

2 Cramér's theorem

2a	A very important notion	6
2b	Surprisingly useful generating functions	10
2c	Independent summands	13

2a A very important notion

You should be puzzled: what about the name of the notion? The answer is that this notion was introduced several times, under different names:

- * “tilted measure” (in the theory of large deviations);¹
- * “canonical ensemble” (in statistical physics);
- * “exponential family” (in statistics);
- * “Esscher transform” (mostly, in financial mathematics and actuarial science).

Here is the simplest nontrivial example. Consider n independent copies X_1, \dots, X_n of a random variable X that takes three values $-1, 0, +1$ with equal probabilities ($1/3$). The frequencies $\nu_{-1} = \frac{1}{n} \cdot \#\{k : X_k = -1\}$, $\nu_0 = \frac{1}{n} \cdot \#\{k : X_k = 0\}$, $\nu_{+1} = \frac{1}{n} \cdot \#\{k : X_k = +1\}$ are random; together they are the so-called empirical distribution $(\nu_{-1}, \nu_0, \nu_{+1})$; and the sample mean $\frac{1}{n}(X_1 + \dots + X_n) = \nu_{+1} - \nu_{-1}$ is also random.

For large n the event $E = \{\frac{1}{n}(X_1 + \dots + X_n) \geq \frac{3}{7}\}$ is of exponentially small probability, and nevertheless, let us consider the conditional distribution of $(\nu_{-1}, \nu_0, \nu_{+1})$ given E . We'll see that, given E ,

$$(\nu_{-1}, \nu_0, \nu_{+1}) \rightarrow \left(\frac{1}{7}, \frac{2}{7}, \frac{4}{7}\right) \quad \text{as } n \rightarrow \infty$$

in probability; that is, for every $\varepsilon > 0$,

$$\mathbb{P}\left(|\nu_{-1} - \frac{1}{7}| \leq \varepsilon, |\nu_0 - \frac{2}{7}| \leq \varepsilon, |\nu_{+1} - \frac{4}{7}| \leq \varepsilon\right) \rightarrow 1 \quad \text{as } n \rightarrow \infty.$$

A wonder, isn't it?

As you may guess, more generally, for arbitrary $a \in (1, \infty)$ the condition $E_a = \{\frac{1}{n}(X_1 + \dots + X_n) \geq \frac{a^2-1}{a^2+a+1}\}$ ² leads in the limit to

$$(\nu_{-1}, \nu_0, \nu_{+1}) = \left(\frac{1}{a^2 + a + 1}, \frac{a}{a^2 + a + 1}, \frac{a^2}{a^2 + a + 1}\right),$$

¹Also, “Cramér-transform” (Hollander, p. 7).

²Thus, $E = E_2$.

that is, to

$$\frac{1}{n}(X_1 + \cdots + X_n) = \nu_{+1} - \nu_{-1} = \frac{a^2 - 1}{a^2 + a + 1} \quad \text{and} \quad \frac{\nu_{+1}}{\nu_0} = \frac{\nu_0}{\nu_{-1}}.$$

It is easy to realize that any violation of the equality $\nu_{+1} - \nu_{-1} = \frac{a^2-1}{a^2+a+1}$ leads to an event $\frac{1}{n}(X_1 + \cdots + X_n) \geq \frac{a^2-1}{a^2+a+1} + \varepsilon$ of probability exponentially smaller than that of E_a . It is less evident that any violation of the equality $\frac{\nu_{+1}}{\nu_0} = \frac{\nu_0}{\nu_{-1}}$ leads also to exponentially smaller probability. But it does, as we'll see.

This fact illustrates a key principle in large deviation theory:

ANY LARGE DEVIATION IS DONE IN THE LEAST UNLIKELY
OF ALL THE UNLIKELY WAYS!

(Quoted from: *Hollander, p. 10.*)

For $a \in (0, 1)$ the same holds under the condition $\frac{1}{n}(X_1 + \cdots + X_n) \leq \frac{a^2-1}{a^2+a+1}$ ($= -\frac{b^2-1}{b^2+b+1}$ for $b = 1/a \in (0, \infty)$).

It is hardly possible to observe in practice the convergence $(\nu_{-1}, \nu_0, \nu_{+1}) \rightarrow (\frac{1}{7}, \frac{2}{7}, \frac{4}{7})$, since for large n it is not feasible to see condition E satisfied even once in a long run.

Now consider a large system of n so-called spin-1 particles, described by the configuration space $\{-1, 0, 1\}^n$. The average spin $\frac{1}{n}(X_1 + \cdots + X_n)$ has practically no chance to reach $3/7$ spontaneously, but can be forced by an external magnetic field. If a measurement shows that the average spin is (close to) $3/7$, then¹ a physicist knows that $(\nu_{-1}, \nu_0, \nu_{+1})$ is (close to) $(\frac{1}{7}, \frac{2}{7}, \frac{4}{7})$; and in particular, $\frac{1}{n}(X_1^2 + \cdots + X_n^2)$ is (close to) $5/7$.

The transition from the distribution $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ on the three-point set $\{-1, 0, 1\}$ to the distribution $(\frac{1}{a^2+a+1}, \frac{a}{a^2+a+1}, \frac{a^2}{a^2+a+1})$ on the same set is a simple example of tilting (called also² twisting, or exponential change of measure, etc).

2a1 Definition. (a) Let μ be a probability measure on \mathbb{R} . For every $t \in \mathbb{R}$ such that $\int e^{tx} \mu(dx) = M_\mu(t) < \infty$ we define the *tilted measure* μ_t by

$$\frac{d\mu_t}{d\mu}(x) = \frac{1}{M_\mu(t)} e^{tx};$$

that is,

$$\int f(x) \mu_t(dx) = \frac{1}{M_\mu(t)} \int f(x) e^{tx} \mu(dx)$$

¹Assuming thermal equilibrium in the external field.

²Buckle, p. 13.

for all bounded continuous functions $f : \mathbb{R} \rightarrow \mathbb{R}$ (and therefore also for all bounded μ -measurable f). Also, the function $M_\mu : \mathbb{R} \rightarrow (0, \infty]$, $M_\mu(t) = \int e^{tx} \mu(dx)$, is called the *moment generating function* (MGF) of μ ; and its logarithm $\Lambda_\mu : \mathbb{R} \rightarrow (-\infty, +\infty]$, $\Lambda_\mu(t) = \ln M_\mu(t)$, is called the *cumulant generating function*.¹

(b) More generally, let μ be a probability measure on a measurable space. For every measurable function² u on this space, satisfying $\int e^u d\mu = M_\mu(u) < \infty$, we define the *tilted measure* μ_u by

$$\frac{d\mu_u}{d\mu}(\cdot) = \frac{1}{M_\mu(u)} e^{u(\cdot)}.$$

Also, M is called the *moment generating functional*, and $\Lambda = \ln M$ is called the *cumulant generating functional*.

Note that the tilted measure is a probability measure.

2a2 Example (Standard normal distribution). $\frac{\mu(dx)}{dx} = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$; $M_\mu(t) = e^{t^2/2}$ (check it); $\Lambda_\mu(t) = t^2/2$; μ_t is just μ shifted by t , that is, $\int f(x-t) \mu_t(dx) = \int f(x) \mu(dx)$ (check it).

2a3 Example (Fair coin). $\mu(\{-1\}) = 1/2 = \mu(\{+1\})$; $M_\mu(t) = \cosh t$; $\mu_t(\{-1\}) = \frac{e^{-t}}{e^t + e^{-t}}$, $\mu_t(\{+1\}) = \frac{e^t}{e^t + e^{-t}}$ (unfair coin).

2a4 Example (Exponential distribution). $\frac{\mu(dx)}{dx} = e^{-x}$ for $x > 0$; $M_\mu(t) = \frac{1}{1-t}$ for $-\infty < t < 1$, otherwise $+\infty$ (check it); μ_t is homothetic to μ , that is, $\int f((1-t)x) \mu_t(dx) = \int f(x) \mu(dx)$ for $-\infty < t < 1$ (check it).

2a5 Example (Discontinuous generating function). $\frac{\mu(dx)}{dx} = \frac{1}{2} \exp(-\sqrt{x})$ for $x > 0$; $M_\mu(t) \leq 1$ for $-\infty < t \leq 0$, but $M_\mu(t) = +\infty$ for $t > 0$, even though $\int x^k \mu(dx) < \infty$ for all k .

2a6 Exercise.

- (a) If $M_\mu(s) < \infty$ then $\forall t \quad M_\mu(s)M_{\mu_s}(t) = M_\mu(s+t)$;
- (b) if both are finite then $(\mu_s)_t = \mu_{s+t}$;
- (c) if $M_\mu(u) < \infty$ then $\forall v \quad M_\mu(u)M_{\mu_u}(v) = M_\mu(u+v)$;
- (d) if both are finite then $(\mu_u)_v = \mu_{u+v}$.

Prove it.

¹Or the logarithmic MGF; also denoted by K_μ .

²Real-valued.

In statistical physics (as was noted in Sect. 1b) probabilities are proportional to $\exp(-\beta H(\cdot))$, where $H(\cdot)$ is the energy, and β the inverse temperature.¹ Thus, if we add a function $h(\cdot)$ to the energy $H(\cdot)$ (which is a usual description of an external field, or another influence) then probabilities are multiplied by $\exp(-\beta h(\cdot))$ and a normalizing constant. It means that the initial probability measure μ is replaced with the tilted measure $\mu_{-\beta h}$. Such a measure is called “canonical ensemble” (or “Gibbs measure”) corresponding to H and β (or $H + h$ and β). A change of the temperature leads to tilting, too.

Tilting on \mathbb{R}^d is a slight generalization of 2a1(a) toward 2a1(b); x and t run over \mathbb{R}^d , and $\langle t, x \rangle$ replaces tx ; M and Λ are defined on \mathbb{R}^d (but still real-valued, or $+\infty$).

The general case 2a1(b) boils down to the tilting on \mathbb{R}^d (sometimes even to $d = 1$, that is, 2a1(a)) as follows. Given μ and u as in 2a1(b), we consider the distribution of u under μ , that is, the pushforward probability measure ν on \mathbb{R} (denoted often $u_*(\mu)$ or $\mu \circ u^{-1}$) defined by

$$\nu([a, b]) = \mu(u^{-1}([a, b])) = \mu(\{\omega : a \leq u(\omega) \leq b\}) \quad \text{for } -\infty < a < b < +\infty,$$

that is,²

$$\int f \, d\nu = \int (f \circ u) \, d\mu$$

for all bounded continuous functions $f : \mathbb{R} \rightarrow \mathbb{R}$ (and therefore also for all ν -integrable f). Then the distribution of u under μ_u is

$$\nu_1 = u_*(\mu_u)$$

since

$$\begin{aligned} M_\nu(1) &= \int e^x \nu(dx) = \int e^u \, d\mu = M_\mu(u) < \infty; \\ \int f(x) \nu_1(dx) &= \frac{1}{M_\nu(1)} \int f(x) e^x \nu(dx) = \\ &= \frac{1}{M_\mu(u)} \int (f \circ u) e^u \, d\mu = \int (f \circ u) \, d\mu_u. \end{aligned}$$

Likewise, given μ and u as above, and another measurable function v on the same measurable space, we consider the joint distribution $\nu = (u, v)_*(\mu)$ of u and v under μ , that is,

$$\nu([a, b] \times [c, d]) = \mu(\{\omega : a \leq u(\omega) \leq b, c \leq v(\omega) \leq d\}),$$

¹See “Canonical ensemble” and “Gibbs measure” in Wikipedia.

²See “Pushforward measure” in Wikipedia.

and get the joint distribution of u and v under μ_u ,

$$\nu_{(1,0)} = (u, v)_*(\mu_u).$$

Also,

$$\nu_{(s,t)} = (u, v)_*(\mu_{su+tv}) \quad \text{whenever } M_\mu(su + tv) < \infty.$$

Change of temperature is a special case. Let μ be the canonical ensemble for H and β_1 ; then the canonical ensemble for H and β_2 is $\mu_{(\beta_1-\beta_2)H}$; and if ν is the joint distribution of H and u at β_1 (that is, under μ), then $\nu_{(\beta_1-\beta_2,0)}$ is the joint distribution of H and u at β_2 (that is, under $\mu_{(\beta_1-\beta_2)H}$). And the distribution of u at β_2 is the corresponding marginal distribution (one-dimensional projection of the two-dimensional distribution).

Likewise, the change of the joint distribution of u_1, \dots, u_k when $\beta_1 H_1$ is replaced with $\beta_2 H_2$ boils down to tilting in \mathbb{R}^{k+2} .

In statistics, a natural exponential family on \mathbb{R} consists of probability measures, parametrized by $\theta \in \mathbb{R}$, with the density $f(\cdot|\theta)$ of the form¹

$$f(x|\theta) = h(x) \exp(\theta x - A(\theta)).$$

Clearly, $f(\cdot|\theta)$ is the tilted $f(\cdot|0)$, and $A(\cdot)$ is (up to an additive constant) the cumulant generating function.

Also the Esscher transform² is another name of tilting.

2b Surprisingly useful generating functions

The generating functions M_μ and $\Lambda_\mu = \ln M_\mu$, defined in 2a1, are surprisingly useful. So much useful that physicists often calculate in terms of these functions only,³ without mentioning tilted measures!

First, let μ be a *compactly supported* probability measure on \mathbb{R} . Then Λ_μ is finite on the whole \mathbb{R} , and for all t ,

$$M_\mu(t) = \int e^{tx} \mu(dx) = \int \left(\sum_{k=0}^{\infty} \frac{1}{k!} t^k x^k \right) \mu(dx) = \sum_{k=0}^{\infty} \frac{1}{k!} t^k \int x^k \mu(dx),$$

which justifies the name “moment generating function”.

¹See “Exponential family” and “Natural exponential family” in Wikipedia.

²See “Esscher transform” in Wikipedia.

³Physicists call M_μ the partition function and denote it $Z_n(\beta)$; they also denote Λ_μ by $\varphi(\beta)$ and call either $\varphi(\beta)/\beta$ or $\varphi(\beta)$ the (canonical) free energy. (See page 30 in “The large deviation approach to statistical mechanics” by H. Touchette, Physics Reports 2009, 478 1–69.)

In particular, for small t ,

$$\begin{aligned}\Lambda_\mu(t) &= \ln M_\mu(t) = \ln \left(1 + tM'_\mu(0) + \frac{1}{2}t^2M''_\mu(0) + O(t^3) \right) = \\ &= tM'_\mu(0) + \frac{1}{2}t^2M''_\mu(0) - \frac{1}{2}t^2(M'_\mu(0))^2 + O(t^3) = \\ &= t \int x \mu(dx) + \frac{1}{2}t^2 \left(\int x^2 \mu(dx) - \left(\int x \mu(dx) \right)^2 \right) + O(t^3),\end{aligned}$$

that is, $\Lambda'_\mu(0)$ is the expectation of μ , and $\Lambda''_\mu(0)$ is the variance of μ . In fact, the derivatives $\Lambda^{(m)}(0)$ are the so-called cumulants of μ .

The equality

$$\Lambda_\mu(t) + \Lambda_{\mu_t}(s) = \Lambda_\mu(t + s)$$

follows from 2a6. Differentiating it in s at $s = 0$ we get

$$\Lambda_\mu^{(k)}(t) = \Lambda_{\mu_t}^{(k)}(0);$$

in particular, $\Lambda'_\mu(t)$ and $\Lambda''_\mu(t)$ are the expectation and the variance of μ_t .

The variance cannot be negative, therefore Λ_μ is convex. Moreover, it is strictly convex, unless μ is a single atom. Another proof of the convexity uses Hölder's inequality: for $s, t \in \mathbb{R}$ and $\alpha, \beta > 0$ with $\alpha + \beta = 1$,

$$\begin{aligned}M_\mu(\alpha s + \beta t) &= \int (e^{sx})^\alpha (e^{tx})^\beta \mu(dx) \leq \\ &\leq \left(\int e^{sx} \mu(dx) \right)^\alpha \left(\int e^{tx} \mu(dx) \right)^\beta = M_\mu^\alpha(s) M_\mu^\beta(t);\end{aligned}$$

take the logarithm.

In general, a probability measure μ on \mathbb{R} need not be compactly supported. Rather, $\mu_k \uparrow \mu$ for some compactly supported subprobability measures μ_k . Accordingly, $\Lambda_{\mu_k} \uparrow \Lambda_\mu$. Convexity of Λ_{μ_k} implies convexity of Λ_μ , therefore, convexity of the set $\{t : \Lambda_\mu(t) < \infty\}$. This set is an interval, containing 0, but not always of the form (a, b) ; it can be $[a, b]$, $[a, b)$, $(a, b]$; it can be unbounded from below, from above, or both; and it can be $\{0\}$.

2b1 Exercise. Find examples (of μ) for all these possibilities.

Consider the interior

$$G = \{t : \Lambda_\mu(t) < \infty\}^\circ = (a, b), \quad -\infty \leq a \leq 0 \leq b \leq +\infty.$$

Leaving aside the trivial case $a = 0 = b$, we get a convex $\Lambda_\mu : (a, b) \rightarrow \mathbb{R}$.

2b2 Lemma. Λ_μ is real-analytic¹ on (a, b) .

Proof. It is sufficient to prove that M_μ is real-analytic.² Let $a < t - \varepsilon < t < t + \varepsilon < b$; we'll prove that M_μ on $[t - \varepsilon, t + \varepsilon]$ is the sum of a power series. By 2a6(a), $M_{\mu_t}(\pm\varepsilon) < \infty$, and it is sufficient to prove that M_{μ_t} on $[-\varepsilon, \varepsilon]$ is the sum of a power series. Now we forget the original μ and rename μ_t into μ . We need to prove that M_μ on $[-\varepsilon, \varepsilon]$ is the sum of a power series, given that $M_\mu(\pm\varepsilon) < \infty$. We have

$$\sum_{k=0}^{\ell} \frac{t^k x^k}{k!} \rightarrow e^{tx} \quad \text{as } \ell \rightarrow \infty,$$

$$\forall \ell \quad \left| \sum_{k=0}^{\ell} \frac{t^k x^k}{k!} \right| \leq e^{|tx|} \leq e^{-\varepsilon x} + e^{\varepsilon x}$$

for $|t| \leq \varepsilon$. By the dominated convergence theorem,

$$M_\mu(t) = \int e^{tx} \mu(dx) = \lim_{\ell} \int \sum_{k=0}^{\ell} \frac{t^k x^k}{k!} \mu(dx) = \sum_{k=0}^{\infty} \frac{t^k}{k!} \int x^k \mu(dx).$$

□

Given a measure μ on \mathbb{R}^d we have for $a \in \mathbb{R}^d$ and $t \in \mathbb{R}$

$$M_\mu(ta) = \int e^{\langle ta, x \rangle} \mu(dx) = \int e^{ty} \nu(dy) = M_\nu(t)$$

where ν is the distribution of $\langle a, x \rangle$ under μ . Assuming that M_μ is finite on some neighborhood of 0 we see that³

$$\left. \frac{d}{dt} \right|_{t=0} M_\mu(ta) = M'_\nu(0) = \int y \nu(dy) = \int \langle a, x \rangle \mu(dx).$$

Similarly,

$$\left. \frac{d^k}{dt^k} \right|_{t=0} M_\mu(ta) = \int \langle a, x \rangle^k \mu(dx).$$

Lemma 2b2 gives

$$M_\mu(a) = \sum_{k=0}^{\infty} \frac{t^k}{k!} \int \langle a, x \rangle^k \mu(dx)$$

¹That is, locally a sum of a power series.

²However, the radius of convergence for Λ_μ may be smaller because of zeros of M_μ on the complex plane.

³This derivative is a linear function of a , but be careful: this fact itself does not ensure that M_μ is differentiable at 0.

for all a in a neighborhood of 0,¹ which shows that M_μ (and therefore also Λ_μ) is real-analytic near 0.

By 2a6, $M_\mu(a+b) = M_\mu(a)M_{\mu_a}(b)$ for $a, b \in \mathbb{R}^d$ such that $M_\mu(a) < \infty$. Thus, all said about M_μ and Λ_μ around 0 applies also to $M_\mu(a+\cdot)/M_\mu(a)$ and $\Lambda_\mu(a+\cdot) - \Lambda_\mu(a)$. The functions M_μ and Λ_μ are real-analytic on the interior G of the set $\{a : M_\mu(a) < \infty\}$. For all $a \in G$ and $b \in \mathbb{R}^d$,

$$\begin{aligned} \frac{d}{dt}\Big|_{t=0} \Lambda_\mu(a+tb) &= \int \langle b, x \rangle \mu_a(dx), \\ \frac{d^2}{dt^2}\Big|_{t=0} \Lambda_\mu(a+bt) &= \int \langle b, x \rangle^2 \mu(dx) - \left(\int \langle b, x \rangle \mu(dx) \right)^2, \end{aligned}$$

the expectation and the variance of $\langle b, \cdot \rangle$ under μ_a . It follows that Λ_μ is convex on G , and moreover, strictly convex, unless μ sits on some affine subspace of dimension $d-1$ (or less). On the other hand, by Hölder's inequality, Λ_μ is convex on the whole \mathbb{R}^d , thus, the set $\{a : \Lambda_\mu(a) < \infty\}$ is convex, and its interior G is also convex.

2c Independent summands

For two independent random variables X and Y , the distribution μ_{X+Y} of the sum $X+Y$ is the convolution $\mu_X * \mu_Y$ of their distributions;

$$(2c1) \quad \int f(z) (\mu_X * \mu_Y)(dz) = \iint f(x+y) \mu_X(dx) \mu_Y(dy).$$

Taking $f(z) = e^{tz}$ we have $f(x+y) = f(x)f(y)$, thus,

$$(2c2) \quad M_{X+Y}(t) = M_X(t)M_Y(t), \quad \Lambda_{X+Y}(t) = \Lambda_X(t) + \Lambda_Y(t).$$

In terms of convolution,

$$(2c3) \quad M_{\mu*\nu}(t) = M_\mu(t)M_\nu(t); \quad \Lambda_{\mu*\nu}(t) = \Lambda_\mu(t) + \Lambda_\nu(t).$$

(However, this does not apply to $M(u), \Lambda(u)$.) Here is the tilted convolution:

$$(2c4) \quad (\mu * \nu)_t = \mu_t * \nu_t \quad \text{whenever } M_\mu(t), M_\nu(t) < \infty,$$

since

$$\int f d(\mu * \nu)_t = \frac{1}{M_{\mu*\nu}(t)} \int f(z) e^{tz} (\mu * \nu)(dz) =$$

¹In fact, in every ball (centered at 0) on which $M_\mu < \infty$; recall the proof of 2b2.

$$\begin{aligned}
&= \frac{1}{M_\mu(t)M_\nu(t)} \iint f(x+y)e^{t(x+y)} \mu(dx)\nu(dy) = \\
&= \iint f(x+y)\mu_t(dx)\nu_t(dy) = \int f d(\mu_t * \nu_t)
\end{aligned}$$

for all bounded continuous f .

We turn to the sum of independent, identically distributed (i.i.d.) random variables; its distribution is

$$\mu^{*n} = \underbrace{\mu * \cdots * \mu}_n.$$

By (2c4),

$$(2c5) \quad (\mu^{*n})_t = (\mu_t)^{*n},$$

thus we need not hesitate writing just μ_t^{*n} .

The Legendre transform Λ_μ^* of Λ_μ will be very useful:

$$\Lambda_\mu^*(x) = \sup_{t \in \mathbb{R}} (tx - \Lambda_\mu(t)) \in [0, \infty].$$

If $x = \Lambda'_\mu(t)$ for some $t \in G$ (that is, $\Lambda_\mu < \infty$ near t), then

$$\Lambda_\mu^*(x) = tx - \Lambda_\mu(t)$$

by convexity of Λ_μ .

2c6 Example (Standard normal distribution, see 2a2). $\Lambda_\mu(t) = \frac{1}{2}t^2$; $x = \Lambda'_\mu(t) = t$; $\Lambda_\mu^*(x) = x \cdot x - \frac{1}{2}x^2 = \frac{1}{2}x^2$.

2c7 Example (Fair coin, see 2a3). $\Lambda_\mu(t) = \ln \cosh t$; $x = \Lambda'_\mu(t) = \tanh t$; note that $\tanh^2 t + \frac{1}{\cosh^2 t} = 1$, thus $\cosh t = \frac{1}{\sqrt{1-x^2}}$ and $\Lambda_\mu(t) = -\frac{1}{2} \ln(1-x^2)$. Also, $t = \operatorname{artanh} x = \frac{1}{2} \ln \frac{1+x}{1-x}$, thus

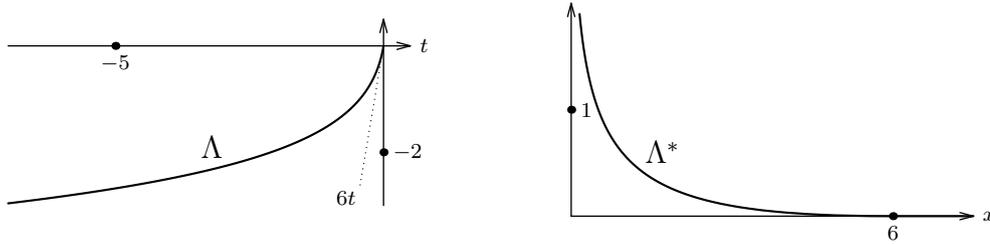
$$\Lambda_\mu^*(x) = \frac{x}{2} \ln \frac{1+x}{1-x} + \frac{1}{2} \ln(1-x^2) = \frac{1}{2}(1+x) \ln(1+x) + \frac{1}{2}(1-x) \ln(1-x)$$

for $x \in [-1, 1]$ (otherwise, ∞); just the function γ of (1a1).

2c8 Example (Exponential distribution, see 2a4). $\Lambda_\mu(t) = -\ln(1-t)$; $x = \Lambda'_\mu(t) = \frac{1}{1-t}$; $t = 1 - \frac{1}{x}$; $\Lambda_\mu^*(x) = (1 - \frac{1}{x})x - \ln x = x - 1 - \ln x$.

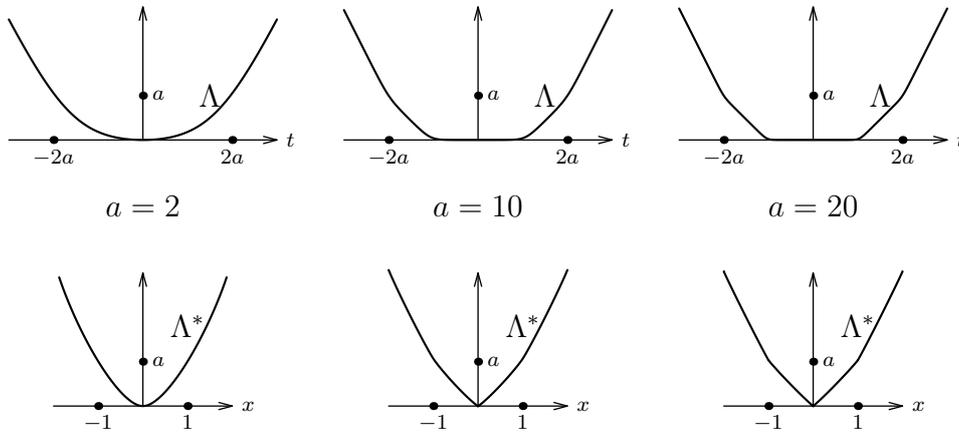
2c9 Example (Discontinuous generating function, see 2a5). $\Lambda_\mu(t) \leq 0$ for $-\infty < t \leq 0$, but $+\infty$ for $t > 0$. Nevertheless, $\int x \mu(dx) = 6 < \infty$, and $\Lambda'_\mu(0-) = 6$. In fact,

$$M_\mu(t) = \frac{1}{-2t} \left(1 - \sqrt{\frac{2\pi}{-2t}} \Phi \left(-\frac{1}{\sqrt{-2t}} \right) \exp \left(\frac{1}{-4t} \right) \right) \quad \text{for } t < 0.$$



$\Lambda^*(x) = \infty$ for $x \in (-\infty, 0]$; $0 < \Lambda^*(x) < \infty$ for $x \in (0, 6)$; and $\Lambda^*(x) = 0$ for $x \in [6, \infty)$. Thus, Λ^* fails to be real-analytic near 6.

2c10 Example (Multiscale case). $\mu(\{-1\}) = \frac{1}{2}e^{-a} = \mu(\{+1\})$, $\mu(\{-2\}) = \frac{1}{2}e^{-3a} = \mu(\{+2\})$, $\mu(\{0\}) = 1 - e^{-a} - e^{-3a}$; $\Lambda_\mu(t) = 1 + e^{-a}(-1 + \cosh t) + e^{-3a}(-1 + \cosh 2t)$.



For large a both functions are approximately piecewise linear. Note that the variance (and higher moments) of μ_t is not at all monotone in t .

2c11 Lemma. Let $t \in G$, $x = \Lambda'_\mu(t)$, and $\varepsilon > 0$; then $\mu_t^{*n}([nx - \varepsilon, nx + \varepsilon]) > 0$ if and only if $\mu^{*n}([nx - \varepsilon, nx + \varepsilon]) > 0$, and in this case

$$\left| \ln \frac{\mu_t^{*n}([nx - \varepsilon, nx + \varepsilon])}{\mu^{*n}([nx - \varepsilon, nx + \varepsilon])} - n\Lambda_\mu^*(x) \right| \leq \varepsilon|t|.$$

Proof. Follows from the inequality

$$\left| \ln \frac{d\mu_t^{*n}}{d\mu^{*n}}(y) - n\Lambda_\mu^*(x) \right| \leq \varepsilon|t| \quad \text{for all } y \in [nx - \varepsilon, nx + \varepsilon],$$

checked easily:

$$\ln \frac{d\mu_t^{*n}}{d\mu^{*n}}(y) = ty - \Lambda_{\mu^{*n}}(t) = t(y - nx) + tnx - n\Lambda_\mu(t) = t(y - nx) + n\Lambda_\mu^*(x).$$

□

We see that $\ln \mu^{*n}([nx - \varepsilon, nx + \varepsilon])$ is $\varepsilon|t|$ -close to $\ln \mu_t^{*n}([nx - \varepsilon, nx + \varepsilon]) - n\Lambda_\mu^*(x)$. An upper bound follows immediately:

$$(2c12) \quad \ln \mu^{*n}([nx - \varepsilon, nx + \varepsilon]) \leq -n\Lambda_\mu^*(x) + \varepsilon|t|.$$

A lower bound needs more effort. The measure μ_t^{*n} has the expectation $n\Lambda'_\mu(t) = nx$ and the variance $n\Lambda''_\mu(t)$; by Chebyshev's inequality,

$$\mu_t^{*n}([nx - \varepsilon, nx + \varepsilon]) \geq 1 - \frac{n\Lambda''_\mu(t)}{\varepsilon^2},$$

which leads to the lower bound

$$(2c13) \quad \ln \mu^{*n}([nx - \varepsilon, nx + \varepsilon]) \geq -n\Lambda_\mu^*(x) - \varepsilon|t| + \ln \left(1 - \frac{n\Lambda''_\mu(t)}{\varepsilon^2} \right).$$

2c14 Theorem. Let $t \in G$, $x = \Lambda'_\mu(t)$, and $\varepsilon_n > 0$ satisfy

$$\frac{\varepsilon_n}{n} \rightarrow 0, \quad \frac{\varepsilon_n}{\sqrt{n}} \rightarrow \infty.$$

Then

$$\frac{1}{n} \ln \mu^{*n}([nx - \varepsilon_n, nx + \varepsilon_n]) \rightarrow -\Lambda_\mu^*(x) \quad \text{as } n \rightarrow \infty.$$

Proof. The upper limit is at most $-\Lambda_\mu^*(x)$ by (2c12). Taking into account that $\frac{n\Lambda''_\mu(t)}{\varepsilon_n^2} \rightarrow 0$ we see that the lower limit is at least $-\Lambda_\mu^*(x)$ by (2c13). □

Here is the same result in a slightly different language.

2c15 Theorem. Let $t \in \mathbb{R}$, and X_1, X_2, \dots be i.i.d. random variables such that $\ln \mathbb{E} \exp \lambda X_1 = \Lambda(\lambda) < \infty$ for all λ close enough to t . Let $\varepsilon_n > 0$ satisfy

$$\varepsilon_n \rightarrow 0, \quad \sqrt{n} \varepsilon_n \rightarrow \infty.$$

Denote $x = \Lambda'(t)$. Then

$$\mathbb{P} \left(x - \varepsilon_n \leq \frac{X_1 + \dots + X_n}{n} \leq x + \varepsilon_n \right) = \exp \left(-n(tx - \Lambda(t)) + o(n) \right) \quad \text{as } n \rightarrow \infty.$$

Now we turn to moderate deviations. Here we assume that $0 \in G$, and in addition, $\Lambda'_\mu(0) = 0$, $\Lambda''_\mu(0) = 1$ (otherwise, use a linear transformation).

2c16 Theorem. Let $x_n \rightarrow 0$, $\sqrt{n}|x_n| \rightarrow \infty$, and $\varepsilon_n > 0$ satisfy

$$\frac{\varepsilon_n}{n|x_n|} \rightarrow 0, \quad \frac{\varepsilon_n}{\sqrt{n}} \rightarrow \infty.$$

Then

$$\ln \mu^{*n}([nx_n - \varepsilon_n, nx_n + \varepsilon_n]) = -\frac{1}{2}nx_n^2(1 + o(1)) \quad \text{as } n \rightarrow \infty.$$

Proof. We take $t_n \rightarrow 0$ such that $x_n = \Lambda'_\mu(t_n)$ and note that $x_n \sim t_n$ (that is, their ratio converges to 1), $\Lambda_\mu(t_n) \sim \frac{1}{2}t_n^2$, and $\Lambda_\mu^*(x_n) = t_n x_n - \Lambda_\mu(t_n) \sim \frac{1}{2}x_n^2$.

By (2c12), $\ln(\dots) \leq -n\Lambda_\mu^*(x_n) + \varepsilon_n|t_n| = -n \cdot \frac{1}{2}x_n^2(1 + o(1)) + \varepsilon_n|x_n|(1 + o(1)) = -\frac{1}{2}nx_n^2(1 + o(1))$, since $\varepsilon_n|x_n| \ll nx_n^2$.

By (2c13), taking into account that $\frac{n\Lambda''_\mu(t_n)}{\varepsilon_n^2} \sim \frac{n}{\varepsilon_n^2} \rightarrow 0$, we get $\ln(\dots) \geq -n\Lambda_\mu^*(x_n) - \varepsilon_n|t_n| + o(1) = -\frac{1}{2}nx_n^2(1 + o(1)) + o(1) = -\frac{1}{2}nx_n^2(1 + o(1))$, since $nx_n^2 \rightarrow \infty$. \square

And the same result in the slightly different language.

2c17 Theorem. Let X_1, X_2, \dots be i.i.d. random variables such that $\ln \mathbb{E} \exp \lambda X_1 < \infty$ for all λ close enough to 0, and $\mathbb{E} X_1 = 0$, $\mathbb{E} X_1^2 = 1$. Let $x_n \in \mathbb{R}$ and $\varepsilon_n > 0$ satisfy

$$|x_n| \rightarrow \infty, \quad x_n = o(\sqrt{n}), \quad \varepsilon_n \rightarrow 0, \quad |x_n|\varepsilon_n \rightarrow \infty.$$

Then

$$\mathbb{P}\left(x_n(1 - \varepsilon_n) \leq \frac{X_1 + \dots + X_n}{\sqrt{n}} \leq x_n(1 + \varepsilon_n)\right) = \exp\left(-\frac{1}{2}x_n^2(1 + o(1))\right)$$

as $n \rightarrow \infty$.

2c18 Exercise. Generalize these results (2c11, 2c14–2c17) to probability measures on \mathbb{R}^d ; in the other language, to i.i.d. random vectors.

The condition “ $\frac{\varepsilon_n}{\sqrt{n}} \rightarrow \infty$ ” in Theorem 2c14 may be replaced with $\frac{\varepsilon_n}{\sqrt{n}} \geq \text{const}$ with an appropriate absolute constant (think, why). The same applies to “ $\sqrt{n}\varepsilon_n \rightarrow \infty$ ” in 2c15, “ $\frac{\varepsilon_n}{\sqrt{n}} \rightarrow \infty$ ” in 2c16, and “ $|x_n|\varepsilon_n \rightarrow \infty$ ” in 2c17.

Much better bounds are obtained via the Berry-Esseen bound for the central limit theorem (CLT). By CLT, the distribution of $\frac{X_1 + \dots + X_n}{\sqrt{n}}$ converges

weakly to the standard normal distribution (assuming $\mathbb{E} X_1 = 0$ and $\mathbb{E} X_1^2 = 1$). That is,

$$\sup_{-\infty < a < b < +\infty} |\mu^{*n}([\sqrt{n}a, \sqrt{n}b]) - (\Phi(b) - \Phi(a))| \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$

By the Berry-Esseen bound, this supremum never exceeds $\text{const} \cdot \mathbb{E} |X_1|^3 / \sqrt{n}$, the constant being absolute.¹ Using the fact that $(\mathbb{E} |X_1|^3)^{1/3} \leq (\mathbb{E} X_1^4)^{1/4}$ we get

$$\left| \mu_t^{*n}([nx - \varepsilon, nx + \varepsilon]) - \left(2\Phi\left(\frac{\varepsilon}{\sqrt{n\Lambda_\mu''(t)}}\right) - 1 \right) \right| \leq \frac{\text{const}}{\sqrt{n}} \cdot \left(\frac{\Lambda_\mu^{(4)}(t)}{(\Lambda_\mu''(t))^2} + 3 \right)^{3/4}.$$

Thus, we may take ε that depends on $\Lambda_\mu^{(4)}(t)$ and $\Lambda_\mu''(t)$ but does not depend on n , and get for $\mu_t^{*n}([nx - \varepsilon, nx + \varepsilon])$ a lower bound of order $1/\sqrt{n}$, which leads to the lower bound

$$\ln \mu^{*n}([nx - \varepsilon, nx + \varepsilon]) \geq -n\Lambda_\mu^*(x) - \mathcal{O}(|t|) - \frac{1}{2} \ln n - \mathcal{O}(1).$$

The same Berry-Esseen bound gives $\mu_t^{*n}([nx - \varepsilon, nx + \varepsilon]) = \mathcal{O}(1/\sqrt{n})$, which leads to the upper bound

$$\ln \mu^{*n}([nx - \varepsilon, nx + \varepsilon]) \leq -n\Lambda_\mu^*(x) + \mathcal{O}(|t|) - \frac{1}{2} \ln n + \mathcal{O}(1).$$

In both cases, LDP and MDP, t is bounded (in n); also $\Lambda_\mu^{(4)}(t)$ and $\Lambda_\mu''(t)$ are bounded; thus,

$$\mu^{*n}([nx - \varepsilon, nx + \varepsilon]) = \frac{1}{\sqrt{n}} \exp(-n\Lambda_\mu^*(x) + \mathcal{O}(1)).$$

Here are sLD-counterparts of the LD-theorems 2c14, 2c15.

2c19 Theorem. Let $t \in G$ and $x = \Lambda'_\mu(t)$. Then for every $\varepsilon > 0$ large enough,

$$\mu^{*n}([nx - \varepsilon, nx + \varepsilon]) = \frac{1}{\sqrt{n}} \exp(-n\Lambda_\mu^*(x) + \mathcal{O}(1)) \quad \text{as } n \rightarrow \infty.$$

2c20 Theorem. Let $t \in \mathbb{R}$, and X_1, X_2, \dots be i.i.d. random variables such that $\ln \mathbb{E} \exp \lambda X_1 = \Lambda(\lambda) < \infty$ for all λ close enough to t . Denote $x = \Lambda'(t)$. Then for every $\varepsilon > 0$ large enough,

$$\mathbb{P}\left(x - \frac{\varepsilon}{n} \leq \frac{X_1 + \dots + X_n}{n} \leq x + \frac{\varepsilon}{n}\right) = \frac{1}{\sqrt{n}} \exp\left(-n(tx - \Lambda(t)) + \mathcal{O}(1)\right)$$

as $n \rightarrow \infty$.

¹See ‘‘Berry-Esseen theorem’’ in Wikipedia.

Think, what happens for the fair coin case, if $\varepsilon < 1/2$.

It is possible to get an approximation up to equivalence (that is, $o(1)$ instead of $\mathcal{O}(1)$ under $\exp(\dots)$), but not easily. To this end, first of all, one has to separate lattice and non-lattice distributions, and not only in proofs but also in formulations.

Now, what about sMD? Here we assume (as before) that $0 \in G$, $\Lambda'_\mu(0) = 0$, $\Lambda''_\mu(0) = 1$, and $x_n \rightarrow 0$, $\sqrt{n}|x_n| \rightarrow \infty$. We take (again) $t_n \rightarrow 0$ such that $x_n = \Lambda'_\mu(t_n)$; still, $x_n \sim t_n$, $\Lambda_\mu(t_n) \sim \frac{1}{2}t_n^2$, and $\Lambda_\mu^*(x_n) \sim \frac{1}{2}x_n^2$. However, now this relation does not satisfy us! Now we need $\Lambda_\mu^*(x_n) = \frac{1}{2}x_n^2 + \mathcal{O}(\frac{1}{n})$ in order to get the normal approximation $\frac{1}{\sqrt{n}} \exp(-\frac{n}{2}x_n^2 + \mathcal{O}(1))$.

The function Λ_μ^* is real-analytic near 0, which follows from the equality $\Lambda_\mu^*(\Lambda'_\mu(t)) = t\Lambda'_\mu(t) - \Lambda_\mu(t)$, since Λ_μ is real-analytic near 0, $\Lambda'_\mu(0) = 0$, and $\Lambda''_\mu(0) = 1 \neq 0$ (indeed, the inverse function to Λ'_μ is real-analytic near 0). For small x we have $\Lambda_\mu^*(x) \sim \frac{1}{2}x^2$, thus,

$$\Lambda_\mu^*(x) = \frac{1}{2}x^2 - a_0x^3 - a_1x^4 - \dots$$

The numbers a_0, a_1, \dots are called the coefficients of the Cramer series.¹ In particular,²

$$a_0 = \frac{1}{6}\Lambda_\mu^{(3)}(0); \quad a_1 = \frac{1}{24}(\Lambda_\mu^{(4)}(0) - 3(\Lambda_\mu^{(3)}(0))^2).$$

If $x_n = \mathcal{O}(n^{-1/3})$ then indeed $\Lambda_\mu^*(x_n) = \frac{1}{2}x_n^2 + \mathcal{O}(\frac{1}{n})$, and we get sMD-counterparts of Theorems 2c16, 2c17.

2c21 Theorem. Let $x_n = \mathcal{O}(n^{-1/3})$. Then for every $\varepsilon > 0$ large enough,

$$\mu^{*n}([nx_n - \varepsilon, nx_n + \varepsilon]) = \frac{1}{\sqrt{n}} \exp(-\frac{1}{2}nx_n^2 + \mathcal{O}(1)) \quad \text{as } n \rightarrow \infty.$$

2c22 Theorem. Let X_1, X_2, \dots be i.i.d. random variables such that $\ln \mathbb{E} \exp \lambda X_1 < \infty$ for all λ close enough to 0, and $\mathbb{E} X_1 = 0$, $\mathbb{E} X_1^2 = 1$. Let $x_n \in \mathbb{R}$ satisfy

$$x_n = \mathcal{O}(n^{1/6}).$$

¹Some authors define the Cramer series as $a_0 + a_1x + \dots$ (V.V. Petrov and J. Robinson 2008, "Large deviations for sums of independent non identically distributed random variables", Communications in Statistics **37** 2984–2990); others define it as $a_0x^3 + a_1x^4 + \dots$ (L.V. Rozovsky 1999, "On the Cramér series coefficients", Theory Probab. Appl. **43** 152–157).

²For a_2, a_3 and a formula for a_k see Rozovsky 1999.

Then for every $\varepsilon > 0$ large enough,

$$\mathbb{P}\left(x_n - \frac{\varepsilon}{\sqrt{n}} \leq \frac{X_1 + \cdots + X_n}{\sqrt{n}} \leq x_n + \frac{\varepsilon}{\sqrt{n}}\right) = \frac{1}{\sqrt{n}} \exp\left(-\frac{1}{2}x_n^2 + \mathcal{O}(1)\right)$$

as $n \rightarrow \infty$.

If $a_0 = 0$, that is, $\mathbb{E} X_1^3 = 0$ (in particular, for all symmetric distributions, for example, the fair coin), then “ $n^{-1/3}$ ” in Theorem 2c21 may be replaced with “ $n^{-1/4}$ ”, and “ $n^{1/6}$ ” in Theorem 2c22 with “ $n^{1/4}$ ”. In general, under these conditions we get “ $-\frac{1}{2}nx_n^2 + a_0nx_n^3$ ” instead of “ $-\frac{1}{2}nx_n^2$ ” in Theorem 2c21, and “ $-\frac{1}{2}x_n^2 + \frac{a_0}{\sqrt{n}}x_n^3$ ” instead of “ $-\frac{1}{2}x_n^2$ ” in Theorem 2c22. The new factor, being $\exp(\mathcal{O}(n^{1/4}))$, matters for SMD but does not matter for MD.

That is, under $n^{1/6}$ (in terms of 2c22) all distributions μ are served by a single, normal approximation. Between $n^{1/6}$ and $n^{1/4}$ they are not; a one-parameter family of approximations is needed. Likewise, between $n^{1/4}$ and $n^{3/10}$, two parameters are needed (a_0 and a_1 ; or $\mathbb{E} X_1^3$ and $\mathbb{E} X_1^4$). And generally, k parameters work between $n^{k/(2(k+2))}$ and $n^{(k+1)/(2(k+3))}$. Somehow, $k = \infty$ means $n^{1/2}$, — the LD territory; and indeed, LD uses a function Λ_μ^* that depends on all μ (rather than several parameters of μ).

In contrast, in the framework of MD (rather than SMD) the normal approximation works in the whole domain $o(n^{1/2})$; the dependence on μ appears at once when $\mathcal{O}(n^{1/2})$ is reached.

Index

Berry-Esseen, 18	moment generating functional, 8
canonical ensemble, 9	real-analytic, 12
convolution, 13	
Cramer series, 19	tilted measure, 7
cumulant, 11	
cumulant generating function, 8	$\mu * \nu$, 13
cumulant generating functional, 8	μ^{*n} , 14
	G , 13
i.i.d., 14	Λ^* , 14
Legendre transform, 14	Λ_μ , 8
	M_μ , 8
MGF, 8	μ_t , 7
moment generating function, 8	μ_u , 8