

Common distributions, MGFs and entropy

Sofia Olhede



September 23, 2020

1 Important Distributions

2 Entropy

Moment Generating Functions III

n → number of trials p → prob. of each trial S → prob. of taking n trials with success

- The mean, variance and moment generating function of

$$X \sim \text{Bin}(n, p) \text{ are given by } f(x; p) = \binom{n}{x} p^x (1-p)^{n-x}$$

$$E(X) = E(X_1 + \dots + X_n) = E[X_1] + \dots + E[X_n] = np$$

$$\mathbb{E}(X) = np, \quad \text{Var}(X) = np(1-p), \quad M_X(t) = (1 - p + pe^t)^n.$$

$$\text{Var}(X) = \text{Var}[X_1 + \dots + X_n] = \text{Var}[X_1] + \dots + \text{Var}[X_n] = n p(1-p)$$

$$\text{Mgf} = \mathbb{E}[e^{tX}] = \sum_x e^{tx} f(x; p)$$

- If $X = \sum_{i=1}^n Y_i$ where $Y_i \stackrel{\text{iid}}{\sim} \text{Bern}(p)$ then $X \sim \text{Bin}(n, p)$.
- A random variable X is said to follow the Geometric distribution with parameter $p \in (0, 1)$ denoted $X \sim \text{Geom}(p)$, if *number of failed trials until first Bernoulli success!*
 - $\mathcal{X} = \{0\} \cup \mathbb{N}$.
 - $f(x; p) = (1-p)^x p$.

- The mean, variance and moment generating of $X \sim \text{Geom}(p)$ are given by $E[X] = \sum_{x=0}^{\infty} x (1-p)^x p = p \sum_{x=1}^{\infty} (1-p)^{x-1} = p(1-p) \left[\frac{d}{dp} - \sum (1-p)^x \right] = p(1-p) \frac{d}{dp} \left(\frac{1}{1-p} \right) = \frac{p}{p-1}$

$$\mathbb{E}(X) = \frac{1-p}{p}, \quad \text{Var}(X) = \frac{1-p}{p^2}, \quad M_X(t) = \frac{p}{1-(1-p)e^t},$$

- the latter for $t < -\log(1-p)$.

NOTE: the moment generating function can be used to compute a distribution's moments: the n th moment about 0 is the n th derivative of the moment-generating function, evaluated at 0.

Moment Generating Functions IV

- Let $\{Y_i\}_{i \geq 1}$ be an infinite collection of random variables, where

$$Y_i \stackrel{\text{iid}}{\sim} \text{Bern}(p). \quad \text{Let } T = \min\{k \in \mathbb{N} : Y_k = 1\} - 1$$

Then $T \sim \text{Geom}(p)$. 
 ie a geometric random variable is equivalent to the number of steps required to get $Y=1$ (success), minus 1.

- A random variable X is said to follow the Negative Binomial distribution with parameter $p \in (0, 1)$ and $r > 0$, denoted $X \sim \text{NegBin}(r, p)$ if

$$\begin{aligned} * \quad & \mathcal{X} = \{0\} \cup \mathbb{N}. \\ * \quad & f(x; p) = \binom{x+r-1}{x} (1-p)^x p^r. \end{aligned}$$

The geometric distribution describes the probability of "x trials are made before a success", and the negative binomial distribution describes that of "x trials are made before r successes are obtained"

Thus, a geometric distribution can be thought of a negative binomial distribution with parameter r=1

Moment Generating Functions V

- The mean, variance and moment generating function of $X \sim \text{NegBin}(r, p)$ are given by

$$\mathbb{E}(X) = r \frac{1-p}{p}, \quad \text{Var}(X) = r \frac{1-p}{p^2}, \quad M_X(t) = \frac{p^r}{(1 - (1-p)e^t)^r},$$

- the latter for $t < -\log(1-p)$.
- If $X = \sum_{i=1}^r Y_i$ where $Y_i \sim \text{Geom}(p)$ then $X \sim \text{NegBin}(r, p)$.
- A random variable X is said to follow a Poisson distribution with parameter $\lambda > 0$ denoted $X \sim \text{Poisson}(\lambda)$ if

- * $\mathcal{X} = \{0\} \cup \mathbb{N}$.
- * $f(x; \lambda) = e^{-\lambda} \frac{\lambda^x}{x!}$.

sum of geometrics give a negative binomial;
also
sum of bernoullis give a binomial

Poisson is a discrete probability distribution that expresses the probability of a given number of events occurring in a fixed interval of time or space if these events occur with a known constant mean rate and independently of the time since the last event

Moment Generating Functions VI

- The mean, variance and moment generating function of $X \sim \text{Poisson}(\lambda)$ are given by

(solution
on the notes)

$$\mathbb{E}(X) = \lambda, \quad \text{Var}(X) = \lambda, \quad M_X(t) = \exp\{\lambda(e^t - 1)\}.$$

- Let $\{X_n\}_{n \geq 1}$ be a sequence of $\text{Binom}(n, p_n)$ random variables such that $p_n = \lambda/n$ for some constant $\lambda > 0$. Then $f_{X_n} \xrightarrow{n \rightarrow \infty} f_Y$ where $Y \sim \text{Poisson}(\lambda)$. *this is a long time going between Binomial and Poisson dist.*
- Let $X \sim \text{Poisson}(\lambda)$ and let $Y \sim \text{Poisson}(\mu)$ be independent. The conditional distribution of X given $X + Y = k$ is $\text{Binom}(k, \frac{\lambda}{\lambda + \mu})$. *as the Poisson generates the Binomial*
- A random vector \mathbf{X} in \mathbb{R}^k is said to follow the Multinomial distribution with parameters $n \in \mathbb{N}$ and $p = (p_1, \dots, p_k) \in (0, 1)^k$, such that $\sum_{i=1}^k p_i = 1$, denoted $\mathbf{X} \sim \text{Multi}(n, p_1, \dots, p_k)$ if

similar to binomial but instead of having 2 classes (success/insuccess, or coin), we have k classes (ie a 6-side die) and we can now see the outcome of all classes. i.e in a die we could count the number of 1s, 2s, 3s ... for $p_i = 1/6$ ans $\sum_i p_i = 1$.

also a generalization of the binomial

Moment Generating Functions VII

- Take

$$\begin{aligned} * \quad & \mathcal{X} = \{0, 1, \dots, n\}^k, \text{ and} \\ * \quad & f(x_1, \dots, x_k; n, \{p_i\}_{i=1}^k) = \\ & \frac{n!}{x_1! \dots x_k!} p_1^{x_1} \dots p_k^{x_k} I\left\{\sum_{i=1}^k x_i = n\right\}. \end{aligned}$$

- The mean, variance, covariance and moment generating function are:

$$\mathbb{E}(X_i) = np_i, \quad \mathbb{V}\text{ar}(X_i) = np_i(1 - p_i), \quad \mathbb{C}\text{ov}(X_i, X_j) = -np_i p_j,$$

$$M_{\mathbf{X}}(\mathbf{u}) = \left(\sum_{i=1}^k p_i e^{u_i} \right)^n.$$

- The multinomial generalizes the binomial distribution: n independent trials, with k possible outcomes.

Moment Generating Functions VIII

- Lemma (Poisson and Multinomial)

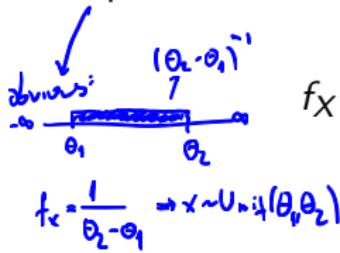
→ If $X_i \sim \text{Poisson}(\lambda_i)$, $i = 1, \dots, k$ are independent, then the conditional distribution of $\mathbb{X} = (X_1, \dots, X_k)^T$ given $\sum_{i=1}^k X_i = n$ is $\text{Multi}(n; p_1, \dots, p_k)$ with

$$p_i = \frac{\lambda_i}{\lambda_1 + \dots + \lambda_k}.$$

discrete prob.

continuous prob.

- A random variable X is said to follow the uniform distribution with parameters $-\infty < \theta_1 < \theta_2 < \infty$ denoted $X \sim \text{Unif}(\theta_1, \theta_2)$ if



Moment Generating Functions IX

- The mean, variance and moment generating function of $X \sim \text{Unif}(\theta_1, \theta_2)$ are given by

calculations on the notes

$$\mathbb{E}(X) = \frac{\theta_1 + \theta_2}{2}, \quad \text{Var}(X) = \frac{(\theta_2 - \theta_1)^2}{12}, \quad M_X(t) = \frac{e^{t\theta_2} - e^{t\theta_1}}{t(\theta_2 - \theta_1)},$$

$$\mathbb{E}(X) = \int_{\theta_1}^{\theta_2} x \cdot \frac{1}{\theta_2 - \theta_1} dx = \frac{1}{\theta_2 - \theta_1} \left[x^2 \right]_{\theta_1}^{\theta_2} = \frac{1}{\theta_2 - \theta_1} \cdot \frac{\theta_2^2 - \theta_1^2}{2} = \frac{\theta_1 + \theta_2}{2}$$

t ≠ 0. We also specify M(0) = 1.

- A random variable X is said to follow the exponential distribution with parameter $\lambda > 0$ denoted $X \sim \text{Expl}(\lambda)$ if

$$f_X(x; \lambda) = \begin{cases} \lambda e^{-\lambda x} & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}.$$

Moment Generating Functions X

- The mean, variance and moment generating function of $X \sim \text{Exp}(\lambda)$ are given by Integration by parts: $\int x \frac{du}{dx} dx = uv - \int u \frac{dv}{dx} dx$ or $\int ab = a \int b - \int a' \int b$

calculation in notes

$$\mathbb{E}(X) = \int_0^\infty x \lambda e^{-\lambda x} dx = [x(-e^{-\lambda x})]_0^\infty + \int_0^\infty \lambda x e^{-\lambda x} dx = 0 - 0 + \left[-\frac{1}{\lambda} e^{-\lambda x} \right]_0^\infty = -\frac{1}{\lambda} (e^{-\lambda \infty} - e^0) = 1/\lambda$$

- If X and Y are independent exponential random variables with rates λ_1 and λ_2 , then $Z = \min(X, Y)$ are also exponential with the rate $\lambda_1 + \lambda_2$. proof in notes
- Lack of memory characterisation.

* Let $X \sim \text{Exp}(\lambda)$. Then

$$\Pr(X \geq x + t | X \geq t) = \Pr(X \geq x).$$

* Conversely: if X is a random variable such that

$$\Pr(X \geq 0) > 0 \text{ and}$$

$$\Pr(X > x + s | X > t) = \Pr(X > s), \quad \forall t, s \geq 0,$$

then there exists a $\lambda > 0$ such that $X \sim \text{Exp}(\lambda)$.

Moment Generating Functions XI

the exponential dist is a special case of the gamma distribution
(when $\alpha=1$ and $\beta=\lambda$)

- A random variable X has a gamma distribution with parameters α and β (the shape and rate of the distribution respectively), written as $X \sim \text{Gamma}(\alpha, \beta)$ if

$$f_X(x; \alpha, \beta) = \begin{cases} \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x} & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$$

- The mean, variance and moment generating function of $X \sim \text{Gamma}(\alpha, \beta)$ are given by

?

$$\mathbb{E}(X) = \frac{\alpha}{\beta}, \quad \text{Var}(X) = \frac{\alpha}{\beta^2}, \quad M_X(t) = \left(\frac{\beta}{\beta - t} \right)^\alpha, \quad t < \beta.$$

Moment Generating Functions XII

- If $Y_1, \dots, Y_\alpha \stackrel{\text{iid}}{\sim} \text{Exp}(\beta)$ then $Y = Y_1 + \dots + Y_\alpha \sim \text{Gamma}(\alpha, \beta)$
→ (also see the Erlang distribution) → a sum of exponential dists. is a gamma dist
- The special case of $X \sim \text{Gamma}(\frac{k}{2}, \frac{1}{2})$ is the chi-square distribution on k degrees of freedom written as χ_k^2 .
- A random variable X follows the normal distribution with parameters $\mu \in \mathbb{R}$ and $\sigma^2 > 0$ denoted $X \sim N(\mu, \sigma^2)$ if

$$f_X(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right), \quad x \in \mathbb{R}.$$

- The mean, variance and moment generating function of $X \sim N(\mu, \sigma^2)$ are given by

$$\mathbb{E}(X) = \mu, \quad \text{Var}(X) = \sigma^2, \quad M_X(t) = \exp(t\mu + t^2\sigma^2/2).$$

↳ variance is the square of the std dev.

- In the case of $Z \sim N(0, 1)$ (standard normal density) we use

$$f_Z(z) = \varphi(z) \text{ and } F_Z(z) = \Phi(z). \quad \begin{cases} \varphi: \text{density of std normal density} \\ \Phi: \text{distribution of std normal} \end{cases}$$

Moment Generating Functions XIII

- Lemma: Let $X \sim N(\mu, \sigma^2)$ and assume $a \neq 0$. Then $aX + b \sim N(a\mu + b, a^2\sigma^2)$. Furthermore

$$F_X(x) = \Phi\left(\frac{x - \mu}{\sigma}\right),$$

where $\Phi()$ is the standard normal CDF or

$$\Phi(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}u^2\right) du.$$

- Corollary: let X_1, \dots, X_n be independent random variables and let

 $X_i \sim N(\mu_i, \sigma_i^2)$. Take S_n as the sum of the X_i . Then

$$\begin{aligned}
 M(t) &= \exp(t\mu + t^2\sigma^2/2) \\
 M_S &= \prod_i M_i(t) = \prod_i \exp(t\mu_i + t^2\sigma_i^2/2) \\
 &= \exp\left(t\sum_i \mu_i + t^2\sum_i \sigma_i^2/2\right) \\
 \Rightarrow S_n &\sim N\left(\sum_i \mu_i, \sum_i \sigma_i^2\right)
 \end{aligned}$$

If sum of Gaussians is a Gaussian

Moment Generating Functions XIII

- First note that the MGF of a Gaussian is

$$M_X(t) = \int_{\mathbb{R}} \frac{e^{tu}}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(u-\mu)^2\right) du \quad (1)$$

$$= \exp(\mu t + \sigma^2 t^2/2). \quad (2)$$

We assume X_1, \dots, X_n are independent $X_i \sim N(\mu_i, \sigma_i^2)$, and so for $Y = \sum_i X_i$ is

$$\begin{aligned} M_Y(t) &= \mathbb{E}_Y\{e^{tY}\} \\ &= \mathbb{E}_{X_1, \dots, X_n}\{e^{t\sum X_i}\} = \prod_i \mathbb{E}\{e^{tX_i}\} \\ &= \prod_i \exp(\mu_i t + \sigma_i^2 t^2/2) \\ &= \exp\left(\sum_i \mu_i t + \sum_i \sigma_i^2 t^2/2\right). \end{aligned}$$

Entropy etc

the more the entropy, the least predictable a random variable is!

- The entropy is used to measure the disorder of a random variable.
- The entropy of a random variable X is defined as

$$\rightarrow H(X) = -\mathbb{E}\{\log f_X(X)\}$$

$$= \begin{cases} -\sum_{x \in \mathcal{X}} f_X(x) \log\{f_X(x)\} & \text{if } X \text{ discrete} \\ -\int_{x \in \mathcal{X}} f_X(x) \log\{f_X(x)\} dx & \text{if } X \text{ continuous} \end{cases}$$

- The entropy is a measure of intrinsic disorder or unpredictability of a random system. The entropy can be thought of as a measure of the uncertainty of the random variable X . One can think of this as the missing information: the larger the entropy the less we know about X .
- It is related to the variance, but is not equivalent to the variance.

Entropy etc

- It can be shown that when X is a discrete random variable then
 - * $H(X) \geq 0$.
 - * $H(g(X)) \leq H(X)$ for any deterministic function g .
- Entropy is expressed in the unit bits. If log is replaced by \lg then the unit is nats.
 \downarrow
base 10
- Can we then use entropy to compare distributions?

Entropy etc

- Let $p(x)$ and $q(x)$ be two probability density (probability mass) functions on \mathbb{R} . We define the Kullback-Leibler divergence or relative entropy of q with respect to p as



$$\text{KL}(q||p) \equiv \int_{\mathbb{R}} p(x) \log \left(\frac{p(x)}{q(x)} \right) dx. \quad (3)$$

- By Jensen's inequality for $X \sim p(\cdot)$ we have

$$\text{KL}(q||p) = \mathbb{E}_p \left\{ -\log \frac{q(X)}{p(X)} \right\} \geq -\log \mathbb{E}_p \left\{ \frac{q(X)}{p(X)} \right\} = 0, \quad (4)$$

i.e. KL is an expectation

as q is unit norm.

- $p = q \Leftrightarrow \text{KL}(q||p) = 0$.
- The KL divergence is not a metric as it both lacks symmetry and violates the triangle inequality, though symmetrized versions exist.