized $w$ as weights of neighbors
- After computing all weights create sparse matrix $M = (I - W)^T (I - W)$ and get bottom $d+1$ eigenvector, discard the bottom (ergval 0). This gives the embed of original points
- using kNN means we take a fine-grained look when there is a lot of data, and coarse-grained view when there's little data
- LLE tends to handle non-uniform sample densities poorly because there is no fixed unit to prevent weights from drifting as various regions differ in sample densities

---

Linearization trick:
$\exp(b^2 / 2a^2) = \int_{-\infty}^\infty (a/\sqrt{2\pi}) \exp(-a^2 x^2 / 2 + bx) dx$
$\exp(-b^2 / 2a^2) = \int_{-\infty}^\infty (a/\sqrt{2\pi}) \exp(a^2 x^2 / 2 - bx) dx$
Gaussian integral: $\int_{-\infty}^\infty e^{-a(x+t)^2} dx = \sqrt{\pi/a}$
- Sum of norm. distr. vars: $X \sim N(\mu_X, \Sigma_X), Y \sim N(\mu_Y, \Sigma_Y)$
  $X + Y \sim N(\mu_X + \mu_Y, \Sigma_X + \Sigma_Y), Y = c + BX \sim N(c + B\mu_X, B\Sigma_X B^T)$
- $I(A;B) = E_{A,B}[\log(P(A,B)/P_A)P(B)]$
- $\frac{1}{2}(|x| \pm x) = \max\{0, \pm x\}$ (shifted correlation clusters)
- $\text{cut}(A, B) = \sum_{i \in A, j \in B} W_{ij}$.   $\text{assoc}(A, V) = \sum_{i \in A, j \in V} W_{ij}$
- assoc measures the total connection strength from nodes in $A$ to all nodes in the graph
- Glivenko-Cantelli: $P[\sup_{x \in \mathbb{R}} |F(x) - \hat{F}_n(x)| \xrightarrow{n \to \infty} 0] = 1$
- four axioms of surprise $S: [0, 1] \to [0, \infty)$
  - certainty: $S(1) = 0$, no surprise for certain event
  - anti-monotonicity: $p \leq q \Rightarrow S(p) \geq S(q)$
  - additivity: $S(p \cdot q) = S(p) + S(q)$
  - continuity: $S(p)$ is a cont. func. of $P$, no surprise jumps for infinitesimal changes in $p$

$f(x) = \frac{1}{\sqrt{2\pi \sigma^2}} \exp\left[-\frac{1}{2} \frac{(x - N)^2}{\sigma^2}\right] \to N(x | \mu, \sigma)$

$f(x) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left[-\frac{1}{2}(x - N)^T \Sigma^{-1} (x - N)\right] \to N(x | \mu, \Sigma)$

- binomial: $f(k, n, p) = Pr(X = k) = \binom{n}{k} p^k (1-p)^{n-k}$
- $Var[X] = \int_x (x - N)^2 p(x) dx$
- $Var[X] = E[(x - E(x))^2] = E[X^2] - E[x]^2$
- $Var[X + Y] = Var[X] + Var[Y] + 2\text{cov}[X, Y]$
- $Cov[X, Y] = E[(x - E(x))(Y - E(Y))]$
- $Cov[aX, bY] = ab \, Cov[X, Y]$
- $Cov(X, Y) = E[XY^T] - E[X] E[Y]$

Jensen's inequality: $X$ is r.v., $\varphi$ is convex →
$\varphi(E[X]) \leq E[\varphi(X)]$

Cauchy-Schwarz: $|E[X, Y]|^2 \leq E[X]^2 E[Y]^2$

$H(X) = -\sum_{\hat{x}}^{x} P(x_i) \log P(x_i)$   evaluating $X, Y$ simult.
$H(X, Y) = H(X | Y) + H(Y) = H(Y | X) + H(X)$ ↗
$H(X, Y) \leq H(Y) + H(X)$ → equal when $X, Y$ indep
cross-ent: $H(p, q) = -\sum_x p(x) \log q(x)$
cond. ent: $H(X | Y) = -\sum_{i,j} p(x_i, y_j) \log \frac{p(x_i, y_j)}{p(y_j)}$
cross-ent: $H(p, q) = E_p[-\log q] = H(p) + D_{KL}(p \| q)$
$I(X; Y) = D_{KL}(P(X, Y) \| P(X) P(Y))$   $H(p) \leq H(p, q)$ ✓
$I(X; Y) \geq 0, \quad I(X; Y) = I(Y; X)$
$I(X; Y) = H(X) - H(X | Y) = H(Y) - H(Y | X)$
$I(X; Y) = H(X, Y) - H(X | Y) - H(Y | X) = H(Y) + H(X) - H(X, Y)$
cond. mut. inf.: $I(X; Y | Z) = H(X | Z) - H(X | Y, Z)$
chain rule: $I(X; Y, Z) = I(X; Z) + I(X; Y | Z)$
$I(X; Y) = \sum_Y \sum_X P(X, Y) \log \frac{P(X, Y)}{P(X) P(Y)}$
$D_{KL}$ (relative ent.): $D_{KL}(p(x) \| q(x)) = -\sum_x p(x) \log \frac{q(x)}{p(x)}$
$\qquad = \sum_x P(x) \frac{p(x)}{q(x)}$
$D_{KL}(p(x) \| q(x)) \geq 0$

---

# Clustering distr. data - Histogram clustering

We partition probsimplex rather than Euclidean space
- hc: lacks topology as KL-div is invariant to permut. of feats, problem in ordered feature spaces (natural topology: frequency ordering, edge directions, color circle etc. - feature similarity)
- heuristic HBM: smoothing: distribute percentages to neighbors
- pdc: no longer permut. invariant, optimize through EM, no closed form soln for Gauss. weights or means, optimize $\mu_a$ through interval bisection, optimize weights pairwise: pick randomly two Gaussians shift weight from one to the other
- rectified Gaussians: feature domain has finite support, cut off the mass that goes outside and scale it up (model assmpt)
- hc notation ⟹ $i$:object, $j$:feature, $l$:#observations, $z$:dataset, $r$:observation, $q(j | a) = $ cluster feat. distr.
- pdc: feats depend only on cluster idx, not explicitly on site
- pdc update → E-step: $q_{iv} = E[\mathbb{I}_{c(i) = v}] \propto \exp(-h_{iv}/T)$
  $h_{iv} = -\log p_v - \sum_j \frac{1}{n} \hat{p}(y_j | x_i) \log(\sum_{v} p(\alpha | v) \hat{q}_a(j))$   $\sum_{v} p(v | u) = 1$
  M-step: $p_v = \frac{1}{n} \sum_i q_{iv}$, numerical solution for $p(\alpha | v)$
- IB: maximize $I(X; c)$: generic clusters, peculiarities of $x$ should be forgotten and what's generic of all the objects assigned to particular clusters should remain $\infty$ assigning all objects to one cluster, entropy death, no structure, counterforce: representations informative as possible relative to a feature space (quantization objective + distortion constr.)
  $R^{IB} = I(X; c) - \lambda I(C; Y)$, where $c_{opt} = \arg\min_c R^{IB}$
  $R(D) = \min_{\{P(c | x) : E_{x, c}(x, c) \leq D\}} I(X; c)$
  $F(p(c | x)) = I(X; c) + \beta E_{x, c}[d(x, c)]$

---

# Pairwise clustering (graph clustering)
- clustering algos always impose structure on data
- sim vs dissim: use large values only when you are pretty certain that these large values are doing the job, distances are difficult to measure accurately in most applics, very large mistakes with very large distances which are very rare might still influence your estimators
- dissim to sim: exponential scaling: $S_{ij} = \exp(-D_{ij}/\Delta)$ maps dissimilarities $D_{ij} \geq \Delta$ into interval $[0, e^{-1}]$, beneficial as large dists hard to measure - alt-conversion: linear mapping: $S_{ij} = \max_j D_{ij} - D_{ij}$
- cc: seems to put wrong emphasis in evaluation, large clusters have more weight (it is less expensive to add new nodes to big clusters than the small ones)
- shifted cc: count sims only relative to a threshold, $S \in [-1, +1]$
  $R^{scc}(c; D) = -\frac{1}{2} \sum_{v \in k} \sum_{(i, j) \in S_v} (|S_{ij} + u| + S_{ij} + u) + \frac{1}{2} \sum_{v \in k} \sum_{\substack{u \neq v}} \sum_{(i, j) \in S_v} (|S_{ij} + u| - S_{ij} - u)$
- graph cut: strong balance for very unbalanced (very small/very large) clusters due to lack of normalization
  $R^{gc}(c; D) = \sum_{v \in k} \sum_{(i) \in S_v, (j) \in S_v} D_{ij} = \sum_{v \in k} D_{ij} - \sum_{v \in k} \sum_{(i, j) \in S_v} D_{ij}$
  $= \text{const} - \sum_{v \in k} \text{cut}(S_v(D), V \setminus S_v(D)) = \text{const} + \sum_{v \in k} \text{cut}(S_v(S), V \setminus S_v(S))$
  $R = \sum_{v \in k} \text{cut}(S_v, V \setminus S_v)/\text{norm} \to \text{normcut} = \text{assoc}(S_v, V)$
  $\qquad \to \text{avg} = |S_v|$, $\text{minmaxcut} = \text{assoc}(S_v, S_v)$
- MM-max cut has severe bias towards equipartitions
- PC shift-invariance: $\tilde{D}_{ij} = D_{ij} + D_0(1 - \delta_{ij}) \to R^{pc}(c; \tilde{D}) = R^{pc}(c; D) + \beta n$
- symmetry-invariance: $R(c, \frac{1}{2}(D_{ij} + D_{ji})) = R(c; D)$   $\frac{1}{\text{off-diagonal}}$
- $R^{pc} = R^{km}$ if $D_{ij} = \|x_i - x_j\|^2$
- CSE: $\tilde{D} = D + D_0(e_n e_n^T - I_n)$ corresponds to squared Euclidean distance where $D_0 = -2\lambda_n(S^c)$ is the minimal constant

## Model selection validation for clustering

- Structure can be invalid in two ways: wrong model order and/or inappropriate model type
- validation methods can be external (= comp. with ground truth) or internal
- we want to use convergence guarantees for sums (e.g. law of large numbers) so we convert the log product term to a big sum over exponent
- Every method for valid. introduces a bias
- you may be grossly misled by naively counting the number of params (cf. VC-dim of sine)
- General approach: measure quality for diff. $k$ with some discount (re. complexity penalty) — will more $k$ use points more bits, pays on fit better
- MDL: minimize $\{-\log p(X|\hat{Q}_k) - \log p(\hat{Q}_k)\}$  $\hat{Q}_k$: MLE of $Q$
  $= \{-\log p(X|\hat{Q}_k) + 1/2 \ell \log n$  $\ell$: #params of $Q_k$
- Bayes factor: $P(X|M_k)/P(X|M_l)$ where
  $p(X|M_k) = \int p(X|Q_k, M_k) p(Q_k|M_k) dQ_k$, use Lapl. approx.
  $\log p(X|M_k) = \log p(X|\hat{Q}_k, M_k) - \ell/2 \cdot \log n + O(1)$
- MDL and BIC are formally equivalent, consistent
- In modern statistics more d.o.f. we have as we have more data — BIC exploits finite dim. param space and growing data (samples $n$) — classical limit of statistics, we understand it well, but not useful (parametric statistics)
- They are well-motivated, they do likelihood based optimization instead of posterior based, which is intrinsically hard
- Gap statistic: find "knee" in costs, it uses max. discrepancy between actual data cost and that of unstructured ref. data (re. which cannot be clustered, so-called "Null-model"), if the data are too noisy, the method may fail
  $gap_n(k) = E_n^*[\log(W_k)] - \log(W_k)$, $D_{ij} = \|x_i - x_j\|^2$,
  $W_k = \sum_{r=1}^{k} 1/2n_r \sum_{(i,j) \in C_r} D_{ij}$, $E_n^*$ is the expect. value w.r.t. sample of size $n$ from null-model
  $\hat{k} = \min \{k | gap_n(k) \geq gap_n(k+1) - b_{k+1}$
- In practice approx. $E_n^*[\log(W_k)]$ by bootstrap.
- gap stat. works satisfactorily for spherical and well-separated clusters (it assumes compact clusters), a structural bias, it is a good heuristic for k-means-like criteria

Stability based valid: idea → solutions on two datasets from same source should be similar, signal in both should be same but fluctuations should be different, model mismatch in soln transfer produces unstable clustering solutions
- type of classifier has a large influence, has to be selected with care (e.g. path-based clust.)
- data with a nearest centroid classifier
- two disjoint sets of equal size: overlap could determine group structure (stat. dependence), algo. should be able to find similar structure in both, cannot find if one is too small and such structure is no longer visible to algo.
- good performance in experimental datasets, has systematic bias towards too simplistic solutions — it is principled but with a deficit: it only tells you how reliable under repeated experimentation the algo. would react, but it doesn't take into account the informativeness of solution (disregards the tradeoff) — A tolerable decrease in stability of inferred patterns might be compensated by a substantial increase of their information content

---

## Stability valid. procedure:

1- Transfer via prediction: construct classifier $\phi$ trained on $(X, Y)$ where $Y$ is the clustering solution to $X$ as $Y = A_k(X)$. We consider labeling $\phi(X') := (\phi(X'_i))_{i \in Y}$ as the extension of the clustering solution $Y$ on dataset $X$ to dataset $X'$. These predicted labels can be compared to those generated on the clusters algo $A_k(X')$

2- Compare solutions: a very natural distance measure for comparing labelings $\phi(X)$ and $Y'$ is Hamming dist. [0-1 loss], this can be treated as the empirical misclassification risk of $\phi$ with regard to the training set

3- Permute solutions: to overcome non-uniqueness we optimally permute solutions to maximize agreement. Hungarian method does this in $O(n+k^3)$ by minimum bipartite matching
- Dissim.: $d_{\phi_k}(\phi(X'), Y') := \min_{\pi \in S_k} \frac{1}{n} \sum_{i=1}^{n} 1\{\pi(\phi(X'_i)) \neq Y'_i\}$
- Stability: $S(A_k) := E_{X, X'} d_{\phi_k}(\phi(X'), Y')$
  smaller the values of the index $S(A_k) \in [0,1]$, the more stable are the solutions

4- Calculate relative stability w.r.t. stability achieved by random label guessing:
  $\tilde{S}(A_k) := S(A_k)/S(R_k)$, $R_k$ assigns labels to objs. with prob. 1/k

- $D_{KL}$ penalizes approximations that underestimate likely events
- entropy of $X \approx$ uncertainty of $X \approx$ amount of info. of $X \approx$ minimum expected number of questions to guess $x$
- Forward KL, difference between $p(x)$ and $q(x)$ is weighted by $p(x)$. During the optimization process then, wherever $p(x) = 0$, $q(x)$ would be ignored. The difference between $p(x)$ and $q(x)$ will be minimized if $p(x) > 0$. It is known as zero-avoiding, as it is avoiding $q(x) = 0$ whenever $p(x) > 0$
  → $P_{M,\theta}$ as mass-seeker: $\arg\min_{P_{M,\theta}} D_{KL}(P \| P_{M,\theta})$
- Reverse KL, as we switch the two distr. position in the equation, now $q(x)$ is the weight. Here, it is better to fit some portion of $p(x)$ as long as that approx. is good. Consequently reverse KL will try to avoid spreading the approximate. As those properties suggest, this form of KL-div is known as zero-forcing as it forces $q(x)$ to be 0 on some areas, even if $p(x) > 0$
  → $P_{M,\theta}$ as mode-seeker: $\arg\min_{P_{M,\theta}} D_{KL}(P_{M,\theta} \| P)$

$D_{KL}(P \| P_{M,\theta})$:


mass-seeker

$D_{KL}(P_{M,\theta} \| P)$:


mode-seeker

---

## Model valid by info. theory

Typical set: Asymptotic equipartition property (AEP):
$-\frac{1}{n} \log P(X_1, ..., X_n) \xrightarrow{\text{in prob}} H(x)$, for i.i.d. $X_1, ..., X_n \sim P(X)$

A typical set $A_\epsilon^{(n)}$ wrt $P(x)$ is a set of sequences $(x_1, ..., x_n) \in X^n$ with property
$2^{-n(H(x) + \epsilon)} \leq P(x_1, ..., x_n) \leq 2^{-n(H(x) - \epsilon)}$
- If $(x_1, ..., x_n) \in A_\epsilon^n$ then $H(x) - \epsilon \leq -\frac{1}{n} \log P(x_1, ..., x_n) \leq H(x) + \epsilon$
- $P\{A_\epsilon^{(n)}\} > 1 - \epsilon$ for $n$ sufficiently large
- $|A_\epsilon^{(n)}| \leq 2^{n(H(x) + \epsilon)}$, where $|A_\epsilon^{(n)}|$ is the cardinality of the typical set
- $|A_\epsilon^{(n)}| \geq (1-\epsilon) 2^{n(H(x) - \epsilon)}$ for $n$ sufficiently large

The typical set $A_\epsilon^{(n)}$ has prob almost 1, all elements of the typical set are nearly equiprobable, and the number of typical solutions is nearly $2^{nH}$

Approx. weights $w: C \times X \times R_+ \to [0,1]$, s.t. $(c, X, \beta) \mapsto w_\beta(c, x)$ weights are non-neg, maximal weight allocated to global minimizer at and it is normalized to one: $w_\beta(c^\perp, X) = 1$. Solutions with large approx. weights $w_\beta(c, x) \geq 1-\epsilon$, $\epsilon \ll 1$ can be accepted as substitutes of the global minimizers. Posterior becomes
$p(c|x) = \frac{w_\beta(c, x)}{\sum_{c'} w_\beta(c', x)}$  inverse order constraints:
$R(c|x) \leq R(\tilde{c}, x) \Leftrightarrow w_\beta(c, x) \geq w_\beta(\tilde{c}, x)$

Example (unnormalized) weights:
- Boltzmann: $w_\beta(c, X) = \exp(-\beta R(c, x))$
- Fermi: $w_{\beta,\gamma}(c, x) = (1 + \exp(-\beta(R(c, x) - \gamma)))^{-1}$
- Approx: $w_\gamma(c, x) = \begin{cases} 1 & \text{if } R(c, x) \leq R(c^\perp, x) + \gamma \\ 0 & \text{otherwise} \end{cases}$

Normalized Boltzmann weights:
$w_\beta(c, X) = \exp(-\beta \Delta R(c, X))$, $\Delta R(c, X) = R(c, X) - R(c^\perp, X)$
- $\beta = 0 \to$ all weights $w_\beta(c, X) = 1$ ind.-of-costs. $Z_\beta = |c(X)|$ indicates the size of the hypothesis space
- high $\beta \to$ all weights are small compared to $w_\beta(c^\perp, x)$ $Z_\beta$ essentially counts the number of globally opt. solns
- intermediate $\beta$: $Z_\beta$ is the effective number of patterns that approx. fit the dataset $X$, where $\beta$ defines the precision of approx. — noise in measurements $X$ reduces the resolution and thus coarsens hypo. class
- weight sum: measures total weight of hypotheses with low costs, aka partition function when we use Boltzmann weights

Equivariance transformations ↓ idea → shifting posterior
Assume transformation $T \in \tau$  → posterior invariant
$\forall T', T'' \in \tau$, $\|P(c|T' \circ x) - P(c|T'' \circ x)\|_1 = 0$
implies $|\tau| \leq |C|$ for discrete hypo. spaces
$\sum_{T \in \tau} P(c|T \circ x) \in [\frac{|\tau|}{|C|}(1-P), \frac{|\tau|}{|C|}]$  $P$ measures inhomogeneity of transformation
To $c(X) = c(T \circ X) \to$ equivariance of algorithm examples: permutation for graph clustering, rotation for SVD, translation for mean estimation, scaling for linear regression, permutation and scaling for sparse linear regression

---

## Communication scenario

Sender and receiver both receive an instance $X'$ from proper generator, calculate $P(c|x')$ and agree on a set of $M$ randomly drawn transform $T = \{T_1, ..., T_M\}$ with $P(T) = |T|^{-1}$. Here posteriors $P_c(c|T_s \circ x')$ play the role of codewords in Shannon's random coding theory.
Sender selects a transformation $T_s \in T$ as message and sends it to the problem generator. PG generates new instance $X'' \sim p(X)$ and applies $T_s$, which yields $\tilde{X} = T_s \circ X''$. PG sends $\tilde{X}$ to receiver without revealing either $T_s$ or $x'$. So receiver lacks both the knowledge of $T_s$ and suffers from the stochastic variability of $X$. Then receiver calculates the expected posterior $P_c(c|\tilde{X})$ and decodes the message $\hat{s}$:
$$\hat{s} \in \arg\max_{T \in T} E_{c \sim P_c(c|T \circ x')} P_c(c|\tilde{X})$$
we introduce kernel function to be maxed in decoding:
$$k_{T_s}(x', x'') = E_{c \sim P_c(c|\tilde{X})} P_c(c|T_s \circ x')$$
$$= \sum_{c \in C} P_c(c|T_s \circ x') P_c(c|T_s \circ x'') \in [0,1]$$
Posterior agreement kernel for $T_j = T_s$:
$$k(x', x'') = \sum_{c \in C} P_c(c|x') P_c(c|x'')$$ measures the similarity of $x'$ and $x''$ that is induced by the posterior distribution of width $\beta$
- essentially, the posterior specifies a sampling procedure how to choose hypotheses $c$ that are highly likely, given data $V$

$$Z_q = Z(X^{(q)}) = \sum_{c \in c(X^{(q)})} \exp(-\beta R(c, X^{(q)})), \quad q = 1, 2$$
$$Z_{12} = Z(X^{(1)}, X^{(2)}) = \sum_{c \in c(X^{(2)})} \exp(-\beta(\Delta R(c, X^{(1)}) + \Delta R(c, X^{(2)})))$$
$$\hat{s} \in \arg\max_{T \in T} \sum_{c \in c(X^{(1)})} \exp(-\beta(R(c, T \circ x^{(1)}) + R(c, \tilde{X})))$$

vanishing error rate: $P \leq I_\beta(T_s, \hat{s}) = \frac{1}{n} \log(\frac{|\Sigma T_s| Z_{12}}{Z_1 Z_2})$

approx. capacity: $CAP(T_s, \hat{s}) = \max_\beta I_\beta(T_s, \hat{s})$

WASC application: randomly split $X$ into $X^{(1)}$ and $X^{(2)}$, for each candidate cost function $R(c, X)$ ER compute mutual info and max. it wrt $\beta$ then select $R$ that achieves highest capacity at the best resolution $\beta^*$
$$I_\beta(T_s, \hat{s}) = \frac{1}{n}(\log\frac{|\Sigma T_s|}{Z_1} + \log\frac{|c^{(2)}|}{Z_2} - \log\frac{|c^{(2)}|}{Z_{12}})$$
$|\Sigma T_s|$: card. of set of possible transformations
$$P(\hat{s} \neq T_s | T_s) = P(\max_{j \neq s} k_{T_j T_s}(x', x'') \geq k(x', x'') | T_s)$$
$$(\text{union}) \leq \sum_{j \neq s} P(k_{T_j T_s} \cdots)$$
$$= \sum_{j \neq s} E_{x', x''} [P(k_{T_j T_s} \geq k | T_s, x', x'')]$$
$$(\text{Markov}) \leq \sum_{j \neq s} E_{x', x''} \frac{E_{T_j} k_{T_j T_s}(x', x'')}{k(x', x'')}$$