

Distribution-hard

zwq

2022.10.22

本文介绍比较复杂但是在应用中非常重要的5个分布: Γ 分布, Beta分布, χ^2 分布, t 分布, F 分布。另外介绍了一下Dirichlet分布. 本文收集了其概率密度函数(probability density function, pdf)的证明(在大多数的工科教科书中并没有).

此外, 没有依赖现有的概率密度函数计算库, 仅根据每个分布pdf的数学定义, 分别绘制了这些分布的概率密度曲线.

预备知识

Γ 函数和Beta函数

由Euler总结的2个著名的反常积分:

$$\begin{aligned}\Gamma(\alpha) &= \int_0^{+\infty} x^{\alpha-1} e^{-x} dx \\ \text{Beta}(\alpha, \beta) &= \int_0^1 x^{\alpha-1} (1-x)^{\beta-1} dx\end{aligned}\tag{1}$$

这两个反常积分有很多变体(换元积分法). 例如Beta函数可以方便地计算三角函数的积分(类似于Wallis公式)

$$B\left(\frac{m+1}{2}, \frac{n+1}{2}\right) = 2 \int_0^{\frac{\pi}{2}} \cos^m \theta \cdot \sin^n \theta d\theta$$

此外, 还有一些基本结论:

$$\begin{aligned}\Gamma(\alpha+1) &= \alpha\Gamma(\alpha), \alpha \in \mathbb{R}^+ \\ \Gamma(n+1) &= n!, \Gamma\left(n+\frac{1}{2}\right) = \frac{\sqrt{\pi}}{2^{2n}} \frac{\Gamma(2n+1)}{\Gamma(n+1)}, n \in \mathbb{N}^+ \\ \Gamma(1) &= 1, \Gamma\left(\frac{1}{2}\right) = \sqrt{\pi} \\ \text{Beta}(\alpha, \beta) &= \text{Beta}(\beta, \alpha) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)} \\ \text{Beta}\left(\frac{1}{2}, \frac{1}{2}\right) &= \pi, \text{Beta}(1, 1) = \frac{1}{2}\end{aligned}$$

1 Γ 分布

$$\begin{aligned}\frac{\Gamma(\alpha)}{\Gamma(\alpha)} &= \frac{1}{\Gamma(\alpha)} \int_0^{+\infty} t^{\alpha-1} e^{-t} dt \\ &\stackrel{t=\beta x}{=} \int_0^{+\infty} \frac{1}{\Gamma(\alpha)} \beta^\alpha x^{\alpha-1} e^{-\beta x} dx \\ &= 1\end{aligned}\tag{2}$$

那么就得到了 Γ 分布的概率密度函数:

$$f(x; \alpha, \beta) = \begin{cases} \frac{1}{\Gamma(\alpha)} \beta^\alpha x^{\alpha-1} e^{-\beta x}, & x > 0 \\ 0, & \text{otherwise} \end{cases} \quad (3)$$

其中两个参数 $\alpha, \beta > 0$, α 为形状参数, β 为尺度(逆尺度)参数.

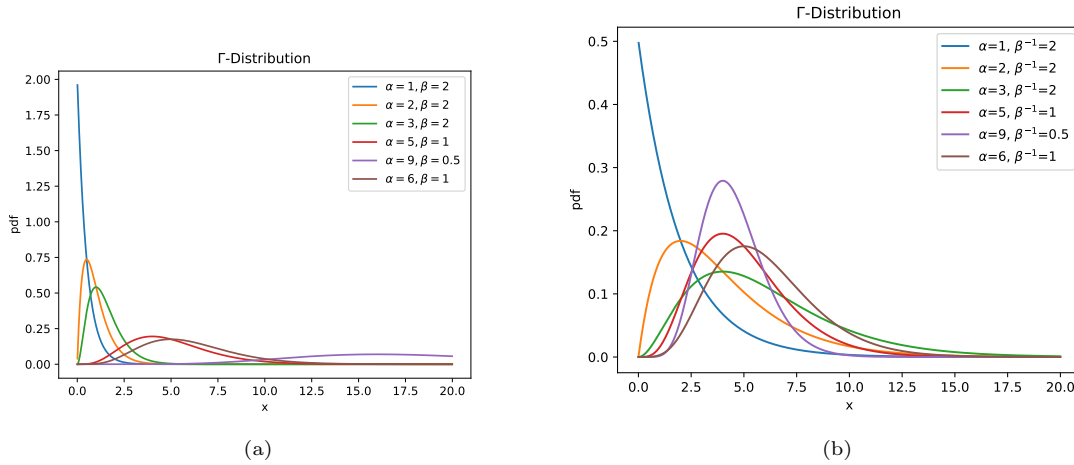
就是说, 也可以写成如下的形式(这种形式也更常用, 下图(b)就是基于这个pdf绘制的):

$$f(x; \alpha, \beta) = \begin{cases} \frac{1}{\Gamma(\alpha) \beta^\alpha} x^{\alpha-1} e^{-\frac{1}{\beta} x}, & x > 0 \\ 0, & \text{otherwise} \end{cases} \quad (4)$$

当 $\alpha < 1$, $f(x; \alpha, \beta)$ 为递减函数;

当 $\alpha = 1$, $f(x; \alpha, \beta)$ 为递减函数;

当 $\alpha > 1$, $f(x; \alpha, \beta)$ 为单峰函数;



Gamma 分布

在绘制Gamma分布的pdf曲线时, 有一个非常有趣现象: 如果按照 $f(x; \alpha, \beta)$ 绘图, 得到的结果不好看Fig. 1a, 但是如果用 $\lambda = \beta^{-1}$ 代替Gamma分布的 β 时, 可得到 $f(x; \alpha, \beta^{-1})$ 的曲线更美观, 也更常用. Fig. 1b. 这可能是 β 被称为为尺度(逆尺度)参数的原因.

2 Beta分布

n-Bernouli试验, 每一次事件发生的概率为 p 且互相独立, 那么这个试验 X 服从n重伯努利分布: $X \sim B(n, p)$

$$Pr(x = k) = \binom{n}{k} p^k (1-p)^{n-k}$$

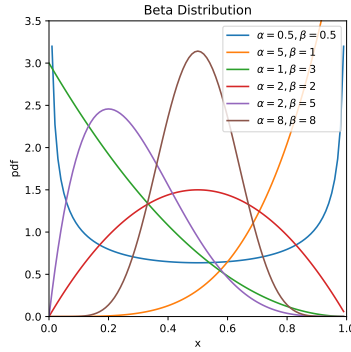
又因为在实验之前, 没有任何先验知识, 我们不知道 p 的值是什么。因此假设 p 服从均匀分布 $p \sim U(0, 1)$, 那么可得 p 的概率密度函数 $p(x) \equiv 1$. 那么 X 的概率累计函数为: 当 $x \leq 0$ ($x \geq 1$)时, $F_X(x) = 0$ (1). 当 $0 < x < 1$ 时, 有

$$F_X(x) = Pr(p \leq x | k, n) = \frac{Pr(k, n, p \leq x)}{Pr(k, n)}$$

$$\begin{aligned}
Pr(k, n) &= \int_0^1 Pr(k, n, p = x) p(x) dx \\
&= \int_0^1 \binom{n}{k} x^k (1-x)^{n-k} dx \\
&= C, \text{ is a constant.}
\end{aligned} \tag{5}$$

因此

$$\begin{aligned}
F_X(x) &= \frac{\int_0^x \binom{n}{k} t^k (1-t)^{n-k} dt}{Pr(k, n)} \\
f_X(x) &= F'_X(x) \\
&= \left(\frac{\int_0^x \binom{n}{k} t^k (1-t)^{n-k} dt}{Pr(k, n)} \right)' \\
&= \frac{\binom{n}{k} x^k (1-x)^{n-k}}{Pr(k, n)} \\
&= \frac{x^k (1-x)^{n-k}}{\int_0^1 u^k (1-u)^{n-k} du}
\end{aligned} \tag{6}$$



Beta 分布

令 $k = \alpha - 1, n - k = \beta - 1$, 那么就得到Beta分布的概率密度函数

$$f(x; \alpha, \beta) = \begin{cases} \frac{1}{\text{Beta}(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}, & 0 < x < 1 \\ 0, & \text{otherwise} \end{cases} \tag{7}$$

因为 $\text{Beta}(\alpha, \beta) = \text{Beta}(\beta, \alpha) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$, 因此可以替换为

$$f(x; \alpha, \beta) = \begin{cases} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1}, & 0 < x < 1 \\ 0, & \text{otherwise} \end{cases}$$

Dirichlet分布

Beta分布是基于二项分布Binomial (n重伯努利分布), 即: n次实验, 每次的结果只有两种(0,1; 发生或不发生; 掷硬币).

推而广之, 多项分布Multinomial, 即多次实验, 每次的结果只有n种情况(掷骰子6种情况).

类比于Beta分布，可得到Dirichlet分布的pdf大致如下：

$$f(x; \alpha) = \begin{cases} \frac{\Gamma(\sum_{i=1}^n \alpha_i)}{\prod_{i=1}^n \Gamma(\alpha_i)} \prod_{i=1}^n x_i^{\alpha_i-1}, & 0 < x_i < 1, \sum_{i=1}^n x_i = 1 \\ 0, & \text{otherwise} \end{cases} \quad (8)$$

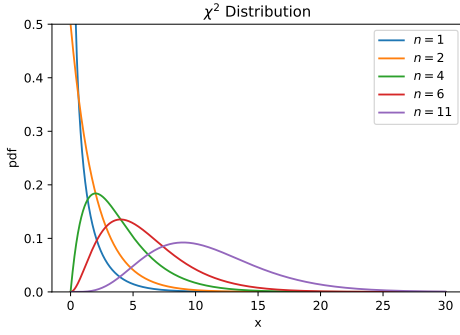
多维空间，不好绘制图像

3 χ^2 分布

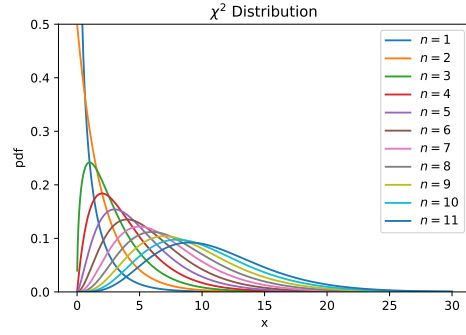
假设随机变量 X_1, X_2, \dots, X_n 独立同分布于标准正态分布 $N(0, 1)$.

那么随机变量 $X = \sum_{i=1}^n X_i^2$ 满足 $X \sim \chi^2(n)$.其概率密度函数为:

$$f(x; n) = \begin{cases} \frac{1}{2^{\frac{n}{2}} \Gamma(\frac{n}{2})} x^{\frac{n}{2}-1} e^{-\frac{x}{2}}, & x > 0 \\ 0, & \text{otherwise} \end{cases} \quad (9)$$



(a)



(b)

χ^2 分布

证明:

$$Y = X_1^2 + X_2^2 + \dots + X_n^2$$

$$\begin{aligned} F_Y(y) &= Pr(Y \leq y) = Pr(X_1^2 + X_2^2 + \dots + X_n^2 \leq y) \\ &= \int \dots \int_{x_1^2 + x_2^2 + \dots + x_n^2 \leq y} \frac{1}{(2\pi)^{\frac{n}{2}}} e^{-\frac{x_1^2 + x_2^2 + \dots + x_n^2}{2}} dx_1 dx_2 \dots dx_n \\ &= \int_0^{\sqrt{y}} \int_{\theta_1=\psi_1}^{\varphi_1} \dots \int_{\theta_{n-1}=\psi_{n-1}}^{\varphi_{n-1}} \frac{1}{(2\pi)^{\frac{n}{2}}} e^{-\frac{r^2}{2}} r^{n-1} dr d\theta_1 \dots d\theta_{n-1} \\ &= C_n \int_0^{\sqrt{y}} e^{-\frac{r^2}{2}} r^{n-1} dr \end{aligned} \quad (10)$$

求导得概率密度函数:

$$\begin{aligned} f_Y(y) &= (F_Y(y))' \\ &= \frac{\partial(C_n \int_0^{\sqrt{y}} e^{-\frac{r^2}{2}} r^{n-1} dr)}{\partial y} \\ &= C_n e^{-\frac{y}{2}} y^{\frac{n-1}{2}} (\sqrt{y})' \\ &= \frac{C_n}{2} e^{-\frac{y}{2}} y^{\frac{n}{2}-1} \end{aligned} \quad (11)$$

又因为概率密度积分为1,可得

$$\begin{aligned}
& \int_0^{+\infty} f_Y(y) dy \\
&= \frac{C_n}{2} \int_0^{+\infty} e^{-\frac{y}{2}} y^{\frac{n}{2}-1} dy \\
&\stackrel{x=\frac{y}{2}}{=} \frac{C_n}{2} \int_0^{+\infty} e^{-x} (2x)^{\frac{n}{2}-1} 2dx \\
&= C_n 2^{\frac{n}{2}-1} \int_0^{+\infty} e^{-x} x^{\frac{n}{2}-1} dx \\
&= \frac{C_n}{2} 2^{\frac{n}{2}} \Gamma\left(\frac{n}{2}\right) = 1
\end{aligned} \tag{12}$$

即 $\frac{C_n}{2} = \frac{1}{2^{\frac{n}{2}} \Gamma(\frac{n}{2})}$

代入上式可得 χ^2 分布的pdf:

$$f_Y(y) = \frac{1}{2^{\frac{n}{2}} \Gamma(\frac{n}{2})} e^{-\frac{y}{2}} y^{\frac{n}{2}-1}, \quad y > 0$$

引理 1 二维连续型随机变量商的分布 二维随机变量 (X, Y) 的联合概率密度为 $p(x, y)$. 那么 $Z = \frac{X}{Y}$ 的概率密度 $f_Z(z) = \int_{-\infty}^{+\infty} p(zy, y)|y|dy$.

证明:

$$\begin{aligned}
F_Z(z) &= Pr(Z \leq z) = Pr\left(\frac{X}{Y} \leq z\right) \\
&= \int_{-\infty}^0 dy \int_{yz}^{+\infty} p(x, y) dx + \int_0^{+\infty} dy \int_{-\infty}^{yz} p(x, y) dx \\
&\stackrel{x=uy}{=} \int_{-\infty}^0 dy \int_z^{+\infty} p(uy, y) y du + \int_0^{+\infty} dy \int_{-\infty}^z p(uy, y) y du \\
&= \int_z^{+\infty} du \int_{-\infty}^0 p(uy, y) y dy + \int_{-\infty}^z du \int_0^{+\infty} p(uy, y) y dy
\end{aligned} \tag{13}$$

求导得:

$$\begin{aligned}
f_Z(z) &= (F_Z(z))' \\
&= - \int_{-\infty}^0 p(zy, y) y dy + \int_0^{+\infty} p(zy, y) y dy \\
&= \int_{-\infty}^{+\infty} p(zy, y) |y| dy
\end{aligned} \tag{14}$$

证毕.

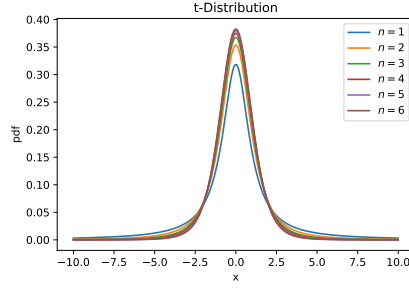
4 t分布

$X \sim N(0, 1)$, $Y \sim \chi^2(n)$, 那么随机变量 $Z = \frac{X}{\sqrt{Y/n}}$ 服从自由度为n的t分布, 即: $Z \sim t(n)$

$$f(t; n) = \frac{\Gamma(\frac{n+1}{2})}{\sqrt{n\pi} \Gamma(\frac{n}{2})} \left(1 + \frac{t^2}{n}\right)^{-\frac{n+1}{2}}, \quad -\infty < t < \infty \tag{15}$$

因为 $X \sim N(0, 1)$, $Y \sim \chi^2(n)$. 记随机变量 $W = \sqrt{\frac{Y}{n}}$, 显然W恒大于零. 当 $w < 0$ 时, $f_W(w) = 0$. 当 $w > 0$ 时, 因为Y服从卡方分布, 由卡方分布的pdf有:

$$\begin{aligned}
F_W(w) &= Pr(W \leq w) = Pr(Y \leq nw^2) \\
&= \int_0^{nw^2} \frac{1}{2^{\frac{n}{2}} \Gamma(\frac{n}{2})} x^{\frac{n}{2}-1} e^{-\frac{x}{2}} dx
\end{aligned} \tag{16}$$



t分布

求导得 $W = \sqrt{\frac{Y}{n}}$ 的pdf:(在第5节-F分布中,还会求得 $W = \frac{Y}{n}$ 的pdf)

$$\begin{aligned}
 f_W(w; n) &= (F_W(w))' \\
 &= \frac{1}{2^{\frac{n}{2}} \Gamma(\frac{n}{2})} (nw^2)^{\frac{n}{2}-1} e^{-\frac{nw^2}{2}} 2nw \\
 &= \frac{1}{\Gamma(\frac{n}{2})} 2^{1-\frac{n}{2}} n^{\frac{n}{2}} w^{n-1} e^{-\frac{nw^2}{2}}
 \end{aligned} \tag{17}$$

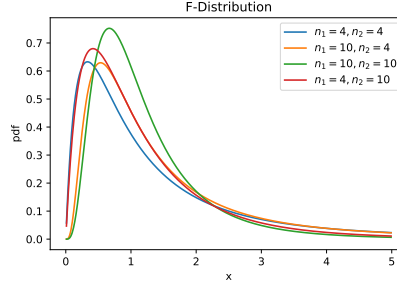
$Z = \frac{X}{\sqrt{Y/n}}$, 可记为 $T = \frac{X}{W}$. 已知X的pdf(标准正态分布)和W的pdf, 又因为X和W相互独立, 那么 $(X, W) \sim p(x, w) = p(x) \cdot f_W(w)$. 那么由引理 1 可求出两个随机变量商的分布:

$$\begin{aligned}
 f_T(t) &= \int_{-\infty}^{+\infty} p(tw, w) |w| dw = \int_{-\infty}^{+\infty} p_X(tw) p_W(w) |w| dw \\
 &= \int_0^{+\infty} p_X(tw) p_W(w) w dw \\
 &= \int_0^{+\infty} \left(\frac{1}{\sqrt{2\pi}} e^{-\frac{t^2 w^2}{2}} \right) \left(\frac{1}{\Gamma(\frac{n}{2})} 2^{1-\frac{n}{2}} n^{\frac{n}{2}} w^{n-1} e^{-\frac{nw^2}{2}} \right) w dw \\
 &= \left(\frac{1}{\sqrt{\pi} \Gamma(\frac{n}{2})} 2^{\frac{1-n}{2}} n^{\frac{n}{2}} \right) \int_0^{+\infty} e^{-\frac{(n+t^2)}{2} w^2} w^n dw \\
 &\stackrel{w=\sqrt{\frac{2z}{n+t^2}}}{=} \left(\frac{1}{\sqrt{\pi} \Gamma(\frac{n}{2})} 2^{\frac{1-n}{2}} n^{\frac{n}{2}} \right) \int_0^{+\infty} e^{-z} \frac{1}{n+t^2} \left(\frac{2}{n+t^2} \right)^{\frac{n-1}{2}} z^{\frac{n-1}{2}} dz \\
 &= \left(\frac{1}{\sqrt{\pi} \Gamma(\frac{n}{2})} 2^{\frac{1-n}{2}} n^{\frac{n}{2}} \right) \frac{1}{2} \left(\frac{2}{n+t^2} \right)^{\frac{n+1}{2}} \int_0^{+\infty} e^{-z} z^{\frac{n-1}{2}} dz \\
 &= \left(\frac{1}{\sqrt{\pi} \Gamma(\frac{n}{2})} 2^{\frac{1-n}{2}} n^{\frac{n}{2}} \right) \frac{1}{2} \left(\frac{2}{n+t^2} \right)^{\frac{n+1}{2}} \Gamma\left(\frac{n+1}{2}\right) \\
 &= \frac{\Gamma(\frac{n+1}{2})}{\sqrt{\pi} \Gamma(\frac{n}{2})} n^{\frac{n}{2}} \left(\frac{1}{\sqrt{n}} \right) \left(\frac{1}{n+t^2} \right)^{\frac{n+1}{2}} \\
 &= \frac{\Gamma(\frac{n+1}{2})}{\sqrt{n\pi} \Gamma(\frac{n}{2})} \left(1 + \frac{t^2}{n} \right)^{-\frac{n+1}{2}}
 \end{aligned} \tag{18}$$

5 F分布

$X \sim \chi^2(n_1)$, $Y \sim \chi^2(n_2)$, 那么随机变量 $Z = \frac{X/n_1}{Y/n_2}$ 服从自由度为 n_1, n_2 的F分布, 即: $Z \sim F(n_1, n_2)$

$$f(t; n_1, n_2) = \begin{cases} \frac{\Gamma(\frac{n_1+n_2}{2})}{\Gamma(\frac{n_1}{2})\Gamma(\frac{n_2}{2})} n_1^{\frac{n_1}{2}-1} n_2^{\frac{n_2}{2}-1} t^{\frac{n_1}{2}-1} (n_1 t + n_2)^{-\frac{n_1+n_2}{2}}, & t > 0 \\ 0, & \text{otherwise} \end{cases} \quad (19)$$



F分布

因为 $X \sim \chi^2(n_1)$, $Y \sim \chi^2(n_2)$, 记新的随机变量 $X' = X/n_1$, $Y' = Y/n_2$, 那么有 $T = \frac{X/n_1}{Y/n_2} = \frac{X'}{Y'}$, 由卡方分布的pdf, 可得 X, Y 的pdf为:

$$f_X(x) = \frac{1}{2^{\frac{n_1}{2}} \Gamma(\frac{n_1}{2})} e^{-\frac{x}{2}} x^{\frac{n_1}{2}-1}, \quad x > 0$$

$$f_Y(y) = \frac{1}{2^{\frac{n_2}{2}} \Gamma(\frac{n_2}{2})} e^{-\frac{y}{2}} y^{\frac{n_2}{2}-1}, \quad y > 0$$

可得 X', Y' 的pdf为:

$$f_{X'}(x) = \frac{n_1^{\frac{n_1}{2}}}{2^{\frac{n_1}{2}} \Gamma(\frac{n_1}{2})} e^{-\frac{n_1 x}{2}} x^{\frac{n_1}{2}-1}, \quad x > 0$$

$$f_{Y'}(y) = \frac{n_2^{\frac{n_2}{2}}}{2^{\frac{n_2}{2}} \Gamma(\frac{n_2}{2})} e^{-\frac{n_2 y}{2}} y^{\frac{n_2}{2}-1}, \quad y > 0$$

证明

$X \sim \chi^2(n)$, 则 $W = \frac{X}{n}$ 的pdf为(证明如下):

$$\begin{aligned} F_W(w) &= Pr(W \leq w) = Pr(X \leq nw) \\ &= \int_0^{nw} \frac{1}{2^{\frac{n}{2}} \Gamma(\frac{n}{2})} x^{\frac{n}{2}-1} e^{-\frac{x}{2}} dx \end{aligned} \quad (20)$$

求导得 $W = \frac{X}{n}$ 的pdf:

$$\begin{aligned} f_W(w; n) &= (F_W(w))' \\ &= \frac{1}{2^{\frac{n}{2}} \Gamma(\frac{n}{2})} (nw)^{\frac{n}{2}-1} e^{-\frac{nw}{2}} n \\ &= \frac{1}{2^{\frac{n}{2}} \Gamma(\frac{n}{2})} n^{\frac{n}{2}} w^{\frac{n}{2}-1} e^{-\frac{nw}{2}} \end{aligned} \quad (21)$$

证毕

记随机变量 $T = \frac{X'}{Y'}$, 显然 T 恒大于零. 当 $t \leq 0$ 时, $f_T(t) = 0$. 当 $t > 0$ 时, 有

$$F_T(t) = Pr(T \leq t) = Pr\left(\frac{X'}{Y'} \leq t\right) = Pr(X' \leq Y't)$$

由于X,Y相互独立, 因此 X', Y' 也相互独立,根据引理 1有:

$$\begin{aligned}
f_T(t) &= \int_{-\infty}^{+\infty} p(ty, y)|y|dy = \int_0^{+\infty} p_{X'}(ty)p_{Y'}(y)ydy \\
&= \int_0^{+\infty} \left(\frac{n_1^{\frac{n_1}{2}}}{2^{\frac{n_1}{2}}\Gamma(\frac{n_1}{2})} e^{-\frac{n_1 ty}{2}} (ty)^{\frac{n_1}{2}-1} \right) \left(\frac{n_2^{\frac{n_2}{2}}}{2^{\frac{n_2}{2}}\Gamma(\frac{n_2}{2})} e^{-\frac{n_2 y}{2}} y^{\frac{n_2}{2}-1} \right) ydy \\
&= \left(\frac{n_1^{\frac{n_1}{2}} n_2^{\frac{n_2}{2}}}{2^{\frac{n_1+n_2}{2}}\Gamma(\frac{n_1}{2})\Gamma(\frac{n_2}{2})} \right) t^{\frac{n_1}{2}-1} \int_0^{+\infty} e^{-\frac{(n_1 t + n_2)y}{2}} y^{\frac{n_1+n_2}{2}-1} ydy \\
&\stackrel{y=\frac{2z}{n_1 t + n_2}}{=} \left(\frac{n_1^{\frac{n_1}{2}} n_2^{\frac{n_2}{2}}}{2^{\frac{n_1+n_2}{2}}\Gamma(\frac{n_1}{2})\Gamma(\frac{n_2}{2})} \right) \int_0^{+\infty} e^{-z} \left(\frac{2}{n_1 t + n_2} \right)^{\frac{n_1+n_2}{2}} z^{\frac{n_1+n_2}{2}-1} t^{\frac{n_1}{2}-1} dz \\
&= \left(\frac{n_1^{\frac{n_1}{2}} n_2^{\frac{n_2}{2}}}{2^{\frac{n_1+n_2}{2}}\Gamma(\frac{n_1}{2})\Gamma(\frac{n_2}{2})} \right) \left(\frac{2}{n_1 t + n_2} \right)^{\frac{n_1+n_2}{2}} t^{\frac{n_1}{2}-1} \int_0^{+\infty} e^{-z} z^{\frac{n_1+n_2}{2}-1} dz \\
&= \left(\frac{n_1^{\frac{n_1}{2}} n_2^{\frac{n_2}{2}}}{2^{\frac{n_1+n_2}{2}}\Gamma(\frac{n_1}{2})\Gamma(\frac{n_2}{2})} \right) \left(\frac{2}{n_1 t + n_2} \right)^{\frac{n_1+n_2}{2}} t^{\frac{n_1}{2}-1} \Gamma\left(\frac{n_1+n_2}{2}\right) \\
&= \frac{\Gamma(\frac{n_1+n_2}{2})}{\Gamma(\frac{n_1}{2})\Gamma(\frac{n_2}{2})} n_1^{\frac{n_1}{2}} n_2^{\frac{n_2}{2}} \left(\frac{1}{n_1 t + n_2} \right)^{\frac{n_1+n_2}{2}} t^{\frac{n_1}{2}-1} \\
&= \frac{\Gamma(\frac{n_1+n_2}{2})}{\Gamma(\frac{n_1}{2})\Gamma(\frac{n_2}{2})} n_1^{\frac{n_1}{2}} n_2^{\frac{n_2}{2}} t^{\frac{n_1}{2}-1} (n_1 t + n_2)^{-\frac{n_1+n_2}{2}}
\end{aligned} \tag{22}$$