

Useful theorems and the sketch of proof

Saki Kuzushima

2019-05-09

Contents

Sum of exponential random variables is a gamma random variable	1
Linear combination of Gaussian random variables is also Gaussian	2
Distribution of maximum values	2
The Dirichlet distribution belongs to the exponential family	2
The expectation of the sufficient statistic of the exponential family is the derivative of the log normalizing factor	3
KL divergence	4
Mutual Information	4
F1 score	5
Reference	5

Sum of exponential random variables is a gamma random variable

Theorem: If $X_i \stackrel{i.i.d}{\sim} Expo(\lambda)$, $i = 1 \dots n$, then $Y = \sum_{i=1}^n X_i \sim Gamma(n, \lambda)$.

Proof:

First, recall that if $X \sim Expo(\lambda)$,

$$M_X(t) = \frac{\lambda}{\lambda - t}, \quad (t < \lambda)$$

if $Z \sim Gamma(\alpha, \beta)$,

$$M_Z(t) = \left(\frac{1}{1 - \beta t}\right)^{-\alpha}$$

Observe that

$$\begin{aligned} M_Y(t) &= \mathbb{E}[e^{tY}] \\ &= \mathbb{E}[e^{t \sum_{i=1}^n X_i}] \\ &= \mathbb{E}\left[\prod_{i=1}^n e^{tX_i}\right] \\ &= \prod_{i=1}^n \mathbb{E}[e^{tX_i}] \\ &= \prod_{i=1}^n \frac{\lambda}{\lambda - t} \\ &= \left(\frac{\lambda}{\lambda - t}\right)^n \\ &= \left(\frac{\lambda - t}{\lambda}\right)^{-n} \\ &= \left(1 - \frac{t}{\lambda}\right)^{-n} \end{aligned}$$

This is an mgf of Gamma distribution. So, $Y \sim \text{Gamma}(n, \lambda)$.

Linear combination of Gaussian random variables is also Gaussian

Theorem: if $X \sim N(\mu, \sigma^2)$, and $Y = aX + b$, ($a \neq 0$), then $Y \sim N(a\mu + b, a^2\sigma^2)$.

Proof: Recall that if $X \sim N(\mu, \sigma^2)$

$$M_X(t) = e^{\mu t + \frac{\sigma^2}{2} t^2}$$

Then,

$$\begin{aligned} M_Y(t) &= \mathbb{E}[e^{t(aX+b)}] \\ &= \mathbb{E}[e^{taX+tb}] \\ &= e^{tb} \mathbb{E}[e^{taX}] \\ &= e^{tb} M_X(ta) \\ &= e^{tb} e^{\mu ta + \frac{\sigma^2}{2} t^2 a^2} \\ &= \exp(tb + t\mu a + \frac{\sigma^2}{2} t^2 a^2) \\ &= \exp((a\mu + b)t + \frac{a^2\sigma^2}{2} t^2) \end{aligned}$$

This is also an mgf of Gaussian. $Y \sim N(a\mu + b, a^2\sigma^2)$.

Distribution of maximum values

Theorem: If $X_1 \dots X_n \stackrel{i.i.d}{\sim} f_X(x)$, and $Y \equiv \max\{X_1 \dots X_n\}$, then $f_Y(y) = nF_X^{n-1}(y)f_X(y)$.

Proof:

$$\begin{aligned} F_Y(y) &= P(Y \leq y) \\ &= P(\max\{X_1 \dots X_n\} \leq y) \\ &= P(X_1 \leq y) \dots P(X_n \leq y) \\ &= F_X(y) \dots F_X(y) \\ &= F_X^n(y) \end{aligned}$$

Then, because $f_Y(y) = \frac{\partial}{\partial y} F_Y(y)$,

$$f_Y(y) = nF_X^{n-1}(y)f_X(y)$$

The Dirichlet distribution belongs to the exponential family

The fact that the Dirichlet distribution belongs to the exponential family is useful. While the expression of the Dirichlet in terms of the exponential family is widely available, the derivation is not. The following sketches the conversion from the usual expression of the pdf of Dirichlet distribution to the form in terms of the exponential family.

The pdf of the Dirichlet distribution is

$$p(\theta|\alpha) = \frac{1}{B(\alpha)} \prod_{i=1}^K \theta_i^{\alpha_i-1} \quad (1)$$

where $\theta = [\theta_1 \dots \theta_K]^\top$, and $\sum_{i=1}^K \theta_i = 1$, $\theta_i \geq 0, \forall i \in \{1, 2, \dots, K\}$.

The exponential family has the following general pdf

$$p(\theta|\eta) = h(\theta) \exp(\eta^\top t(\theta) - A(\eta)) \quad (2)$$

where $t(\theta)$ is the sufficient statistic, η is called the natural parameter, $A(\eta)$ is the log normalization factor, and $h(x)$ is the base measure.

It is known that the Dirichlet distribution belongs to the exponential family, and the parameters are

- The natural parameter: $\eta = \alpha = [\alpha_1 \dots \alpha_K]^\top$
- The sufficient statistic: $t(\theta) = \log \theta = \log[\theta_1, \dots, \theta_K]^\top$.
- The log normalization factor: $A(\eta) = \sum_{i=1}^K \log \Gamma(\eta_i) - \log \Gamma(\sum_{i=1}^K \eta_i)$
- The base measure: $h(x) = \frac{1}{\prod_{i=1}^K \theta_i}$

The following sketches why the parameters are expressed in these ways.

First, recall that $B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$. Then the pdf of the Dirichlet distribution is

$$\begin{aligned} p(\theta|\alpha) &= \frac{1}{B(\alpha)} \prod_{i=1}^K \theta_i^{\alpha_i-1} \\ &= \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\prod_{i=1}^K \Gamma(\alpha_i)} \cdot \frac{\prod_{i=1}^K \theta_i^{\alpha_i}}{\prod_{i=1}^K \theta_i} \end{aligned}$$

Taking the exponential and log,

$$\begin{aligned} p(\theta|\alpha) &= \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\prod_{i=1}^K \Gamma(\alpha_i)} \cdot \frac{\prod_{i=1}^K \theta_i^{\alpha_i}}{\prod_{i=1}^K \theta_i} \\ &= \frac{1}{\prod_{i=1}^K \theta_i} \exp \left[\log \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\prod_{i=1}^K \Gamma(\alpha_i)} \prod_{i=1}^K \theta_i^{\alpha_i} \right] \\ &= \frac{1}{\prod_{i=1}^K \theta_i} \exp \left[\log \Gamma(\sum_{i=1}^K \alpha_i) - \log \prod_{i=1}^K \Gamma(\alpha_i) + \log \prod_{i=1}^K \theta_i^{\alpha_i} \right] \\ &= \frac{1}{\prod_{i=1}^K \theta_i} \exp \left[\log \Gamma(\sum_{i=1}^K \alpha_i) - \sum_{i=1}^K \log \Gamma(\alpha_i) + \sum_{i=1}^K \alpha_i \log \theta_i \right] \\ &= \frac{1}{\prod_{i=1}^K \theta_i} \exp \left[\alpha^\top \log \theta - \left\{ \sum_{i=1}^K \log \Gamma(\alpha_i) - \log \Gamma(\sum_{i=1}^K \alpha_i) \right\} \right] \end{aligned}$$

Comparison of this expression with the general pdf of the exponential family yields the parameters listed above.

The expectation of the sufficient statistic of the exponential family is the derivative of the log normalizing factor

Note that this is useful in the derivation of the original variational inference for LDA (in Blei et al 2003).

Theorem:

$$\mathbb{E}_\theta[t(\theta)] = \frac{\partial}{\partial \eta} A(\eta) \quad (3)$$

Proof:

Using the fact that the pdf must integrate to one, we first express $a(\eta)$ by the other parameters.

$$\begin{aligned} 1 &= \int_{\theta} p(\theta|\eta) d\theta = \int_{\theta} h(\theta) \exp(\eta t(\theta) - A(\eta)) d\theta \\ \iff 1 &= \frac{1}{\exp A(\eta)} \int_{\theta} h(\theta) \exp(\eta t(\theta)) d\theta \\ \iff \exp A(\eta) &= \int_{\theta} h(\theta) \exp(\eta t(\theta)) d\theta \\ \iff A(\eta) &= \log \int_{\theta} h(\theta) \exp(\eta t(\theta)) d\theta \end{aligned}$$

Let $g(\eta) = \frac{1}{\exp A(\eta)}$. Then this is also equivalent to

$$1 = g(\eta) \int_{\theta} h(\theta) \exp(\eta t(\theta)) d\theta \quad (4)$$

Taking the derivative of Eq (3) with respect to η ,

$$\begin{aligned} 0 &= \frac{\partial}{\partial \eta} \left[g(\eta) \int_{\theta} h(\theta) \exp(\eta t(\theta)) d\theta \right] \\ &= g'(\eta) \int_{\theta} h(\theta) \exp(\eta t(\theta)) d\theta + g(\eta) \frac{\partial}{\partial \eta} \int_{\theta} h(\theta) \exp(\eta t(\theta)) d\theta \\ &= \frac{g'(\eta)}{g(\eta)} + g(\eta) \int_{\theta} h(\theta) \exp(\eta t(\theta)) t(\theta) d\theta \\ &= \frac{g'(\eta)}{g(\eta)} + \int_{\theta} g(\eta) h(\theta) \exp(\eta t(\theta)) t(\theta) d\theta \\ &= \frac{g'(\eta)}{g(\eta)} + \mathbb{E}_\theta[t(\theta)] \end{aligned}$$

i.e.

$$\mathbb{E}_\theta[t(\theta)] = -\frac{g'(\eta)}{g(\eta)} = -\frac{\partial}{\partial \eta} \log(g(\eta)) \quad (5)$$

because $\frac{\partial}{\partial \eta} \log(g(\eta)) = \frac{1}{g(\eta)} \cdot \frac{\partial}{\partial \eta} g(\eta)$.

Because $g(\eta) = \frac{1}{\exp(A(\eta))} = -\exp(A(\eta))$, Eq (4) is equivalent to

$$\mathbb{E}[t(\theta)] = -\frac{\partial}{\partial \eta} \log(g(\eta)) = \frac{\partial}{\partial \eta} A(\eta) \quad (6)$$

KL divergence

Mutual Information

The mutual information measures the dependence of two random variables, and defined as follows.

Let X and Y be two (discrete) random variables, then the mutual information is

$$I(X;Y) = D_{KL}(P_{XY}||P_X \otimes P_Y) = \sum_y \sum_x P_{XY}(x,y) \log \left(\frac{P_X Y(x,y)}{P_X(x)P_Y(y)} \right) \quad (7)$$

where D_{KL} is KL divergence. Likewise the normalized mutual information (NMI) is

$$\frac{I(X,Y)}{H(X)H(Y)} \quad (8)$$

where $H(\cdot)$ is entropy.

F1 score

One of the measure to evaluate the performance of classificaiton. F1 score is defined as the harmonic mean of the precision and recall, where precision = $P(truepositive|truepositive + falsepositive)$ and recall = $P(truepositive|truepositive + falsenegative)$. Harmonic mean of $x_1...x_n$ is

$$H = \left\{ \frac{\sum_{i=1}^n x_i^{-1}}{n} \right\}^{-1}$$

Thus, F1 score is

$$F1 = \left[\frac{recall^{-1} + precision^{-1}}{2} \right]^{-1}$$

Reference

Wikipedia on mutla information