

Problem set 3

Math 212a

Thursday, October 14, 2004, Due Oct. 21

The purpose of this problem set is to walk through the proof of the “central limit theorem” of probability theory. Roughly speaking, this theorem asserts that if the random variable S_n is the sum of many independent random variables

$$S_n = X_1 + \cdots + X_n$$

all with mean zero and finite variance then under appropriate additional hypotheses

$$\frac{1}{s_n} S_n$$

is approximately normally distributed where

$$s_n^2 := \text{var}(S_n) = \sigma_1^2 + \cdots + \sigma_n^2, \quad \sigma_i^2 := \text{var}(X_i).$$

The actual condition, equation (12) below, is rather technical. But roughly speaking it says that no one X_i outweighs the others. The condition of mean zero for each of the X_i is just a matter of convenience in formulation. Otherwise we would replace the statement by the the assertion that $\frac{1}{s_n}(S_n - E(S_n))$ is approximately normally distributed. The phrase “approximately normally distributed” will be taken to mean in the sense of weak convergence: “Weak convergence” says that for each bounded continuous function f on \mathbb{R} we have

$$\lim_{n \rightarrow \infty} E \left(f \left(\frac{1}{s_n} S_n \right) \right) \rightarrow E(f(N)) \quad (1)$$

where N is a unit normal random variable:

$$P(N \in [a, b]) = \frac{1}{\sqrt{2\pi}} \int_a^b e^{-x^2/2} dx.$$

It is enough to prove this for a set of functions f which are dense in the uniform norm. In fact, the proof outlined below will assume that f is three times differentiable with a uniform bound on its first three derivatives.

Philosophically, the import of the theorem is that if we observe an effect whose variation is the sum of many small positive and negative random contributions, no one outweighing the others, then the effect will be distributed

according to a normal distribution with a finite mean and variance. So this theorem accounts for the ubiquity of the normal distribution in science.

The idea of the proof presented here is the same as the idea behind the proof that we gave in Problem Set 2 of Poisson's law of small numbers. We replace each X_k successively in the sum giving S_n by a normal random variable Z_k having mean zero and variance σ_k . We then estimate the total error committed by these substitutions.

Let me be clear about the definition of independence: Events e_1, e_2, \dots , are said to be independent if for any subcollection

$$P(e_{i_1} \cap \dots \cap e_{i_k}) = P(e_{i_1}) \cdot P(e_{i_2}) \cdots P(e_{i_k}).$$

The random variables X_1, X_2, \dots are called independent if for all choices of a_i, b_i the events $\mathbf{1}_{a_1 \leq X_1 \leq b_1}, \mathbf{1}_{a_2 \leq X_2 \leq b_2}, \dots$ are independent.

Contents

1 Push Forward.	2
2 Moment generating functions	7
3 Central limit theorem by the coupling method.	9

1 Push Forward.

I am going to use some language of measure theory which we have not yet gotten to in class. You can look ahead in the notes to discover the meaning of these terms or you can blissfully ignore the precise meaning of some of this language: You can simply take the term “ σ -field” to mean a collection of subsets of a given set satisfying some arcane conditions and a “measure space” to mean a set with a chosen σ -field and with a function m (also satisfying some conditions) which assigns to every element of this σ -field a non-negative real number or $+\infty$ called the “measure” of this element (= subset). The computations below should be far more transparent than these definitions.

Let (M, \mathcal{F}, m) be a measure space and (N, \mathcal{G}) a space Y with a σ -field \mathcal{G} . We call a map $f : M \rightarrow N$ “measurable” if $f^{-1}(B) \in \mathcal{F}$ for every $B \in \mathcal{G}$ and then define the push forward f_*m of the measure m by

$$f_*(B) = m(f^{-1}(B)).$$

We will need some formulas for pushforward. In the cases we consider we will have $M = [0, 1]$ (so frequently we will think of push forward as the process of “simulating” a random variable using a random number generator as explained in Problem set 2), or we will have $M = \mathbf{R}^n$ or some nice subset of \mathbf{R}^n . Similarly for N .

In all our examples, measures will either be discrete or have densities relative to Lebesgue measure. (Think of Lebesgue measure for the moment as the rule which assigns to each interval (on the real line) its length and assigns to each Riemann integrable function f its integral $\int_{\mathbb{R}} f(x)dx$. Similarly for Lebesgue measure on \mathbb{R}^n .)

A discrete measure ν integrates functions by the formula

$$\langle \nu, \phi \rangle = \int \phi \nu = \sum \phi(x_k) r_k, \quad r_k = \nu(\{x_k\}).$$

Here the x_k are a finite or countable subset of X . The push forward of the discrete measure ν by a map f is concentrated on the set

$$\{y_\ell\} = \{f(x_k)\} \quad \text{with} \quad (f_*\nu)(\{y_\ell\}) = \sum_{f(x_k)=y_\ell} \nu(\{x_k\}).$$

For functions we have

$$\langle f_*\nu, \phi \rangle = \langle \nu, f^*\phi \rangle.$$

In this equation and in what follows we use the “pull back” notation

$$f^*\phi := \phi \circ f.$$

The measures with density have

$$\langle \nu, \phi \rangle = \int \phi(x) \rho(x) dx$$

where ρ is the density.

For the standard Lebesgue linear measure du (so density one) and for $f : u \mapsto x$ a map of real intervals which is one to one, continuously differentiable and with continuously differentiable inverse we have

$$f_*du = |f'(u(x))|^{-1} dx = |du/dx| dx. \quad (2)$$

Proof. By the change of variables formula

$$\int \psi dx = \int f^*\psi |dx/du| du.$$

Set

$$\psi = \phi |dx/du|^{-1}.$$

The change of variables formula becomes

$$\begin{aligned} \int \phi |dx/du|^{-1} dx &= \int f^*\phi f^*|dx/du|^{-1} |dx/du| du \\ &= \int f^*\phi du \end{aligned}$$

by the chain rule since

$$\frac{dx}{du}(x(u))\frac{du}{dx}(u) \equiv 1.$$

QED

Example:

$$x = -(1/\lambda) \ln u.$$

This maps the interval $(0, 1]$ onto the positive numbers (reversing orientation so 1 goes to 0 and 0 goes to ∞). So

$$u = e^{-\lambda x} \quad \text{and} \quad |du/dx| = \lambda e^{-\lambda x}.$$

The probability law on \mathbf{R}^+ with density $\lambda e^{-\lambda x}$ is known as the exponential law. So the function $-(1/\lambda) \ln u$ simulates the exponential law. The MATLAB command

$$-(1/\lambda) * \log(\text{rand}(M,N))$$

will produce an $M \times N$ matrix whose entries are non-negative numbers independently exponentially distributed with parameter λ .

Example:

$$r = \sqrt{-2 \ln u} \quad \text{so} \quad u = e^{-r^2/2}.$$

The push forward of the uniform measure du is

$$e^{-r^2/2} r dr.$$

The same argument shows that the push forward of $\rho(u)du$ is given by

$$f_*[\rho du] = \rho(u(x))|du/dx|dx. \quad (3)$$

The formula (3) works in n dimensions where du/dx is interpreted as the Jacobian matrix.

Example. The area in polar coordinates in the plane is given by

$$r dr d\theta.$$

If we set $S = r^2$ we can write this as

$$dA = \frac{1}{2} dS d\theta.$$

So, for example, if we want to describe a uniform probability measure on the unit disk, it will be given by

$$\frac{1}{2\pi} dS d\theta.$$

Notice that in this formula S is uniformly distributed on $[0, 1]$ and θ is uniformly distributed on $[0, 2\pi]$ and they are independent. Suppose we start with S and θ and (changing from our previous notation) set

$$r = \sqrt{-2 \ln S}.$$

We get the measure

$$\frac{1}{2\pi} e^{-r^2/2} r dr d\theta$$

for a probability measure on the plane. In particular the total integral of the above expression is one.

If we change from polar coordinates to rectangular coordinates:

$$x = r \cos \theta, \quad y = r \sin \theta$$

the density becomes

$$\frac{1}{2\pi} e^{-(x^2+y^2)/2} dx dy.$$

Notice that the random variables X and Y (projections onto the coordinate axes) are independent, each with density

$$\frac{1}{\sqrt{2\pi}} e^{-x^2/2}.$$

In particular

$$\frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-x^2/2} dx = 1.$$

Example. The rejection method.

Suppose that V_1 and V_2 are independent random variables uniformly distributed on $[-1, 1]$ and that

$$S = V_1^2 + V_2^2.$$

and so

$$\cos \theta = \frac{V_1}{S^{1/2}}, \quad \sin \theta = \frac{V_2}{S^{1/2}}$$

if we use V_1, V_2 as rectangular coordinates.

The vector (V_1, V_2) is uniformly distributed over the square $[-1, 1] \times [-1, 1]$. We would like to have a measure uniformly distributed over the unit disk. So add an instruction which rejects the output and tries again if $S > 1$. Among the accepted outputs, the probability of ending up in any region of the disk is proportional to the area of the region. Thus we have produced a uniform distribution on the unit disk, and we know from the preceding discussion that S , the radius squared, and the angle θ are independent and uniformly distributed, S on $[0, 1]$ and θ on $[0, 2\pi]$.

Substituting into the above we get that

$$X = (-2 \ln S)^{1/2} \frac{V_1}{S^{1/2}}, \quad Y = (-2 \ln S)^{1/2} \frac{V_2}{S^{1/2}}$$

are independent normally distributed random variables. So we can simulate two independent normally distributed random variables (X, Y) by the following algorithm

- Generate two independent random numbers U_1 and U_2 using the random number generator. (In MATLAB $[U_1, U_2] = \text{rand}(1,2)$).

- Set

$$V_1 = 2U_1 - 1, \quad V_2 = 2U_2 - 1, \quad S = V_1^2 + V_2^2.$$

- If $S > 1$ return to step 1. Otherwise

- Set

$$X = V_1(-2 \ln S/S)^{1/2}, \quad Y = V_2(-2 \ln S/S)^{1/2}.$$

Of course, MATLAB has its own built in normal random number generator: the command $\text{randn}(M,N)$ produces an $M \times N$ matrix whose entries are independent normally distributed random variables.

Push forward from higher to lower dimension involves summation or integration. (multiple integral as iterated integral). These may or may not converge. But they will converge if the total measure is finite, for example for a probability measure.

Example. $(x, y) \mapsto x + y$ is a map $f : \mathbf{R}^2 \rightarrow \mathbf{R}$. If m is a density on \mathbf{R}^2 and ϕ is a function on \mathbf{R} , then

$$\int \int f^* \phi m(x, y) dx dy = \int \int \phi(x + y) m(x, y) dx dy.$$

Set $z = x + y$ and $u = y$. Apply the change of variables formula; the Jacobian is identically 1. So the integral becomes

$$\int \int \phi(z) m(z - u, u) dz du = \int \phi(z) \int m(z - u, u) du dz$$

by iterated integration. So

$$f^* m = \rho$$

where

$$\rho(z) = \int m(z - u, u) du.$$

An important special case is where $m(x, y) = r(x)s(y)$ (independence). Then

$$\rho(z) = \int r(z - u)s(u) du$$

is called the convolution of r and s and written as

$$\rho = r \star s.$$

Notice that

$$r \star s = s \star r.$$

This differs slightly from the convention we use in the lecture on the Fourier transform which involves a factor of $1/\sqrt{2\pi}$. But this is the standard probabilists' convention.

2 Moment generating functions

For any random variable X try to define

$$M_X(t) := E(e^{tX}).$$

Of course this may only be defined for a limited range of t due to convergence problems. For example, suppose that X is exponentially distributed with parameter λ . Then

$$M_X(t) = \lambda \int_0^{\infty} x e^{(t-\lambda)x} dx$$

diverges for $t \geq \lambda$. But for $t < \lambda$ the integral converges to

$$\frac{\lambda}{\lambda - t}.$$

For any random variable, if t is interior to the convergence domain of the moment generating function, M_X , then M_X is differentiable to all orders at t and we have

$$M^{(n)}(t) = E(X^n e^{tX}).$$

In particular, if 0 is interior to the domain of convergence, we have

$$E(X^n) = M^{(n)}(0).$$

For example, for the exponential distribution the n -th derivative of the moment generating function is given by

$$\frac{n!\lambda}{(\lambda - t)^{n+1}}$$

and so

$$E(X^n) = n!\lambda^{-n}.$$

So

$$E(X) = \frac{1}{\lambda}, \quad E(X^2) = \frac{2}{\lambda^2}$$

and so

$$\text{Var}(X) = \frac{1}{\lambda^2}.$$

In statistical mechanics the tradition is to study $M(-\beta) =: P(\beta)$ which is called the *partition function*.

1. Let N denote unit normal random variable, so N has density $\frac{1}{\sqrt{2\pi}}e^{-x^2/2}$. Compute its moment generating function $M_N(t)$ and the first four derivatives of $M_N(t)$. Evaluate the first four moments $E(N)$, $E(N^2)$, $E(N^3)$ and $E(N^4)$. In particular verify that $\text{var}(N) = 1$.

Suppose we set

$$X = \sigma N + m. \quad (4)$$

Then

$$E(X) = m \quad \text{and} \quad \text{Var}(X) = \sigma^2 \quad (5)$$

while the change of variables formula implies that the density of X is given by

$$\frac{1}{\sigma\sqrt{2\pi}}e^{-\frac{(x-m)^2}{2\sigma^2}}. \quad (6)$$

Conversely, suppose that Z is a random variable with density

$$Ce^{-q(x)/2}$$

where

$$q(x) = ax^2 + 2bx + c, \quad a > 0,$$

is a quadratic polynomial, and C is a constant. Completing the square allows us to write

$$q(x) = a\left(x + \frac{b}{a}\right)^2 - \frac{b^2}{a^2} + c$$

and the fact that the total integral must be one constrains the constants so that

$$Ce^{c-b^2/a^2} = \frac{1}{\sigma\sqrt{2\pi}} \quad \text{where} \quad \sigma = \frac{1}{\sqrt{a}}.$$

Hence Z is equal in law to a random variable of the form (4) where

$$\sigma = \frac{1}{\sqrt{a}}, \quad m = -\frac{b}{a}.$$

A random variable whose density is proportional to $e^{-q(x)/2}$ for some quadratic polynomial is called a *Gaussian* random variable. We have seen that such a random variable is completely determined in law by its mean and variance, and is equal in law to a random variable of the form (4).

In particular, if X is a Gaussian with mean m_X and variance σ_X^2 and Y is a Gaussian with mean m_Y and variance σ_Y^2 , and if X and Y are independent, then $X + Y$ is a Gaussian with mean $m_X + m_Y$ and variance $\sigma_X^2 + \sigma_Y^2$.

3 Central limit theorem by the coupling method.

First a preliminary technical lemma:

Lemma 1 *Let Z be a Gaussian random variable with mean zero and variance σ^2 . Then there exists a constant L , independent of σ, s or ϵ such that*

$$E_{|Z|>\epsilon s}(Z^2) \leq L \frac{\sigma^3}{\epsilon s}. \quad (7)$$

We use the notation $E_A(Y)$ to mean the expected value of Y over the set A , i.e. $E(\mathbf{1}_A(Y)Y)$. For example the left hand side of the inequality in the lemma

$$2 \frac{1}{\sigma\sqrt{2\pi}} \int_{\epsilon s}^{\infty} z^2 e^{-\frac{z^2}{2\sigma^2}} dz.$$

2. Prove the lemma. [Hint: Make the change of variables $z = \sigma x$. Let $a = \frac{\epsilon s}{\sigma}$. Show that $a \int_a^{\infty} x^2 e^{-x^2/2} dx$ is bounded.]

Let $X_1, X_2, X_3, \dots, X_n, \dots$ be independent random variables with means zero and with variances σ_i^2 . Let

$$S_n = X_1 + \dots + X_n$$

so that

$$s_n^2 \stackrel{def}{=} \text{Var} S_n = \sigma_1^2 + \dots + \sigma_n^2.$$

We wish to prove that under condition (12) to be stated below we have

$$\lim_{n \rightarrow \infty} E \left(f \left(\frac{1}{s_n} S_n \right) \right) \rightarrow E(f(N)) \quad (8)$$

for any three times differentiable function, f with bounded third derivatives.

The method of proof is to choose Gaussian random variables

$$Z_1, Z_2, \dots, Z_n, \dots$$

all independent of the X_i 's and each other, with

$$\text{Var}(Z_i) = \sigma_i^2$$

and make the successive substitutions

$$\begin{aligned} S_n = X_1 + \dots + X_{n-1} + X_n &\rightarrow X_1 + \dots + X_{n-1} + Z_n \\ &\rightarrow X_1 + \dots + Z_{n-1} + Z_n \\ &\vdots \\ &\rightarrow Z_1 + Z_2 + \dots + Z_{n-1} + Z_n \stackrel{def}{=} Z, \end{aligned}$$

and to estimate the difference in expectation at each stage of substitution. At the end of the substitutions, Z is a Gaussian with variance s_n^2 and hence $(1/s_n)Z$ is equal in law to the (unit) normal distribution, N .

The difference at the $(n - k)$ -th stage is

$$E\left(f\left(\frac{1}{s_n}(T_k + X_k)\right)\right) - E\left(f\left(\frac{1}{s_n}(T_k + Z_k)\right)\right)$$

where

$$T_k = X_1 + \cdots + X_{k-1} + Z_{k+1} + \cdots + Z_n$$

is the sum of all the random variables which are unchanged during this particular substitution. Now we certainly have

$$\left|E\left(f\left(\frac{1}{s_n}(T_k + X_k)\right)\right) - E\left(f\left(\frac{1}{s_n}(T_k + Z_k)\right)\right)\right| \leq \sup_u \left|E\left(f\left(u + \frac{1}{s_n}X_k\right)\right) - E\left(f\left(u + \frac{1}{s_n}Z_k\right)\right)\right|$$

since, if we set

$$h(u) = E\left(f\left(u + \frac{1}{s_n}X_k\right)\right) - E\left(f\left(u + \frac{1}{s_n}Z_k\right)\right)$$

we have

$$E\left(h\left(\frac{1}{s_n}T_k\right)\right) = E\left(f\left(\frac{1}{s_n}(T_k + X_k)\right)\right) - E\left(f\left(\frac{1}{s_n}(T_k + Z_k)\right)\right).$$

So it will be enough for us to get a bound on

$$\left|E\left(f\left(u + \frac{1}{s_n}X_k\right)\right) - E\left(f\left(u + \frac{1}{s_n}Z_k\right)\right)\right|,$$

valid for all u .

Consider the Taylor expansion, with remainder, of f about the point u :

$$f(u + y) = f(u) + yf'(u) + \frac{1}{2}f''(u)y^2 + g(u, y),$$

where

$$g(u, y) = f(u + y) - \left[f(u) + yf'(u) + \frac{1}{2}f''(u)y^2\right] \quad (9)$$

is of the form

$$g(u, y) = \frac{1}{3!}f'''(u^*)y^3 \quad (10)$$

where u^* is some point between u and $u + y$. Since $f(u)$ is a constant as far as y is concerned, taking its expectations either with respect to Z_k or X_k gives the same value, $f(u)$. The expectations $E(X_k)$ and $E(Z_k)$ both vanish, and the expectations

$$E(X_k^2) = E(Z_k^2).$$

So

$$E\left(f(u) + \frac{1}{s_n}X_k f'(u) + \frac{1}{2}f''(u)\left[\frac{1}{s_n}X_k\right]^2\right) = E\left(f(u) + \frac{1}{s_n}Z_k f'(u) + \frac{1}{2}f''(u)\left[\frac{1}{s_n}Z_k\right]^2\right)$$

and hence

$$E(f(u + \frac{1}{s_n}X_k) - E(f(u + \frac{1}{s_n}Z_k)) = E(g(u, \frac{1}{s_n}X_k) - E(g(u, \frac{1}{s_n}Z_k))$$

so

$$|E(f(u + \frac{1}{s_n}X_k)) - E(f(u + \frac{1}{s_n}Z_k))| \leq E(|g(u, \frac{1}{s_n}(X_k))|) + E(|g(u, \frac{1}{s_n}Z_k)|)$$

and we shall estimate each of the terms on the right separately. From (10) we have

$$|g(u, y)| \leq K_1 y^3, \quad K_1 = \frac{1}{3!} \sup_u |f'''(u)|$$

valid for all y while for large $|y|$ the definition (9) implies that

$$|g(u, y)| \leq K_2 |y|^2$$

where K_2 can be expressed in terms of $\sup |f|$, $\sup |f'|$, and $\sup |f''|$. So there is a constant K such that

$$|g(u, y)| \leq K \min(|y|^2, |y|^3)$$

for all u and y . Of course we want to use the $|y|^3$ estimate for small $|y|$ and the $|y|^2$ estimate for large $|y|$.

In any event, given $\epsilon > 0$ we have

$$\begin{aligned} E(|g(u, \frac{1}{s_n}(X_k))|) &= E_{|X_k| \leq \epsilon s_n}(|g(u, \frac{1}{s_n}X_k)|) + E_{|X_k| > \epsilon s_n}(|g(u, \frac{1}{s_n}X_k)|) \\ &\leq K E_{|X_k| \leq \epsilon s_n}(|X_k|^3/s_n^3) + K E_{|X_k| > \epsilon s_n}(X_k^2/s_n^2) \\ &\leq K \epsilon E_{|X_k| \leq \epsilon s_n}(|X_k|^2/s_n^2) + K E_{|X_k| > \epsilon s_n}(X_k^2/s_n^2) \\ &\leq K \frac{\epsilon}{s_n^2} E(|X_k|^2) + K E_{|X_k| > \epsilon s_n}(X_k^2/s_n^2) \\ &= K \epsilon \frac{\sigma_k^2}{s_n^2} + \frac{K}{s_n^2} E_{|X_k| > \epsilon s_n}(X_k^2). \end{aligned}$$

We can make a similar estimate with Z_k instead of X_k to obtain

$$\begin{aligned} E(|g(u, \frac{1}{s_n}Z_k)|) &\leq K \epsilon \frac{\sigma_k^2}{s_n^2} + \frac{K}{s_n^2} E_{|Z_k| > \epsilon s_n}(Z_k^2) \\ &\leq K \epsilon \frac{\sigma_k^2}{s_n^2} + \left(\frac{KL}{\epsilon}\right) \left(\frac{\sigma_k^3}{s_n^3}\right) \end{aligned}$$

where, in passing from the first line to the second we made use of (7). We now add all these inequalities over k , using the facts that

$$\frac{1}{s_n}(Z_1 + \cdots + Z_n) \stackrel{law}{=} N$$

and

$$\sigma_1^2 + \dots + \dots \sigma_n^2 = s_n^2$$

to obtain

$$|E(f(\frac{1}{s_n}S_n)) - E(f(N))| \leq 2K\epsilon + \frac{K}{s_n^2} \sum_k E_{|X_k| > \epsilon s_n}(X_k^2) + \frac{KL}{\epsilon} \sum \frac{\sigma_k^3}{s_n^3}. \quad (11)$$

We can now state the central limit theorem:

Theorem 1 *Suppose that for any $\epsilon > 0$*

$$\lim_{n \rightarrow \infty} \left(\frac{1}{s_n^2} \right) \sum_k E_{|X_k| > \epsilon s_n}(X_k^2) = 0. \quad (12)$$

Then for any three times differentiable function, f , which is bounded together with its first three derivatives in absolute value we have

$$E\left(f\left(\frac{1}{s_n}S_n\right)\right) \rightarrow E(f(N)). \quad (13)$$

where N is the normally distributed random variable.

Proof. We are assuming that the second term in (11) goes to zero with n for any $\epsilon > 0$. So the question is estimating the third term. The sum giving the third term is maximized by $\max_k \sigma_k/s_n$ since the sum with squares instead of cubes adds up to 1.

3. Show that (12) implies that $\max_k \sigma_k/s_n \rightarrow 0$. [Hint: Break the expression $E\left(\frac{X_k^2}{s_n^2}\right)$ up into two parts, one $E_{|X_k| \leq \delta s_n}$ and the other $E_{|X_k| > \delta s_n}$. Show that by choosing $\delta < \epsilon^2$ the right hand side of (11) can be estimated by some multiple of ϵ .]

A loose interpretation of the condition (12) is that it is a little stronger than requiring that no one variance outweighs the others in the sense that

$$\frac{\sigma_k}{s_n} \rightarrow 0.$$

4. Suppose that all the X_k s are identically distributed. This means that they are all equal in law to a common X with variance, say $\sigma^2 > 0$. In this case

$$s_n^2 = n\sigma^2.$$

Show that (12) is satisfied in this case.