# 3   The least unlikely way maximizes entropy

## 3a   Cramér-Varadhan theorem

Let $\mu$ be a probability measure on $\mathbb{R}^d$ that is truly $d$-dimensional, that is, does not sit on an affine subspace of dimension $d-1$ (or less). Its cumulant generating function $\Lambda_\mu$ is real-analytic on the interior $G \subset \mathbb{R}^d$ of the convex set $\{t \in \mathbb{R}^d : \Lambda_\mu(t) < \infty\}$; we assume that $G \neq \emptyset$. For every $t \in G$ the tilted measure $\mu_t$ has the expectation vector $\nabla \Lambda_\mu(t) = \left(\frac{\partial}{\partial t_j} \Lambda_\mu(t)\right)_j$, and the covariance matrix $\left(\frac{\partial^2}{\partial t_i \partial t_j} \Lambda_\mu(t)\right)_{i,j}$, positive definite. The Legendre transform $\Lambda_\mu^*$ of $\Lambda_\mu$ is

$$\Lambda_\mu^*(x) = \sup_{t \in \mathbb{R}^d} \left(\langle t, x \rangle - \Lambda_\mu(t)\right) \in [0, \infty];$$

and if $x = \nabla \Lambda_\mu(t)$ for some $t \in G$, then

$$\Lambda_\mu^*(x) = \langle t, x \rangle - \Lambda_\mu(t).$$

The mapping $\nabla \Lambda_\mu : G \to \mathbb{R}^d$ is one-to-one (since $\Lambda_\mu$ is strictly convex on $G$), real-analytic, and its differential does not degenerate. It follows that the image $\nabla \Lambda_\mu(G)$ is open, and $\nabla \Lambda_\mu : G \to \nabla \Lambda_\mu(G)$ is a diffeomorphism.

**3a1 Theorem.** For every nonempty bounded open set $U \subset \mathbb{R}^d$ such that $\overline{U} \subset \nabla \Lambda_\mu(G)$,

$$\ln \mu^{*n}(nU) = -n \inf_{x \in U} \Lambda_\mu^*(x) + o(n) \quad \text{as } n \to \infty.$$

(Here $nU = \{nx : x \in U\}$.)

**Proof.** Theorem 2c14, generalized to $\mathbb{R}^d$ (according to 2c18), gives for $x \in U$

(3a2)          $\ln \mu^{*n}(nx + [-\varepsilon_n, \varepsilon_n]^d) = -n \Lambda_\mu^*(x) + o(n) \quad \text{as } n \to \infty,$

provided that $\varepsilon_n/n \to 0$ and $\varepsilon_n/\sqrt{n} \to \infty$. We take such a sequence $(\varepsilon_n)_n$.

First, we prove the lower bound:

$$\liminf_{n\to\infty} \frac{1}{n} \ln \mu^{*n}(nU) \geq -\inf_{x\in U} \Lambda_\mu^*(x) \,,$$

that is, $\liminf(\dots) \geq -\Lambda_\mu^*(x)$ for every $x \in U$. To this end we note that, for $n$ large enough,

$$\frac{1}{n}\ln\mu^{*n}(nU) \geq \frac{1}{n}\ln\mu^{*n}\big(nx+[-\varepsilon_n,\varepsilon_n]^d\big) \geq -\Lambda_\mu^*(x) + o(1)\,.$$

For proving the upper bound,

$$\limsup_{n\to\infty} \frac{1}{n}\ln\mu^{*n}(nU) \leq -\inf_{x\in U}\Lambda_\mu^*(x)\,,$$

we note that the $o(n)$ in (3a2) is uniform in $x \in U$ (since $t$ and $\Lambda_\mu''(t)$ are bounded in $x \in U$; inspect the proof of 2c14). We cover[1] $U$ by $\mathcal{O}\big((n/\varepsilon_n)^d\big)$ cubes of the form $x + [-\frac{\varepsilon_n}{n}, \frac{\varepsilon_n}{n}]^d$ with $x \in U$ and get

$$\ln\mu^{*n}(nU) \leq \ln\mathcal{O}\big((n/\varepsilon_n)^d\big) + \Big(-n\inf_{x\in U}\Lambda_\mu^*(x) + o(n)\Big)\,;$$

the upper bound follows, since $n/\varepsilon_n = o(\sqrt{n})$.     □

Note that the lower bound holds whenever $U$ (rather than its closure) is contained in $\nabla\Lambda_\mu(G)$. For arbitrary open $U \subset \mathbb{R}^n$ the lower bound still holds, but I do not prove it now; this is quite a different story, and rarely needed. Upper bounds for larger sets are more useful.

**3a3 Exercise** (upper bound for a half-space)**.** Let $\Lambda_\mu(t) < \infty$ and $c \geq 0$, then[2]
$$\mu^{*n}\big(\{nx : \langle t,x\rangle - \Lambda_\mu(t) \geq c\}\big) \leq \mathrm{e}^{-nc}\,.$$

Prove it.

**3a4 Lemma** (half-space not containing the expectation)**.** Let $0 \in G$ and $c > \int \langle t,x\rangle\,\mu(\mathrm{d}x)$; then

$$\exists \varepsilon > 0\ \forall n\ \ \mu^{*n}(\{nx : \langle t,x\rangle \geq c\}) \leq \mathrm{e}^{-\varepsilon n}\,.$$

---

[1]For $x \notin U$ note that $U \cap \big(x + [-\frac{\varepsilon_n}{n}, \frac{\varepsilon_n}{n}]^d\big)$ can be covered by no more than $2^d$ cubes of the form $y + [-\frac{\varepsilon_n}{n}, \frac{\varepsilon_n}{n}]^d$ with $y \in U$.

[2]No need to assume that $G \neq \emptyset$.

**Proof.** We note that $c > \langle t, \nabla\Lambda_\mu(0)\rangle = \lim_{\delta\to 0}\frac{\Lambda_\mu(\delta t)}{\delta}$ and take $\delta > 0$ such that $\delta t \in G$ and $\Lambda_\mu(\delta t) < \delta c$. By 3a3,

$$\mu^{*n}(\{nx : \langle t, x\rangle \geq c\}) = \mu^{*n}(\{nx : \langle \delta t, x\rangle - \Lambda_\mu(\delta t) \geq \delta c - \Lambda_\mu(\delta t)\}) \leq \mathrm{e}^{-\varepsilon n}$$

where $\varepsilon = \delta c - \Lambda_\mu(\delta t) > 0$. $\qquad\square$

**3a5 Proposition** (exponential concentration near the expectation)**.** If $0 \in G$, and $U \subset \mathbb{R}^d$ is a neighborhood of the point $\int x\,\mu(\mathrm{d}x)$, then

$$\exists \varepsilon > 0 \; \forall n \;\; \mu^{*n}(nU) \geq 1 - \mathcal{O}(\mathrm{e}^{-\varepsilon n})\,.$$

**Proof.** Lemma 3a4 applied to $t = \pm e_1, \ldots, \pm e_d$ (where $(e_1, \ldots, e_d)$ is the usual basis of $\mathbb{R}^d$) gives

$$\forall \delta > 0 \; \exists \varepsilon > 0 \; \forall n \;\; \mu^{*n}\left(n\Big(\int x\,\mu(\mathrm{d}x) + [-\delta, \delta]^d\Big)\right) \geq 1 - 2d\mathrm{e}^{-\varepsilon n}\,.$$

It remains to take $\delta$ such that $\int x\,\mu(\mathrm{d}x) + [-\delta, \delta]^d \subset U$. $\qquad\square$

**3a6 Proposition.** Assume that $K \subset \mathbb{R}^d$ is a compact set, and the restriction $\Lambda_\mu^*|_K$ is continuous. Then

$$\ln \mu^{*n}(\{nx : x \in K, \Lambda_\mu^*(x) \geq c\}) \leq -nc + o(n) \quad \text{as } n \to \infty$$

for every $c \in (0, \infty)$.

**Proof.** We'll prove that

$$\limsup_{n\to\infty} \frac{1}{n} \ln \mu^{*n}(\{nx : x \in K, \Lambda_\mu^*(x) \geq c\}) \leq -(c - \varepsilon)$$

for arbitrary $\varepsilon > 0$.

For $t \in \mathbb{R}^d$ such that $\Lambda_\mu(t) < \infty$ the half-space

$$H_t = \{x \in \mathbb{R}^d : \langle t, x\rangle - \Lambda_\mu(t) > c - \varepsilon\}$$

satisfies

$$\mu^{*n}(nH_t) \leq \exp\big(-n(c - \varepsilon)\big)$$

by 3a3. It is sufficient to cover the set $\{x : x \in K, \Lambda_\mu^*(x) \geq c\}$ by finitely many half-spaces $H_t$. The union of all $H_t$ is an open covering of this compact set (recall the definition of $\Lambda_\mu^*$); it has a finite subcovering. $\qquad\square$

**3a7 Exercise** (minimum of $\Lambda^*$ on half-space)**.** If $t \in G$ and $x = \nabla\Lambda_\mu(t)$, then $\Lambda^*(y) - \Lambda^*(x) \geq \langle t, y\rangle - \langle t, x\rangle$ for all $y$; and therefore

$$\min_{y:\langle t,y\rangle \geq \langle t,x\rangle} \Lambda^*(y) = \Lambda^*(x)\,.$$

Prove it.

## 3b   Relative entropy

Let $\mu, \nu$ be two probability measures on a measurable space $\Omega$, such that

$$\nu = \mu_u \quad \text{for some measurable function } u : \Omega \to \mathbb{R} \,.$$

By the Radon-Nikodym theorem, this happens if and only if $\mu$ and $\nu$ are mutually absolutely continuous (in other words, equivalent), that is, have the same sets of measure zero.

By 2a6, $\mu = \nu_{-u}$, and $\Lambda_\nu(-u) = -\Lambda_\mu(u)$. Note that, for arbitrary $c \in \mathbb{R}$, $\mu_{u+c} = \mu_u$ and $\Lambda_\mu(u+c) = \Lambda_\mu(u) + c$. We may replace the given $u$ with $u - \Lambda_\mu(u)$, getting

$$\nu = \mu_u \,, \quad \mu = \nu_{-u} \,, \quad \Lambda_\mu(u) = 0 = \Lambda_\nu(-u) \,, \quad u = \ln \frac{\mathrm{d}\nu}{\mathrm{d}\mu} = -\ln \frac{\mathrm{d}\mu}{\mathrm{d}\nu} \,.$$

This is closely related to the so-called Neyman-Pearson (statistical) test:[1] given a sample $(\omega_1, \ldots, \omega_n) \in \Omega^n$ from a partially known distribution (either $\mu$ or $\nu$), a statistician needs to decide, is it sampled from $\mu$ or $\nu$. It is well-known (and easy to see) that the optimal test is based on the so-called observed log-likelihood ratio

$$S_n = \ln \frac{\mathrm{d}\nu^n}{\mathrm{d}\mu^n}(\omega_1, \ldots, \omega_n) = \ln \frac{\mathrm{d}\nu}{\mathrm{d}\mu}(\omega_1) + \cdots + \ln \frac{\mathrm{d}\nu}{\mathrm{d}\mu}(\omega_n) \,.$$

The statistician compares $S_n/n$ with a chosen threshold $\gamma_n$; if $S_n/n \geq \gamma_n$, the decision is $\nu$, and if[2] $S_n/n < \gamma_n$, the decision is $\mu$. Of course, the decision can be erroneous. The decision $\nu$ is wrong with probability $\beta_n = \mu^n\{S_n/n \geq \gamma_n\}$;[3] the decision $\mu$ is wrong with probability $\alpha_n = \nu^n\{S_n/n < \gamma_n\}$. If $\nu$ is the null hypothesis and $\mu$ is the alternative hypothesis, then the wrong decision $\mu$ is called the type I error ("false positive", "false alarm"), and the wrong decision $\nu$ is called the type II error ("false negative", "miss"). Often, one chooses a small $\alpha$ (the size of the test) and takes the greatest $\gamma_n$ such that $\alpha_n \leq \alpha$.[4] The corresponding $1 - \beta_n$ is called the power of the test.

The large deviations theory can help in estimating $\beta_n$ for large $n$. To this end, tilting on $\Omega$ may be replaced with tilting on $\mathbb{R}$ (recall Sect. 2a); that is, we replace $\Omega$ with $\mathbb{R}$, $\mu$ with $u_*(\mu)$ (renamed to $\mu$), and $\nu$ with $u_*(\nu)$, getting
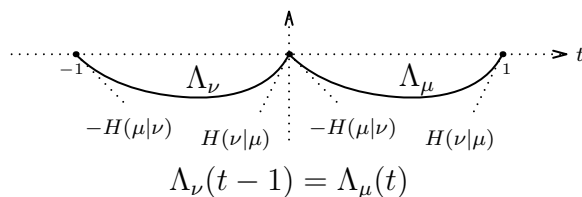
$$\nu = \mu_1 \,, \quad \mu = \nu_{-1} \,, \quad \Lambda_\mu(1) = 0 = \Lambda_\nu(-1) \,, \quad \frac{\mathrm{d}\nu}{\mathrm{d}\mu}(x) = \mathrm{e}^x \,;$$

---

[1]See Sect. 3.4 in Dembo and Zeitouni, and "Neyman-Pearson lemma" in Wikipedia.
[2]When $S_n/n = \gamma_n$, the decision may be randomized. Never mind.
[3]That is, $\mu^n\big(\{(\omega_1, \ldots, \omega_n) \in \Omega^n : S_n(\omega_1, \ldots, \omega_n)/n \geq \gamma_n\}\big)$, of course.
[4]This is the Neyman-Pearson approach. There is also Bayesian approach.

$$\Lambda_\nu(t-1) = \Lambda_\mu(t)$$

$S_n$ is distributed either $\mu^{*n}$ or $\nu^{*n}$, thus,[1]

$$\alpha_n = \nu^{*n}(-\infty, n\gamma_n)\,, \quad \beta_n = \mu^{*n}[n\gamma_n, \infty)\,.$$

**3b1 Exercise.** If $\mu \neq \nu$ then $\int x\,\mu(\mathrm{d}x) = \Lambda_\mu'(0) \in [-\infty, 0)$ and $\int x\,\nu(\mathrm{d}x) = \Lambda_\nu'(0) \in (0, \infty]$.

    Prove it. Show by examples that all cases (finite or infinite) are possible.

    We have

$$\int x\,\nu(\mathrm{d}x) = \int \ln \frac{\mathrm{d}\nu}{\mathrm{d}\mu}\,\mathrm{d}\nu = H(\nu|\mu) \in (0, \infty]\,;$$

$$\int x\,\mu(\mathrm{d}x) = \int \ln \frac{\mathrm{d}\nu}{\mathrm{d}\mu}\,\mathrm{d}\mu = -\int \ln \frac{\mathrm{d}\mu}{\mathrm{d}\nu}\,\mathrm{d}\mu = -H(\mu|\nu) \in [-\infty, 0)\,;$$

$H(\nu|\mu)$ is the well-known *relative entropy* of $\nu$ relative to $\mu$, called also the Kullback-Leibler divergence[2] $D(\eta\|\mu)$. Note that the transition from $\Omega$ to $\mathbb{R}$ does not change the relative entropy.

    By a well-known result of Chernoff,[3] $\beta_n$ is exponentially small, as follows.

**3b2 Proposition.** If $\frac{1}{n}\ln \alpha_n \to 0$ and $\limsup_n \alpha_n < 1$, then $\frac{1}{n}\ln \beta_n \to -H(\nu|\mu)$.

    That is, a large part of the measure $\nu^n$ sits on a set whose measure $\mu^n$ is $\exp(-nH(\eta|\mu) + o(n))$ (and cannot be smaller).

    In particular, when $\mu$ is the uniform distribution on a finite $\Omega$, we have

(3b3) $$H(\nu|\mu) = \ln(\#\Omega) - H(\nu)\,;$$

---

[1] I write just $\nu^{*n}(-\infty, n\gamma_n)$ instead of $\nu^{*n}\big((-\infty, n\gamma_n)\big)$, of course.

[2] See "Kullback-Leibler divergence" in Wikipedia.

[3] Strangely, it is well-known under the name "Chernoff-Stein lemma" or even "Stein's lemma"! Here is a quote from page 85 of "Statistical signal processing" by Don Johnson: "The attribution to statistician Charles Stein is probably incorrect. Herman Chernoff wrote a paper providing a derivation of this result. A reviewer stated that he thought Stein had derived the result in a technical report, which Chernoff had not seen. Chernoff modified his paper to include the reference without checking it. Chernoff's paper provided the link to Stein. However, Stein later denied he had proved the result; so much for not questioning reviewers! Stein's Lemma should be known as Chernoff's Lemma."

here[1]

$$H(\nu) = -\sum_\omega \nu(\omega) \ln \nu(\omega)$$

is the (Shannon) entropy of $\nu$. Thus, a large part of the measure $\nu^n$ sits on a set of $\exp\big(nH(\nu) + o(n)\big)$ points (and cannot be smaller).

**Proof.** First we prove that $\liminf_n \gamma_n \geq H(\nu|\mu)$ (be it finite or not). Given $a < H(\nu|\mu)$, we'll prove that $\gamma_n > a$ for all $n$ large enough. We cannot use Lemma 3a4, since $\Lambda_\nu(0+)$ may be $\infty$, but we can adapt that proof as follows.

We note that $\Lambda'_\nu(0-) = H(\nu|\mu) > a$, and take $\delta > 0$ such that $\Lambda_\nu(-\delta) < -\delta a$. By 3a3, for $\varepsilon = -\delta a - \Lambda_\nu(-\delta) > 0$ (and $t = -\delta$),

$$\mathrm{e}^{-\varepsilon n} \geq \nu^{*n}\big(\{nx : -\delta x - \Lambda_\nu(-\delta) \geq \varepsilon\}\big) =$$
$$= \nu^{*n}\big(\{nx : -\delta x \geq -\delta a\}\big) = \nu^{*n}(-\infty, na]\,,$$

thus, $\frac{1}{n} \ln \nu^{*n}(-\infty, na] \leq -\varepsilon < \frac{1}{n} \ln \alpha_n = \frac{1}{n} \ln \nu^{*n}(-\infty, n\gamma_n)$ for all $n$ large enough, therefore $na < n\gamma_n$, that is, $\gamma_n > a$, which proves that $\liminf_n \gamma_n \geq H(\nu|\mu)$.

Second, we prove the upper bound: $\limsup_n \frac{1}{n} \ln \beta_n \leq -H(\nu|\mu)$ (finite or not). Given $a < H(\nu|\mu)$, we'll prove that $\frac{1}{n} \ln \beta_n \leq -a$, that is, $\beta_n \leq \mathrm{e}^{-an}$, for all $n$ large enough. We have $\gamma_n > a$; by 3a3 (for $t = 1$),

$$\beta_n = \mu^{*n}[n\gamma_n, \infty) \leq \mu^{*n}[na, \infty) = \mu^{*n}\big(\{nx : x - \Lambda_\mu(1) \geq a\}\big) \leq \mathrm{e}^{-an}\,.$$

Third, assuming $H(\nu|\mu) < \infty$ (otherwise all is done already), we prove the lower bound: $\liminf_n \frac{1}{n} \ln \beta_n \geq -H(\nu|\mu)$. Given $a > H(\nu|\mu)$, we'll prove that $\frac{1}{n} \ln \beta_n \geq -a + o(1)$ for $n \to \infty$. We have $a > \int x\,\nu(\mathrm{d}x)$; by the weak law of large numbers, $\nu^{*n}[an, \infty) \to 0$; therefore

$$\beta_n = \mu^{*n}[n\gamma_n, \infty) \geq \mu^{*n}[n\gamma_n, an] \geq$$
$$\geq M_{\mu^{*n}}(1)\mathrm{e}^{-an}\mu_1^{*n}[n\gamma_n, an] = \mathrm{e}^{-an}\nu^{*n}[n\gamma_n, an]\,;$$
$$\liminf_n \nu^{*n}[n\gamma_n, an] = \liminf_n \big(\nu^{*n}[n\gamma_n, \infty) - \nu^{*n}(an, \infty)\big) \geq$$
$$\geq 1 - \limsup_n \alpha_n > 0\,;$$

$\frac{1}{n} \ln \nu^{*n}[n\gamma_n, an] \to 0$; $\frac{1}{n} \ln \beta_n \geq -a + o(1)$.                                    $\square$

---

[1]Strangely, the relative entropy and the entropy differ by the sign; I do not know why.

## 3c   Sanov's theorem

We return to the empirical distribution $(\nu_{-1}, \nu_0, \nu_{+1})$ mentioned in Sect. 2a. More generally, now we consider $d$ possible values of a random variable, with probabilities $p_1, \ldots, p_d > 0$, $p_1 + \cdots + d_d = 1$. The frequencies $\eta_1, \ldots, \eta_d$ satisfy $\eta_1 + \cdots + \eta_d = 1$ and are distributed multinomially,[1]

$$\mathbb{P}\left(\eta_1 = \frac{a_1}{n}, \ldots, \eta_d = \frac{a_d}{n}\right) = \frac{n!}{a_1! \ldots a_d!} p_1^{a_1} \ldots p_d^{a_d},$$

and we could use the Stirling formula (as in Sect. 1a). But instead, we introduce such a probability measure $\mu$ on $\mathbb{R}^d$:

$$\mu\left(\left\{e_j - \frac{e_1 + \cdots + e_d}{d}\right\}\right) = p_j \quad \text{for } j = 1, \ldots, d,$$

where $(e_1, \ldots, e_d)$ is the usual basis of $\mathbb{R}^d$. That is, we use the vectors $e_j - \frac{e_1 + \cdots + e_d}{d}$ as the possible values of the random variable. These vectors span a $(d-1)$-dimensional vector subspace (a hyperplane) $E = \{(x_1, \ldots, x_d) : x_1 + \cdots + x_d = 0\}$ in $\mathbb{R}^d$, and we may apply the Cramer-Varadhan theorem within this subspace.[2] Indeed, the measure $\mu^{*n}$ is the distribution of $\left(n\eta_1 - \frac{n}{d}, \ldots, n\eta_d - \frac{n}{d}\right)$.

The convex hull of the vectors $e_1, \ldots, e_d$ is the simplex $\{(x_1, \ldots, x_d) \in [0, \infty)^d : x_1 + \cdots + x_d = 1\}$, thus, the convex hull of the vectors $e_j - \frac{e_1 + \cdots + e_d}{d}$ is the simplex $\{(x_1 - \frac{1}{d}, \ldots, x_d - \frac{1}{d}) : x_1, \ldots, x_d \geq 0, x_1 + \cdots + x_d = 1\} = E \cap [-\frac{1}{d}, \infty)^d$. The expectation of $\mu$,

$$\int x \, \mu(\mathrm{d}x) = p_1\left(e_1 - \frac{e_1 + \cdots + e_d}{d}\right) + \cdots + p_d\left(e_d - \frac{e_1 + \cdots + e_d}{d}\right) =$$
$$= \left(p_1 - \frac{1}{d}\right)e_1 + \cdots + \left(p_d - \frac{1}{d}\right)e_d \in E \cap \left(-\frac{1}{d}, \infty\right)^d$$

is an interior point of the simplex. And every interior point of the simplex corresponds to some $p_1, \ldots, p_d$, thus, to some $\mu$.

A tilted measure $\mu_t$ also sits on these $d$ vertices of the simplex, with probabilities $q_j = \frac{1}{M_\mu(t)} e^{t_j} p_j$ $(j = 1, \ldots, d)$. Interestingly, when $t$ runs over $E$, the tilted measure runs over *all* strictly positive probability measures on these vertices. Indeed, given $q_1, \ldots, q_d > 0$, $q_1 + \cdots + q_d = 1$, we take

$$s_j = \ln q_j - \ln p_j, \quad t_j = s_j - \frac{s_1 + \cdots + s_d}{d},$$

---

[1]The binomial distribution, treated in Sect. 1a, is a special case.

[2]You may choose an orthonormal basis in $E$, thus turning $E$ into a copy of $\mathbb{R}^{d-1}$.

then $\mu_t$ corresponds to $q_1, \ldots, q_d$ (since $e^{t_j} p_j = \text{const} \cdot q_j$).

We have $M_\mu : E \to (0, \infty)$,

$$M_\mu(t) = \int e^{\langle t, x \rangle} \mu(\mathrm{d}x) = p_1 e^{t_1} + \cdots + p_d e^{t_d};$$

$\nabla \Lambda_\mu : E \to E$,

$$\nabla \Lambda_\mu(t) = \frac{1}{M_\mu(t)} \nabla M_\mu(t) = \Big( \frac{p_1 e^{t_1}}{M_\mu(t)} - \frac{1}{d}, \ldots, \frac{p_d e^{t_d}}{M_\mu(t)} - \frac{1}{d} \Big).$$

We know that $\nabla \Lambda_\mu(t) = \int x \, \mu_t(\mathrm{d}x)$ runs over the interior of the simplex (when $t$ runs over $E$). Thus,

$$\nabla \Lambda_\mu(G) = \nabla \Lambda_\mu(E) = E \cap \Big( -\frac{1}{d}, \infty \Big)^d.$$

The Legendre transform $\Lambda_\mu^*(x)$ for $x \in E \cap \left( -\frac{1}{d}, \infty \right)^d$ is $\langle t, x \rangle - \Lambda_\mu(t)$ when $\nabla \Lambda_\mu(t) = x$; but let us treat it as a function of the tilted measure $\mu_t$, that is, of its probabilities $q_1, \ldots, q_d$ (rather than $x$). We have $x_j = q_j - \frac{1}{d}$,

$$q_j = \frac{1}{M_\mu(t)} e^{t_j} p_j; \quad \ln \frac{q_j}{p_j} = t_j - \Lambda_\mu(t);$$

$$\langle t, x \rangle = \sum_j t_j \Big( q_j - \frac{1}{d} \Big) = \sum_j t_j q_j = \sum_j q_j \ln \frac{q_j}{p_j} + \Lambda_\mu(t); \quad \langle t, x \rangle - \Lambda_\mu(t) = \sum_j q_j \ln \frac{q_j}{p_j};$$

thus,

$$(3c1) \qquad \Lambda_\mu^*(x) = q_1 \ln \frac{q_1}{p_1} + \cdots + q_d \ln \frac{q_d}{p_d} = \int \ln \frac{\mathrm{d}\mu_t}{\mathrm{d}\mu} \, \mathrm{d}\mu_t = H(\mu_t | \mu);$$

the relative entropy again!

Ignoring the distinction between a probability measure $\nu$ on the vertices $e_j - \frac{e_1 + \cdots + e_d}{d}$ of the (closed) simplex $K = E \cap [-\frac{1}{d}, \infty)^d$ and the point $\int x \, \nu(\mathrm{d}x) \in K$ of this simplex, we write

$$\Lambda_\mu^*(\nu) = H(\nu | \mu) \quad \text{for all } \mu \in K.$$

This equality holds on the interior of $K$, and therefore on the whole $K$, since $H(\cdot | \mu)$ is continuous on $K$ ($0 \ln 0 = 0$, of course), and $\Lambda_\mu^*(\cdot)$ is convex and lower semicontinuous (being the supremum of affine functions), therefore continuous on $K$, too. Prop. 3a6 gives an upper bound for the empirical distribution $\eta$:

$$(3c2) \qquad \ln \mu^{*n}(\{n\eta : H(\eta | \mu) \geq c\}) \leq -nc + o(n) \quad \text{as } n \to \infty$$

for every $c \in (0, \infty)$. On the other hand, Theorem 3a1 gives a lower bound:

$$(3c3) \qquad \ln \mu^{*n}(nU) \geq -n \inf_{\nu \in U \cap K} H(\nu|\mu) + o(n) \quad \text{as } n \to \infty$$

for every open set $U$.

**3c4 Exercise** (minimum of relative entropy on half-space)**.** Let $\mu$ be a strictly positive probability measure on a finite set $\Omega$, and $u : \Omega \to \mathbb{R}$. Then $H(\nu|\mu) - H(\mu_u|\mu) \geq \int u \, \mathrm{d}\nu - \int u \, \mathrm{d}\mu_u$ for all probability measures $\nu$ on $\Omega$; and therefore

$$\min_{\nu : \int u \, \mathrm{d}\nu \geq \int u \, \mathrm{d}\mu_u} H(\nu|\mu) = H(\mu_u|\mu) \, .$$

Prove it.[1]

## 3d   Conditioning and tilting

Recall the essential infimum and supremum of a measurable function $u$ on a probability space $(\Omega, \mu)$:

$$\operatorname{ess\,inf} u = \sup\{x : \mu(u^{-1}(-\infty, x)) = 0\} = \inf\{x : \mu(u^{-1}(-\infty, x)) > 0\} \in [-\infty, +\infty) \, ;$$
$$\operatorname{ess\,sup} u = \inf\{x : \mu(u^{-1}(x, \infty)) = 0\} = \sup\{x : \mu(u^{-1}(x, \infty)) > 0\} \in (-\infty, +\infty]$$

(here $\sup \emptyset = -\infty$, $\inf \emptyset = +\infty$).

**3d1 Exercise.** (a) $\lim_{t \to -\infty} \int u \, \mathrm{d}\mu_{tu} = \operatorname{ess\,inf} u$, $\lim_{t \to +\infty} \int u \, \mathrm{d}\mu_{tu} = \operatorname{ess\,sup} u$.
   (b) If $\operatorname{ess\,inf} u < \operatorname{ess\,sup} u$, then the function $t \mapsto \int u \, \mathrm{d}\mu_{tu}$ is strictly increasing.
   Prove it.[2]

Now we turn to measures and functions on a finite set; let it be $\{1, \ldots, d\}$.
   For every $\omega \in \{1, \ldots, d\}^n$ we define frequencies $\eta_{n,\omega}(x) = \frac{1}{n} \cdot \#\{k : \omega_k = x\}$ for $x \in \{1, \ldots, d\}$; these frequencies are a probability measure $\eta_{n,\omega}$ on $\{1, \ldots, d\}$. Thus, $\eta_n$ maps $\{1, \ldots, d\}^n$ to the set of probability measures on $\{1, \ldots, d\}$.
   Let $\mu$ be a strictly positive probability measure on $\{1, \ldots, d\}$; we endow $\{1, \ldots, d\}^n$ with the probability measure $\mu^n$, and treat $\eta_n$ as a *random probability measure* on $\{1, \ldots, d\}$.
   Let $u : \{1, \ldots, d\} \to \mathbb{R}$, and $c \in \mathbb{R}$, $\int u \, \mathrm{d}\mu < c < \max(u(1), \ldots, u(d))$. We consider events $E_n = \{\int u \, \mathrm{d}\eta_n \geq c\}$; that is, $E_n \subset \{1, \ldots, d\}^n$, $E_n = \{\omega : \int u \, \mathrm{d}\eta_{n,\omega} \geq c\}$.

---

[1]Hint: use 3a7.
   [2]Hint: the general case reduces to a measure on $\mathbb{R}$ and the function $u(x) = x$, as noted in Sect. 2a.

**3d2 Proposition.** Conditionally, given $E_n$, the random measure $\eta_n$ converges (as $n \to \infty$) in probability to the tilted measure $\mu_{tu}$ where $t > 0$ is such that $\int u \, d\mu_{tu} = c$.

Why just this measure, $\mu_{tu}$? Because (as we'll see soon) this is the unique minimizer of the relative entropy $H(\nu|\mu)$ among all measures $\nu$ such that $\int u \, d\nu \geq c$. In particular, when $\mu$ is the uniform measure, $\mu_{tu}$ is the unique maximizer of the Shannon entropy $H(\nu)$ (recall (3b3)).[1]

**Proof.** In the compact simplex $K$ of all probability measures on $\{1, \ldots, d\}$ we consider the closed set $F = \{\nu : \int u \, d\nu \geq c\}$, the open set $G = \{\nu : \int u \, d\nu > c\}$, note that $F = \overline{G}$, and $\mu_{tu} \in F \setminus G$. Let $U$ be a neighborhood of $\mu_{tu}$. The lower bound (3c3) gives

$$\ln \mu^n \{\eta_n \in G \cap U\} \geq -n \inf_{\nu \in G \cap U} H(\nu|\mu) + o(n) \geq -n H(\mu_{tu}|\mu) + o(n) \,,$$

since $\mu_{tu}$ belongs to the closure of $G \cap U$. The upper bound (3c2) gives

$$\ln \mu^n (\{H(\eta_n|\mu) > H(\mu_{tu}|\mu) + \varepsilon\}) \leq -n(H(\mu_{tu}|\mu) + \varepsilon) + o(n)$$

for all $\varepsilon > 0$. Thus, the conditional probability of the event $H(\eta_n|\mu) > H(\mu_{tu}|\mu) + \varepsilon$ given $E_n$ (that is, $\eta_n \in F$) tends to 0. We'll prove that

$$F \setminus U \subset \{\nu : H(\nu|\mu) > H(\mu_{tu}|\mu) + \varepsilon\}$$

for some $\varepsilon > 0$. By compactness, it is sufficient to prove that $H(\nu|\mu) > H(\mu_{tu}|\mu)$ for all $\nu \in F \setminus \{\mu_{tu}\}$. The weaker inequality $H(\nu|\mu) \geq H(\mu_{tu}|\mu)$ follows from 3c4. It remains to note that the function $H(\cdot|\mu)$ on $K$ is strictly convex, since the function $x \mapsto x \ln x$ on $[0,1]$ is strictly convex. $\square$

Of course, the set $\{1, \ldots, d\}$ may be replaced with any finite set.

**3d3 Example.** Continuing an example of Sect. 2a, we take $\mu$ uniform on $\{-1, 0, 1\}$, $u(x) = x$, and $c = 3/7$. Then $\mu_{tu}$ is $\left(\frac{1}{a^2+a+1}, \frac{a}{a^2+a+1}, \frac{a^2}{a^2+a+1}\right)$ where $a = e^t$; and $\int u \, d\mu_{tu} = \frac{a^2-1}{a^2+a+1}$ is $3/7$ for $a = 2$.

Clearly, the conditioning on $\eta_n \in F$ may be replaced with conditioning on $\eta_n \in A$ whenever $A \subset F$ and $A \supset F \cap U$ for some neighborhood $U$ of $\mu_{tu}$. In particular, we may condition on $c \leq \int u \, d\nu \leq c + \varepsilon$ for arbitrary $\varepsilon > 0$.

Treating $\omega = (x_1, \ldots, x_n) \in \{1, \ldots, d\}^n$ as a configuration of a physical system of $n$ particles (as in Sections 1b and 2a), and $n \int u \, d\eta_{n,\omega} = u(x_1) +$

---

[1] "MaxEnt" (maximization of entropy) is used widely; just see "Maximum entropy" in Wikipedia. Regretfully, the strange sign of the relative entropy leads to "MinRelEnt" (??)

$\cdots + u(x_n)$ as the energy[1] $H(\omega)$ of the configuration, and taking $\mu$ to be the uniform measure, we observe that most of configurations $\omega$ satisfying

$$c \leq \frac{1}{n} H(\omega) \leq c + \varepsilon \quad \text{(that is, } c \leq \int u \, \mathrm{d}\eta_{n,\omega} \leq c + \varepsilon)$$

satisfy

$$\eta_{n,\omega} \approx \mu_{tu} \,,$$

and $\mu_{tu}$ maximizes the (Shannon) entropy among all measures $\nu$ such that $c \leq \int u \, \mathrm{d}\nu \leq c + \varepsilon$.

In statistical physics, usually, a system performs[2] a kind of random walk on the set $\{\omega : H(\omega) \approx c\}$. It spends most of the time in the subset $\{\omega : \eta_{n,\omega} \approx \mu_{tu}\}$; and no wonder!

# Index

---

[1]Unfortunately, $H$ is widely used for both the Hamiltonian and the (Shannon) entropy. In physics, entropy is usually denoted $S$.

[2]Since $H$ is the Hamiltonian of an isolated system, but there is also a weak interaction with environment.