

# The Gaussian Distribution from Scratch

Karl Stratos

## 1 Definitions

Let  $\mu \in \mathbb{R}^d$  and  $\Sigma \in \mathbb{R}_{>0}^{d \times d}$  where we assert  $\Sigma \succ 0$  (i.e., positive-definite) to avoid handling degenerate cases. We define the **Gaussian distribution**  $\mathcal{N}(\mu, \Sigma) : \mathbb{R}^d \rightarrow [0, 1]$  as

$$\mathcal{N}(\mu, \Sigma)(x) := \frac{1}{(\sqrt{2\pi})^d \sqrt{\det(\Sigma)}} \exp\left(-\frac{1}{2}(x - \mu)^\top \Sigma^{-1}(x - \mu)\right) > 0$$

which integrates to 1 over  $\mathbb{R}^d$  (Lemma D.4). We write  $X \sim \mathcal{N}(\mu, \Sigma)$  to mean that the probability of  $X = x$  is  $\mathcal{N}(\mu, \Sigma)(x)$ . The following statements are equivalent (Lemma D.12):

1.  $X \sim \mathcal{N}(\mu, \Sigma)$ .
2.  $M_X(t) := \mathbf{E}[\exp(t^\top X)] = \exp(t^\top \mu + \frac{1}{2}t^\top \Sigma t)$  for all  $t \in \mathbb{R}^d$  (moment-generating function).
3.  $X = \mu + \Sigma^{1/2}Z$  where  $Z \sim \mathcal{N}(0_d, I_{d \times d})$ .
4.  $a^\top X \sim \mathcal{N}(a^\top \mu, a^\top \Sigma a)$  for all nonzero  $a \in \mathbb{R}^d$ .

If any holds, we say  $X \in \mathbb{R}^d$  is **normally distributed** (or normal) with parameters  $(\mu, \Sigma)$ . These alternative definitions are useful. For instance, definition 3 is the reparameterization trick in GANs where we view the image  $X$  sampled from a Gaussian distribution as a differentiable perturbation of model parameters  $(\mu, \Sigma)$ . Definition 2 can be used to view a point-mass distribution on  $x \in \mathbb{R}^d$  as “normal” with parameters  $(x, 0_{d \times d})$  since its moment-generating function is  $\mathbf{E}[\exp(t^\top X)] = \exp(t^\top x)$ .

**Linear transformation.** A critical property of the Gaussian distribution is that it is closed under linear transformation. Note that definitions 3 and 4 are consistent with this property. For any  $A \in \mathbb{R}^{d' \times d}$  and  $b \in \mathbb{R}^{d'}$  where  $A$  is full-rank with  $d' \leq d$  (so that  $A\Sigma A^\top \succ 0$ ),  $X \sim \mathcal{N}(\mu, \Sigma)$  implies (Lemma C.2):

$$AX + b \sim \mathcal{N}(A\mu + b, A\Sigma A^\top) \tag{1}$$

## 2 Joint Distribution

We say  $X \in \mathbb{R}^d$  and  $Y \in \mathbb{R}^{d'}$  are **jointly normally distributed** (or jointly normal) with parameters  $(\mu, \Sigma)$  if

$$\begin{bmatrix} X \\ Y \end{bmatrix} \sim \mathcal{N}\left(\mu = \begin{bmatrix} \mu_X \\ \mu_Y \end{bmatrix}, \Sigma = \begin{bmatrix} \Sigma_X & \Sigma_{XY} \\ \Sigma_{YX} & \Sigma_Y \end{bmatrix}\right)$$

where  $\mu_X \in \mathbb{R}^d$ ,  $\mu_Y \in \mathbb{R}^{d'}$ ,  $\Sigma_X \in \mathbb{R}_{>0}^{d \times d}$ ,  $\Sigma_Y \in \mathbb{R}_{>0}^{d' \times d'}$ ,  $\Sigma_{XY} \in \mathbb{R}^{d \times d'}$ , and  $\Sigma_{YX} = \Sigma_{XY}^\top$ .<sup>1</sup> Jointly normal is generally not equivalent to individually normal. See Appendix E for an example in which  $X, Y \in \mathbb{R}$  are normal (moreover uncorrelated) but  $(X, Y) \in \mathbb{R}^2$  is not normal. However, if  $X$  and  $Y$  are *independently* normal, then they are jointly normal since we can write

$$\begin{bmatrix} X \\ Y \end{bmatrix} \sim \mathcal{N}\left(\mu = \begin{bmatrix} \mu_X \\ \mu_Y \end{bmatrix}, \Sigma = \begin{bmatrix} \Sigma_X & 0_{d \times d'} \\ 0_{d' \times d} & \Sigma_Y \end{bmatrix}\right)$$

<sup>1</sup> $\Sigma_X, \Sigma_Y \succ 0$  since they are main-diagonal blocks of  $\Sigma \succ 0$  (Lemma D.9) and  $\Sigma_{XY} = \Sigma_{YX}^\top$  since  $\Sigma$  is symmetric.

## 2.1 Linear Combinations

Let  $A \in \mathbb{R}^{p \times d}$ ,  $B \in \mathbb{R}^{p \times d'}$ , and  $b \in \mathbb{R}^p$  where  $A, B$  are full-rank with  $p \leq \min(d, d')$ . If  $X \in \mathbb{R}^d$  and  $Y \in \mathbb{R}^{d'}$  are jointly normal with parameters  $(\mu, \Sigma)$ , we have from (1) that

$$AX + BY + b \sim \mathcal{N}(A\mu_X + B\mu_Y + b, A\Sigma_X A^\top + A\Sigma_{XY} B^\top + B\Sigma_{YX} A^\top + B\Sigma_Y B^\top) \quad (2)$$

In particular, if  $X$  and  $Y$  are independently normal, then their sum is normal:

$$X + Y \sim \mathcal{N}(\mu_X + \mu_Y, \Sigma_X + \Sigma_Y)$$

Note that we need joint normality to guarantee the normality of a linear combination. In general a linear combination of normal variables may not be normal (e.g., (42)).

## 2.2 Conditional Distribution

If  $X \in \mathbb{R}^d$  and  $Y \in \mathbb{R}^{d'}$  are jointly normal with parameters  $(\mu, \Sigma)$ , and if  $\Sigma_Y - \Sigma_{YX}\Sigma_X^{-1}\Sigma_{XY}$  is invertible, then for all  $x \in \mathbb{R}^d$  and  $y \in \mathbb{R}^{d'}$  (Lemma D.10):

$$\begin{aligned} \mathcal{N}(\mu, \Sigma)((x, y)) &= \mathcal{N}(\mu_X, \Sigma_X)(x) \\ &\quad \times \mathcal{N}(\mu_Y + \Sigma_{YX}\Sigma_X^{-1}(x - \mu_X), \Sigma_Y - \Sigma_{YX}\Sigma_X^{-1}\Sigma_{XY})(y) \end{aligned}$$

Therefore the conditional distribution over  $Y|X = x$  is also Gaussian. In particular, the Gaussian distribution is conjugate to itself. This is used in applications like Kalman filtering.

## 3 Entropy

Let  $\mu' \in \mathbb{R}^d$  and  $\Sigma' \in \mathbb{R}_{>0}^{d \times d}$  be parameters of an additional Gaussian distribution over  $\mathbb{R}^d$ . Then (Lemma D.6):

$$H(\mathcal{N}(\mu', \Sigma'), \mathcal{N}(\mu, \Sigma)) = \frac{1}{2}(\mu' - \mu)^\top \Sigma^{-1}(\mu' - \mu) + \frac{1}{2} \text{tr}(\Sigma^{-1}\Sigma') + \frac{1}{2} \log((2\pi)^d \det(\Sigma))$$

It follows that

$$\begin{aligned} H(\mathcal{N}(\mu, \Sigma)) &= \frac{1}{2} \log((2\pi e)^d \det(\Sigma)) \\ D_{\text{KL}}(\mathcal{N}(\mu', \Sigma') || \mathcal{N}(\mu, \Sigma)) &= \frac{1}{2}(\mu' - \mu)^\top \Sigma^{-1}(\mu' - \mu) + \frac{1}{2} \text{tr}(\Sigma^{-1}\Sigma' - I_{d \times d}) + \frac{1}{2} \log\left(\frac{\det(\Sigma)}{\det(\Sigma')}\right) \end{aligned}$$

$\mathcal{N}(\mu, \Sigma)$  has the largest entropy among all distributions over  $\mathbb{R}^d$  with mean  $\mu$  and covariance  $\Sigma$  (Theorem B.1). This is mainly because it standardizes  $x$  inside the exponential function.

### 3.1 Mutual Information

Let  $X \in \mathbb{R}^d$  and  $Y \in \mathbb{R}^{d'}$  be jointly normal with parameters  $(\mu, \Sigma)$ . If  $\Sigma_Y - \Sigma_{YX}\Sigma_X^{-1}\Sigma_{XY}$  is invertible, then (Lemma D.11):

$$\begin{aligned} H(Y|X) &= \frac{1}{2} \log\left((2\pi e)^{d'} \det(\Sigma_Y - \Sigma_{YX}\Sigma_X^{-1}\Sigma_{XY})\right) \\ I(X, Y) &= \frac{1}{2} \log\left(\frac{\det(\Sigma_X) \det(\Sigma_Y)}{\det(\Sigma)}\right) \end{aligned} \quad (3)$$

Note that  $I(X, Y)$  is infinite if  $Y = X$ . By the [noisy-channel coding theorem](#), mutual information is the capacity (highest information rate that can be achieved nearly error-free) of a communication channel between  $X$  and  $Y$ . Below we give some well-known models with controllable mutual information.

**Additive white Gaussian noise channel.** Let  $X \sim \mathcal{N}(0, \sigma^2)$  and  $Z \sim \mathcal{N}(0, \nu^2)$  independently, and define  $Y = X + Z$ .  $X$  and  $Y$  are jointly normal because  $a_1 X + a_2 Y = (a_1 + a_2)X + a_2 Z$  is a sum of independently normal variables and thus normal for all nonzero  $a = (a_1, a_2)$  (definition 4). Since  $\text{Var}(Y) = \sigma^2 + \nu^2$  and  $\text{Cov}(X, Y) = \sigma^2$ ,

$$\begin{bmatrix} X \\ Y \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma^2 & \sigma^2 \\ \sigma^2 & \sigma^2 + \nu^2 \end{bmatrix}\right) \quad \Rightarrow \quad I(X, Y) = \frac{1}{2} \log\left(1 + \frac{\sigma^2}{\nu^2}\right)$$

Thus  $I(X, X + Z)$  grows logarithmically in signal-to-noise ratio  $\frac{\sigma^2}{\nu^2}$ .

**Correlated standard normal channel.** Let  $X, Y \in \mathbb{R}$  be jointly standard normal with correlation  $\rho < 1$ . One way to construct them is to let  $X, Z \stackrel{iid}{\sim} \mathcal{N}(0, 1)$  and set  $Y = \rho X + \sqrt{1 - \rho^2} Z$ . Then

$$\begin{bmatrix} X \\ Y \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix} \right) \quad \Rightarrow \quad I(X, Y) = -\frac{1}{2} \log(1 - \rho^2)$$

By taking the correlation  $\rho \rightarrow 1$  we can arbitrarily increase  $I(X, Y)$ .

## 4 Central Limit Theorem

Let  $\mathbf{Unk}(\mu, \sigma^2)$  denote an unknown distribution over  $\mathbb{R}$  with mean  $\mu$  and variance  $\sigma^2 > 0$ . It is often of interest to consider the sample average  $\bar{X}_N$  defined as

$$X_1 \dots X_N \stackrel{iid}{\sim} \mathbf{Unk}(\mu, \sigma^2) \quad \bar{X}_N := \frac{1}{N} \sum_{i=1}^N X_i$$

The average is itself random: every time we draw  $N$  iid samples from  $\mathbf{Unk}(\mu, \sigma^2)$ , we draw a single sample of  $\bar{X}_N$ . We can easily verify that  $\mathbf{E}[\bar{X}_N] = \mu$  and  $\text{Var}(\bar{X}_N) = \frac{\sigma^2}{N}$ , which states that  $\bar{X}_N$  concentrates around  $\mu$  as  $N \rightarrow \infty$  (this is called the “law of large numbers”). But what is the distribution of  $\bar{X}_N$ ? The **central limit theorem** (CLT) states that  $\bar{X}_N$  is asymptotically normal. More precisely, as  $N \rightarrow \infty$  we have

$$\sqrt{N}(\bar{X}_N - \mu) \stackrel{\text{approx.}}{\sim} \mathcal{N}(0, \sigma^2) \quad (4)$$

or, using the closure under linear transformation,

$$\bar{X}_N \stackrel{\text{approx.}}{\sim} \mathcal{N}\left(\mu, \frac{\sigma^2}{N}\right) \quad (5)$$

which is consistent with but not implied by the law of large numbers. CLT allows us to make probabilistic statements about sample averages regardless of the underlying distribution. For instance, if  $X_1 \dots X_N$  are arbitrary iid samples with mean 42 and variance 7, then approximately  $\bar{X}_N \sim \mathcal{N}(42, \frac{7}{N})$  so that we can calculate quantities like  $\Pr(\bar{X}_N \leq 50)$  (e.g., by consulting a standard normal table).

A proof of CLT shows that the KL divergence between the distribution of  $\sqrt{N}(\bar{X}_N - \mu)$  and  $\mathcal{N}(0, \sigma^2)$  goes to zero as  $N \rightarrow \infty$ . It is nontrivial: we refer to [Marsh \(2013\)](#) for details. CLT generalizes naturally to multivariate. If  $\mathbf{Unk}(\mu, \Sigma)$  is an unknown distribution over  $\mathbb{R}^d$  with mean  $\mu$  and covariance  $\Sigma \succ 0$ , then the average  $\bar{X}_N$  of samples  $X_1 \dots X_N \stackrel{iid}{\sim} \mathbf{Unk}(\mu, \Sigma)$  satisfies as  $N \rightarrow \infty$ :

$$\sqrt{N}(\bar{X}_N - \mu) \stackrel{\text{approx.}}{\sim} \mathcal{N}(0_d, \Sigma) \quad (6)$$

$$\bar{X}_N \stackrel{\text{approx.}}{\sim} \mathcal{N}\left(0_d, \frac{1}{N}\Sigma\right) \quad (7)$$

## References

Marsh, C. (2013). Introduction to continuous entropy. *Department of Computer Science, Princeton University*.

# A Integration for Dummies

## A.1 Single-Variable

An **antiderivative** of  $f : \mathbb{R} \rightarrow \mathbb{R}$  is a function  $F : \mathbb{R} \rightarrow \mathbb{R}$  such that  $F' = f$ . If  $F$  is an antiderivative, then so is  $F + C$  for any constant  $C \in \mathbb{R}$ . For instance,  $(1/3)x^3 + 42$  is an antiderivative of  $x^2$ .

The (definite) **integral** of  $f : \mathbb{R} \rightarrow \mathbb{R}$  over  $[a, b]$  is a scalar  $\int_a^b f(x)dx \in \mathbb{R}$  that represents the signed area of  $f$  on  $[a, b]$ . The quantity  $f(x)dx$  is interpreted as the product of the function value and an infinitesimally small interval. There are different ways to formalize the area. The most common definition is the Riemannn integral which partitions  $[a, b]$  into intervals  $[i\delta, (i + 1)\delta]$  of width  $\delta > 0$  and define

$$\int_a^b f(x)dx := \lim_{\delta \rightarrow 0} \sum_i f(x_i^\delta)\delta \quad (8)$$

where  $x_i^\delta \in [i\delta, (i + 1)\delta]$ . The finite sum  $\sum_i f(x_i^\delta)\delta$  for a given width  $\delta$  is called a **Riemann sum**. Thus an integral is simply the limiting value of a Riemann sum (if it exists it is unique). A more general definition is a Lebesgue integral which partitions the range of  $f$ .

The **fundamental theorem of calculus** (FTC) allows us to evaluate integrals by antiderivatives: if  $F$  is any antiderivative of  $f$ , then

$$\int_a^b f(x)dx = F(x)\Big|_a^b := F(b) - F(a) \quad (9)$$

For instance, the signed area under  $x^2$  over  $[-1, 1]$  is  $2/3$ . Basic properties of integration include

$$\int_a^b \alpha f(x) + \beta g(x)dx = \alpha \int_a^b f(x)dx + \beta \int_a^b g(x)dx \quad (\text{linearity}) \quad (10)$$

$$\int_a^b f(g(x))g'(x)dx = \int_{g(a)}^{g(b)} f(u)du \quad (u\text{-substitution}) \quad (11)$$

$$\int_a^b f(x)G(x)dx = F(x)G(x)\Big|_a^b - \int_a^b F(x)g(x)dx \quad (\text{integration by parts}) \quad (12)$$

(Exercise: verify (11–12) using the chain rule and the product rule in differentiation.)

### A.1.1 Substitution in practice

While (11) is the standard form of  $u$ -substitution, we often use it mechanically as follows. We wish to integrate  $f$  over the interval  $a < b$ . We view  $f$  as a (hopefully simpler) function of  $u = g(x)$  where  $g : \mathbb{R} \rightarrow \mathbb{R}$  is invertible and differentiable with nonzero derivative over  $(a, b)$ . The infinitesimals are related as  $du = g'(x)dx$  by the chain rule, or equivalently  $dx = g'(g^{-1}(u))^{-1}du$ . This yields a “plug-in” version of (11) where we substitute  $g(x) = u$  and  $dx = g'(g^{-1}(u))^{-1}du$ ,

$$\int_a^b f(g(x))dx = \int_{g(a)}^{g(b)} f(u)g'(g^{-1}(u))^{-1}du \quad (13)$$

For instance,

$$\begin{aligned} \int_0^{\sqrt{\frac{\pi}{2}}} 2x \cos(x^2) dx &= \int_0^{\frac{\pi}{2}} 2\sqrt{u} \cos(u) \left(\frac{1}{2\sqrt{u}}\right) du \\ &= \int_0^{\frac{\pi}{2}} \cos(u) du = \sin(u)\Big|_0^{\frac{\pi}{2}} = 1 \end{aligned}$$

where  $2x \cos(x^2) = 2\sqrt{u} \cos(u)$  with  $u = g(x) = x^2$ . Note that  $g$  is invertible on  $(0, \sqrt{\frac{\pi}{2}})$  so that  $x = \sqrt{u}$ ; it is also differentiable with nonzero derivative  $g'(x) = 2x$ . Writing  $dx = (2\sqrt{u})^{-1}du$ , we cancel terms and are finally able to use FTC (9).

**Orientation of region.** Observe that

$$1 = \int_0^1 1dx = \int_0^{-1} (-1)du = \int_{-1}^0 (+1)du$$

The first equality is by FTC. The second equality is by (13) with  $f(x) = 1$  and  $u = g(x) = -x$ . The final equality is again by FTC, simply acknowledging that  $(-x)|_0^{-1} = x|_{-1}^0 = 1$ . More generally, when  $g'(x)^{-1} < 0$  (i.e.,  $u$  is moving in the opposite direction of  $x$ ), we also change the “orientation of region” in integration (right-to-left instead of left-to-right). We can consider an alternative orientation-free formulation of  $u$ -substitution by always assuming integrating left-to-right. Let  $R$  denotes a region  $a < b$ , then

$$\int_R f(g(x))dx = \int_{g(R)} f(u) |g'(g^{-1}(u))^{-1}| du \quad (14)$$

where  $g(R)$  is the output region of  $g$  when applied to  $R$ , integrated from a smaller value to a larger value. This formulation is useful because it generalizes to higher dimensions (16).

## A.2 Multi-Variable

The integral of  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  over a region  $R \subseteq \mathbb{R}^d$  is a scalar  $\int_R f(x)dx \in \mathbb{R}$  that represents the signed hypervolume of  $f$  on  $R$ . Evaluation of such an integral is generally challenging because the region may take complicated forms (high-dimensional curves).

We can greatly simplify the problem by restricting the region to be a hypercube  $R = [a, b]$  where  $a, b \in \mathbb{R}^d$  specify a  $d$ -dimensional bounding box  $[a_1, b_1] \times \dots \times [a_d, b_d]$  (potentially all of  $\mathbb{R}^d$ ). A central tool in this setting is **Fubini’s theorem**, which states that

$$\int_{[a,b]} f(x)dx = \int_{a_{\pi(d)}}^{b_{\pi(d)}} \left( \dots \left( \int_{a_{\pi(1)}}^{b_{\pi(1)}} f(x_1 \dots x_d) dx_{\pi(1)} \right) \dots \right) dx_{\pi(d)}$$

where  $\pi$  is any permutation of  $\{1 \dots d\}$ . Thus we can evaluate a multi-variable integral by iteratively evaluating a single-variable integral in any order.

Many properties of integration carry over (like linearity), but some need to be generalized. One important generalization is **multi-variable  $u$ -substitution**. Let  $R \subseteq \mathbb{R}^d$  and  $g : R \rightarrow \mathbb{R}^d$  such that  $J_g(x) \in \mathbb{R}^{d \times d}$  (Jacobian of  $g$ ) is nonzero for all  $x \in R$ . Then

$$\int_R f(g(x)) |\det(J_g(x))| dx = \int_{g(R)} f(u) du \quad (15)$$

Similar to the single-variable case, we often use substitution mechanically as follows. We integrate  $f$  over a region  $R$  by viewing it as a simpler function of  $u = g(x)$  where  $g : R \rightarrow \mathbb{R}^d$  is assumed to be invertible (i.e.,  $\det(J_g(x)) \neq 0$ ). The infinitesimals are related as  $du = |\det(J_g(x))| dx$  or equivalently  $dx = |\det(J_g(x))|^{-1} du$ . This gives

$$\int_R f(g(x))dx = \int_{g(R)} f(u) |\det(J_g(g^{-1}(u)))|^{-1} du \quad (16)$$

where we “plug in”  $g(x) = u$  and  $dx = |\det(J_g(g^{-1}(u)))|^{-1} du$ . This strictly generalizes (14).

### A.2.1 Applications to probability

Let  $X \in \mathbb{R}^d$  be a random vector with distribution  $p_X$  supported on  $S \subseteq \mathbb{R}^d$  (i.e.,  $p_X(x) \geq 0$  and  $\int_S p_X(x)dx = 1$ ). The probability that  $X$  lies in a region  $R \subseteq S$  is

$$\Pr(X \in R) = \int_R p_X(x)dx$$

Let  $t : S \rightarrow T$  be a smooth invertible function where  $T \subseteq \mathbb{R}^d$ . Define a new random vector  $Y = t(X)$  supported on  $T$ . We claim that  $Y$  has the distribution

$$p_Y(y) = p_X(t^{-1}(y)) |\det(J_{t^{-1}}(y))| \quad \forall y \in T \quad (17)$$

Equivalently,

$$p_Y(t(x)) = p_X(x) |\det(J_{t^{-1}}(t(x)))| \quad \forall x \in S \quad (18)$$

*Proof sketch.* For any  $R \subseteq T$ ,

$$\Pr(Y \in R) = \Pr(X \in t^{-1}(R)) = \int_{t^{-1}(R)} p_X(x) dx = \int_R p_X(t^{-1}(y)) |\det(J_{t^{-1}}(y))| dy$$

where the last equality applies (15) with  $g = t^{-1}$ . This implies (17).

## B Continuous Entropy and KL Divergence

We generalize results in Marsh (2013) to multivariate. The continuous/differential entropy of  $X \in \mathbb{R}^d$  with density  $p_X$  supported on  $S \subseteq \mathbb{R}^d$  is defined as<sup>2</sup>

$$H(X) := - \int_S p_X(x) \log p_X(x) dx \quad (19)$$

It is easily seen that entropy is additive for independent variables. That is, if  $X \in \mathbb{R}^d$  and  $Y \in \mathbb{R}^{d'}$  are independent then the entropy of  $Z = (X, Y) \in \mathbb{R}^{d+d'}$  is  $H(Z) = H(X) + H(Y)$ .

- The uniform distribution  $u_{[a,b]}(x) := \frac{1}{b-a}$  over  $[a, b] \subset \mathbb{R}$  has entropy

$$H(X) = \int_a^b \frac{1}{b-a} \log(b-a) dx = \log(b-a) \quad (20)$$

- The Gaussian distribution  $\mathcal{N}(\mu, \Sigma)$  over  $\mathbb{R}^d$  has entropy (Corollary D.7)

$$H(X) = \frac{1}{2} \log((2\pi e)^d \det(\Sigma))$$

- The exponential distribution  $e_\lambda(x) := \lambda \exp(-\lambda x)$  over  $[0, \infty)$  with parameter  $\lambda > 0$  has entropy (Lemma D.5)

$$H(X) = 1 - \log \lambda \quad (21)$$

Unfortunately, continuous entropy suffers from various shortcomings (reviewed in Section B.1), most notably negativity (e.g., (20) is negative if  $b - a < 1$ , (21) is negative if  $\lambda > e$ ). On the other hand, let  $q_X$  be another density of  $X$  with support  $S$ . Define the continuous KL divergence (aka. relative entropy) between  $p_X$  and  $q_X$  as

$$D_{\text{KL}}(p_X || q_X) := \int_S p_X(x) \log \frac{p_X(x)}{q_X(x)} dx \quad (22)$$

Continuous KL divergence is nonnegative:

$$\begin{aligned} D_{\text{KL}}(p_X || q_X) &= \mathbf{E}_{x \sim p_X} \left[ \log \frac{p_X(x)}{q_X(x)} \right] \\ &= \mathbf{E}_{x \sim p_X} \left[ -\log \frac{q_X(x)}{p_X(x)} \right] \\ &\geq -\log \left( \mathbf{E}_{x \sim p_X} \left[ \frac{q_X(x)}{p_X(x)} \right] \right) \quad (\text{convexity of } -\log) \\ &= -\log \left( \int_S p_X(x) \frac{q_X(x)}{p_X(x)} dx \right) \\ &= -\log \left( \int_S q_X(x) dx \right) = 0 \end{aligned}$$

where  $D_{\text{KL}}(p_X || q_X) = 0$  iff  $p_X = q_X$  almost everywhere. This has useful implications.

<sup>2</sup>We use the term ‘‘density’’ in this section to distinguish continuous vs discrete variables.

- The cross entropy between  $p_X$  and  $q_X$  upper bounds the entropy of  $p_X$ ,

$$H(p_X, q_X) := H(p_X) + D_{\text{KL}}(p_X || q_X) \geq H(p_X) \quad (23)$$

- Mutual information is nonnegative,

$$I(X, Y) := D_{\text{KL}}(p_{XY} || p_X p_Y) \geq 0 \quad (24)$$

The cross entropy upper bound can be used to derive various maximum entropy densities.

**Theorem B.1.**

$$\mathcal{N}(\mu, \Sigma) \in \arg \max_{p_X: \mathbf{E}[X]=\mu, \text{Var}(X)=\Sigma} H(p_X) \quad (25)$$

$$u_{[a,b]} \in \arg \max_{p_X: \text{Support}(p_X)=[a,b]} H(p_X) \quad (26)$$

$$e_\lambda \in \arg \max_{p_X: \text{Support}(p_X)=\mathbb{R}_{\geq 0}^d, \mathbf{E}[X]=\lambda^{-1}} H(p_X) \quad (27)$$

where  $u_{[a,b]}$  denotes the uniform distribution over  $[a, b] \subset \mathbb{R}^d$  and  $e_\lambda$  denotes the product exponential density over  $\mathbb{R}_{\geq 0}^d$  with  $\lambda > 0_d$

*Proof.* (25): Let  $p_X$  with mean  $\mu \in \mathbb{R}^d$  and covariance  $\Sigma \succ 0$ . Then

$$\begin{aligned} H(p_X, \mathcal{N}(\mu, \Sigma)) &= \int_{\mathbb{R}^d} p_X(x) \left( \frac{1}{2} (x - \mu)^\top \Sigma^{-1} (x - \mu) + \frac{1}{2} \log((2\pi)^d \det(\Sigma)) \right) \\ &= \frac{1}{2} \mathbf{E}_{x \sim p_X} [(x - \mu)^\top \Sigma^{-1} (x - \mu)] + \frac{1}{2} \log((2\pi)^d \det(\Sigma)) \\ &= \frac{d}{2} + \frac{1}{2} \log((2\pi)^d \det(\Sigma)) \\ &= \frac{1}{2} \log((2\pi e)^d \det(\Sigma)) = H(\mathcal{N}(\mu, \Sigma)) \geq H(p_X) \end{aligned}$$

(26): Assume  $d = 1$ . Given any  $p_X$  with support  $[a, b]$  we have

$$H(p_X, u_X) = \int_a^b p_X(x) \log(b - a) dx = \log(b - a) = H(u_{[a,b]}) \geq H(p_X)$$

The statement holds for  $d > 1$  since each dimension is independently optimized.

(27): Assume  $d = 1$ . Given any  $p_X$  with support  $[0, \infty)$  and mean  $\lambda^{-1} > 0$  we have

$$H(p_X, e_\lambda) = \int_0^\infty p_X(x) (\lambda x - \log \lambda) dx = \lambda \mathbf{E}_{x \sim p_X} [x] - \log \lambda = 1 - \log \lambda = H(e_\lambda) \geq H(p_X, e_\lambda)$$

The statement holds for  $d > 1$  since each dimension is independently optimized. □

## B.1 Shortcomings of Continuous Entropy

### B.1.1 Inconsistency with Shannon entropy

The Shannon entropy of discrete  $X \in \{x_1 \dots x_n\}$  with distribution  $p_X$  is

$$H(X) := - \sum_{i=1}^n p_X(x_i) \log p_X(x_i) \quad (28)$$

This definition was [derived](#) by Shannon as a solution that satisfies axioms of information (regarding monotonicity, non-negativity, zero information, and independence). (19) appears to be a natural continuous extension of (28) in the sense that both are  $\mathbf{E}_{x \sim p_X} [-\log p_X(x)]$ , but it fails to satisfy the axioms (e.g., it can be negative). One way to

better understand why is to show that (19) is inconsistent with (28) in the limit. Assume  $d = 1$  and let  $p_X$  be a density supported on  $[a, b]$ . By definition

$$\int_a^b p_X(x) dx = \lim_{\delta \rightarrow 0} \sum_i p_X(x_i^\delta) \delta = 1 \quad (29)$$

where  $\sum_i p_X(x_i^\delta) \delta$  is a finite Riemann sum of width  $\delta > 0$ . Thus we can cast the density  $p_X$  as an increasingly fine-grained discrete distribution with probabilities  $p_X(x_i^\delta) \delta$  as  $\delta \rightarrow 0$ . Note that each value of  $\delta > 0$  yields a discrete distribution with a well-defined Shannon entropy. This Shannon entropy, in the limit, is

$$\begin{aligned} \lim_{\delta \rightarrow 0} \left( - \sum_i (p_X(x_i^\delta) \delta) \log(p_X(x_i^\delta) \delta) \right) &= - \lim_{\delta \rightarrow 0} \sum_i (p_X(x_i^\delta) \log p_X(x_i^\delta)) \delta - \lim_{\delta \rightarrow 0} \sum_i p_X(x_i^\delta) \delta \log \delta \\ &= - \int_a^b p_X(x) \log p_X(x) dx - \lim_{\delta \rightarrow 0} \sum_i p_X(x_i^\delta) \delta \log \delta \\ &= H(X) - \left( \lim_{\delta \rightarrow 0} \sum_i p_X(x_i^\delta) \delta \right) \left( \lim_{\delta \rightarrow 0} \log \delta \right) \end{aligned} \quad (30)$$

$$= H(X) + \infty \quad (31)$$

where (30) follows from the [generalized product rule of limits](#) using (29).<sup>3</sup> So the limiting Shannon entropy diverges from the continuous entropy by an infinite offset.

### B.1.2 Variability under change of coordinates

A good measure of information should not depend on the representation of samples from a distribution. For instance, let  $p_X$  be a distribution over finitely many circles, each of which can be specified by its radius or area. Clearly, the Shannon entropy of the circle is the same regardless of the representation. Now let  $p_X$  be a density over all circles. The continuous entropy of the circle under the radius representation is different from that under the area representation. A general statement that implies this result is given below.

**Lemma B.2.** Let  $X \in \mathbb{R}^d$  with density  $p_X$  supported on  $S$ . For any invertible mapping  $t$  on  $S$ ,

$$H(t(X)) = H(X) - \mathbf{E}_{x \sim p_X} [\log |\det(J_{t^{-1}}(t(x)))|]$$

*Proof.*

$$\begin{aligned} H(t(X)) &= - \int_S p_X(x) \log p_X(t(x)) dx \\ &= - \int_S p_X(x) \log p_X(x) dx - \int_S p_X(x) \log |\det(J_{t^{-1}}(t(x)))| dx \quad (\text{by (18)}) \\ &= H(X) - \mathbf{E}_{x \sim p_X} [\log |\det(J_{t^{-1}}(t(x)))|] \end{aligned}$$

□

**Corollary B.3.** For any invertible  $A \in \mathbb{R}^{d \times d}$  and  $b \in \mathbb{R}^d$ ,

$$H(AX + b) = H(X) - \log |\det(A^{-1})| \quad (32)$$

**Corollary B.4.** For  $\alpha > 0$ ,

$$H(\alpha X) = H(X) + d \log \alpha$$

<sup>3</sup>Assume  $\lim_{x \rightarrow a} f(x) \neq 0$ . If  $g(x)$  does not oscillate around  $a$ ,

$$\lim_{x \rightarrow a} f(x)g(x) = \lim_{x \rightarrow a} f(x) \lim_{y \rightarrow a} g(y)$$

If  $g(x)$  oscillates around  $a$ , then so does  $f(x)g(x)$ .



*Proof.*

$$\begin{aligned}
H(\alpha X) &= H(X) - \log |\det(\alpha^{-1} I_{d \times d})| && \text{(by (32))} \\
&= H(X) - \log |\alpha^{-d}| \\
&= H(X) - \log \alpha^{-d} && \text{(since } \alpha > 0) \\
&= H(X) + d \log \alpha
\end{aligned}$$

□

Corollary B.4 states that we can vacuously increase the continuous entropy of  $X \in \mathbb{R}^d$  to infinity by multiplying each value with a scalar  $\alpha$  as we take  $\alpha \rightarrow \infty$ .

## C Moment-Generating Function

Let  $X \in \mathbb{R}^d$  denote a random vector with distribution  $p_X$ . The **moment-generating function** (MGF) of  $X$  is a real-valued mapping  $M_X : \mathbb{R}^d \rightarrow \mathbb{R}$  defined as

$$M_X(t) := \mathbf{E}_{x \sim p_X} [\exp(t^\top x)] \quad (33)$$

Not every distribution has a corresponding MGF (because (33) may diverge). But a classical result in probability theory is that an MGF uniquely determines a probability distribution. More formally, let  $X, Y \in \mathbb{R}^d$  be random vectors with distributions  $p_X, p_Y$  with well-defined MGFs  $M_X, M_Y$ . Then  $p_X = p_Y$  iff  $M_X = M_Y$ . Thus an MGF is an alternative characterization of a random variable.

What makes  $M_X$  special is obviously the exponential function. Since  $e^z = \sum_{n=0}^{\infty} \frac{z^n}{n!}$ ,

$$M_X(t) = 1 + t^\top \underbrace{\mathbf{E}[X]}_{\text{1st moment}} + \frac{1}{2} t^\top \underbrace{\mathbf{E}[XX^\top]}_{\text{2nd moment}} t + \dots$$

so that  $\frac{\partial^n M_X(t)}{\partial t^n} |_{t=0_d}$  is the  $n$ -th moment of  $p_X$  (hence the name).

The MGF of a linear transformation of  $X$  is

$$M_{AX+b}(t) = \mathbf{E}_{x \sim p_X} [\exp(t^\top Ax) \exp(t^\top b)] = \exp(t^\top b) M_X(A^\top t) \quad (34)$$

The *log* of an MGF is convex: by Hölder's inequality  $\mathbf{E}[|XY|] \leq \mathbf{E}[|X|^p]^{1/p} + \mathbf{E}[|Y|^q]^{1/q}$ ,

$$\begin{aligned}
\log M_X(\alpha t + (1 - \alpha)w) &= \log \mathbf{E} [\exp(\alpha t^\top X) \exp((1 - \alpha)w^\top X)] \\
&\leq \log \left( \mathbf{E} [\exp(t^\top X)^\alpha] + \mathbf{E} [\exp(w^\top X)^{(1-\alpha)}] \right) \\
&= \alpha \log M_X(t) + (1 - \alpha) \log M_X(w)
\end{aligned}$$

**Lemma C.1.** Let  $X \sim \mathcal{N}(\mu, \Sigma)$ . Then

$$M_X(t) = \exp \left( t^\top \mu + \frac{1}{2} t^\top \Sigma t \right)$$

*Proof.* We use the same substitution in the proof of Lemma D.4. Let  $\Sigma = U\Lambda U^\top$  denote an orthonormal eigendecomposition. Let  $u = g(x)$  where  $g(x) = \Lambda^{-1/2} U^\top (x - \mu)$ , which implies  $x = U\Lambda^{1/2}u + \mu$ . Thus

$|\det(J_g(x))| = |\det(\Lambda^{-1/2}U^\top)| = \det(\Lambda)^{-1/2}$ , so we have the infinitesimal  $dx = \sqrt{\det(\Lambda)}du$ . Then

$$\begin{aligned}
& \int_{\mathbb{R}^d} \frac{1}{(\sqrt{2\pi})^d \sqrt{\det(\Sigma)}} \exp\left(-\frac{1}{2}(x-\mu)^\top \Sigma^{-1}(x-\mu)\right) \exp(t^\top x) dx \\
&= \int_{\mathbb{R}^d} \frac{\sqrt{\det(\Lambda)}}{(\sqrt{2\pi})^d \sqrt{\det(\Lambda)}} \exp\left(-\frac{1}{2}u^\top u\right) \exp(t^\top U\Lambda^{1/2}u + t^\top \mu) du \\
&= \exp(t^\top \mu) \int_{\mathbb{R}^d} \frac{1}{(\sqrt{2\pi})^d} \exp\left(-\frac{1}{2}u^\top u + t^\top U\Lambda^{1/2}u\right) du \\
&= \exp(t^\top \mu) \int_{\mathbb{R}^d} \frac{1}{(\sqrt{2\pi})^d} \exp\left(-\frac{1}{2}\|u - U\Lambda^{1/2}t\|^2 + \frac{1}{2}t^\top \Sigma t\right) du \\
&= \exp\left(t^\top \mu + \frac{1}{2}t^\top \Sigma t\right) \int_{\mathbb{R}^d} \frac{1}{(\sqrt{2\pi})^d} \exp\left(-\frac{1}{2}\|u - U\Lambda^{1/2}t\|^2\right) du \\
&= \exp\left(t^\top \mu + \frac{1}{2}t^\top \Sigma t\right)
\end{aligned}$$

□

An interesting consequence of the Gaussian MGF in Lemma C.1 is that a point-mass density can be viewed as a degenerate Gaussian distribution with zero variance. That is, if  $X \in \mathbb{R}^d$  takes value  $a \in \mathbb{R}^d$  with probability 1, then  $M_X(t) = \exp(a^\top t)$ , which is equal to the Gaussian MGF with  $\Sigma = 0_{d \times d}$ .

One application of MGF is showing that a linear transformation of a Gaussian random variable is also Gaussian. This applies the fact that a Gaussian MGF is an exponential to the MGF of a linear transformation (34).

**Lemma C.2.** Let  $X \sim \mathcal{N}(\mu, \Sigma)$ . Let  $A \in \mathbb{R}^{d' \times d}$  and  $b \in \mathbb{R}^{d'}$  where  $d' \leq d$  and  $A$  has full rank. Then  $AX + b \sim \mathcal{N}(A\mu + b, A\Sigma A^\top)$ .

*Proof.* For any  $t \in \mathbb{R}^{d'}$ ,

$$\begin{aligned}
M_{AX+b}(t) &= \exp(t^\top b) M_X(A^\top t) && \text{(by (34))} \\
&= \exp(t^\top b) \exp\left(t^\top A\mu + \frac{1}{2}t^\top A\Sigma A^\top t\right) && \text{(by Lemma C.1)} \\
&= \exp\left(t^\top (A\mu + b) + \frac{1}{2}t^\top A\Sigma A^\top t\right)
\end{aligned}$$

The last term is the MGF of a random variable with distribution  $\mathcal{N}(A\mu + b, A\Sigma A^\top)$  where  $A\Sigma A^\top \succ 0$ . The statement follows from the one-to-one correspondence between MGFs and distributions. □

## D Lemmas

**Lemma D.1** (Polar coordinates). For any integrable  $f : \mathbb{R} \rightarrow \mathbb{R}$ ,

$$\int_{\mathbb{R}^2} f(x^2 + y^2) d(x, y) = 2\pi \int_0^\infty f(r^2) r dr$$

*Proof.* Let  $R = [0, \infty) \times [0, 2\pi]$  and define  $g : R \rightarrow \mathbb{R}^2$  by  $g(r, \theta) = (r \cos \theta, r \sin \theta)$ . Note that  $r^2 = x^2 + y^2$  and  $g(R) = \mathbb{R}^2$ . The Jacobian of  $g$  at  $(r, \theta)$  is

$$J_g(r, \theta) = \begin{bmatrix} \frac{\partial r \cos \theta}{\partial r} & \frac{\partial r \cos \theta}{\partial \theta} \\ \frac{\partial r \sin \theta}{\partial r} & \frac{\partial r \sin \theta}{\partial \theta} \end{bmatrix} = \begin{bmatrix} \cos \theta & -r \sin \theta \\ \sin \theta & r \cos \theta \end{bmatrix}$$

Thus  $|\det(J_g(r, \theta))| = |r(\cos^2 \theta + \sin^2 \theta)| = r$ . Thus

$$\begin{aligned}
\int_{\mathbb{R}^2} f(x^2 + y^2) d(x, y) &= \int_R f(g_1(r, \theta)^2 + g_2(r, \theta)^2) |J_g(r, \theta)| d(r, \theta) && \text{(by (15))} \\
&= \int_R f(r^2) r d(r, \theta) \\
&= \int_0^\infty \left( \int_0^{2\pi} \exp(-r^2) r d\theta \right) dr && \text{(Fubini)} \\
&= \int_0^\infty 2\pi \exp(-r^2) r dr && \text{(FTC)} \\
&= 2\pi \int_0^\infty \exp(-r^2) r dr && \text{(linearity)}
\end{aligned}$$

□

**Lemma D.2** (Gaussian integral).

$$\int_{-\infty}^\infty \exp(-x^2) dx = \sqrt{\pi} \tag{35}$$

*Proof.* A standard proof shows that  $(\int_{-\infty}^\infty \exp(-x^2) dx)^2 = \pi$  as follows:

$$\begin{aligned}
\left( \int_{-\infty}^\infty \exp(-x^2) dx \right) \left( \int_{-\infty}^\infty \exp(-y^2) dy \right) &= \int_{-\infty}^\infty \left( \int_{-\infty}^\infty \exp(-x^2) dx \right) \exp(-y^2) dy && \text{(linearity)} \\
&= \int_{-\infty}^\infty \left( \int_{-\infty}^\infty \exp(-x^2) \exp(-y^2) dx \right) dy && \text{(linearity)} \\
&= \int_{\mathbb{R}^2} \exp(-(x^2 + y^2)) d(x, y) && \text{(Fubini)} \\
&= 2\pi \int_0^\infty \exp(-r^2) r dr && \text{(Lemma D.1)} \\
&= 2\pi \left( -\frac{1}{2} \exp(-r^2) \right) \Big|_0^\infty && \text{(FTC)} \\
&= 2\pi \left( 0 + \frac{1}{2} \right) = \pi
\end{aligned}$$

□

**Lemma D.3.** For any  $\mu \in \mathbb{R}$  and  $\sigma^2 > 0$ ,

$$\int_{-\infty}^\infty \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) dx = 1 \tag{36}$$

*Proof.* Let  $u = \frac{x-\mu}{\sqrt{2\sigma}}$  which gives the infinitesimal  $dx = \sqrt{2\sigma} du$ . Then

$$\begin{aligned}
\int_{-\infty}^\infty \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) dx &= \int_{-\infty}^\infty \frac{\sqrt{2\sigma}}{\sqrt{2\pi\sigma}} \exp(-u^2) du && \text{(by (13))} \\
&= \int_{-\infty}^\infty \frac{1}{\sqrt{\pi}} \exp(-u^2) du \\
&= \frac{1}{\sqrt{\pi}} \int_{-\infty}^\infty \exp(-u^2) du && \text{(linearity)} \\
&= 1 && \text{(Lemma D.2)}
\end{aligned}$$

□

**Lemma D.4.**

$$\int_{\mathbb{R}^d} \frac{1}{(\sqrt{2\pi})^d \sqrt{\det(\Sigma)}} \exp\left(-\frac{1}{2}(x-\mu)^\top \Sigma^{-1}(x-\mu)\right) dx = 1$$

*Proof.* Let  $\Sigma = U\Lambda U^\top$  denote an orthonormal eigendecomposition. Let  $u = g(x)$  where  $g(x) = \Lambda^{-1/2}U^\top(x-\mu)$ . Thus  $|\det(J_g(x))| = |\det(\Lambda^{-1/2}U^\top)| = \det(\Lambda)^{-1/2}$ , so we have the infinitesimal  $dx = \sqrt{\det(\Lambda)}du$ . Then

$$\begin{aligned} \int_{\mathbb{R}^d} \frac{1}{(\sqrt{2\pi})^d \sqrt{\det(\Sigma)}} \exp\left(-\frac{1}{2}(x-\mu)^\top \Sigma^{-1}(x-\mu)\right) dx &= \int_{\mathbb{R}^d} \frac{\sqrt{\det(\Lambda)}}{(\sqrt{2\pi})^d \sqrt{\det(\Lambda)}} \exp\left(-\frac{1}{2}u^\top u\right) du \\ &= \int_{\mathbb{R}^d} \prod_{i=1}^d \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{u_i^2}{2}\right) du \end{aligned}$$

By Fubini and linearity,

$$\begin{aligned} \int_{\mathbb{R}^d} \prod_{i=1}^d \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{u_i^2}{2}\right) du &= \int_{-\infty}^{\infty} \left( \cdots \left( \int_{-\infty}^{\infty} \prod_{i=1}^d \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{u_i^2}{2}\right) du_1 \right) \cdots \right) du_d \\ &= \prod_{i=1}^d \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{u_i^2}{2}\right) du_i \\ &= \left( \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right) dx \right)^d = 1 \end{aligned}$$

where the last step applies Lemma D.3 with  $\mu = 0$  and  $\sigma^2 = 1$ . □

**Lemma D.5.** For any  $\lambda > 0$ , the exponential distribution  $e_\lambda(x) := \lambda \exp(-\lambda x)$  over  $[0, \infty)$  has entropy

$$H(X) = 1 - \log \lambda$$

*Proof.*

$$\begin{aligned} H(X) &= - \int_0^\infty \lambda \exp(-\lambda x) \log(\lambda \exp(-\lambda x)) dx \\ &= - \log \lambda - \lambda \int_0^\infty \exp(-\lambda x)(-\lambda x) dx \end{aligned}$$

We evaluate the last integral as follows. Let  $u = g(x) = -\lambda x$ , then  $g'(x) = -\lambda$  so that  $|g'(g^{-1}(u))^{-1}| = 1/\lambda$ . Reorienting the region between  $g(0) = 0$  and  $g(\infty) = -\infty$  and applying (14),

$$\begin{aligned} \lambda \int_0^\infty \exp(-\lambda x)(-\lambda x) dx &= \int_{-\infty}^0 \exp(u)u du \\ &= \exp(u)u \Big|_{-\infty}^0 - \int_{-\infty}^0 \exp(u) du && \text{(integration by parts (12))} \\ &= (0 - 0) - \exp(u) \Big|_{-\infty}^0 && (\lim_{u \rightarrow -\infty} \exp(u)u = 0) \\ &= -1 \end{aligned}$$

□

**Lemma D.6.** Define  $\Delta := \mu' - \mu$ . Then

$$H(\mathcal{N}(\mu', \Sigma'), \mathcal{N}(\mu, \Sigma)) = \frac{1}{2} \Delta^\top \Sigma^{-1} \Delta + \frac{1}{2} \text{tr}(\Sigma^{-1} \Sigma') + \frac{1}{2} \log((2\pi)^d \det(\Sigma))$$

*Proof.*

$$\begin{aligned} H(\mathcal{N}(\mu', \Sigma'), \mathcal{N}(\mu, \Sigma)) &:= \mathbf{E}_{x \sim \mathcal{N}(\mu', \Sigma')} [-\log \mathcal{N}(\mu, \Sigma)(x)] \\ &= \frac{1}{2} \mathbf{E}_{x \sim \mathcal{N}(\mu', \Sigma')} [(x - \mu)^\top \Sigma^{-1} (x - \mu)] + \frac{1}{2} \log((2\pi)^d \det(\Sigma)) \end{aligned}$$

By the cyclic property and the linearity of trace,

$$\begin{aligned} \mathbf{E}_{x \sim \mathcal{N}(\mu', \Sigma')} [(x - \mu)^\top \Sigma^{-1} (x - \mu)] &= \mathbf{E}_{x \sim \mathcal{N}(\mu', \Sigma')} [\text{tr}((x - \mu)^\top \Sigma^{-1} (x - \mu))] \\ &= \mathbf{E}_{x \sim \mathcal{N}(\mu', \Sigma')} [\text{tr}(\Sigma^{-1} (x - \mu)(x - \mu)^\top)] \\ &= \text{tr} \left( \Sigma^{-1} \mathbf{E}_{x \sim \mathcal{N}(\mu', \Sigma')} [(x - \mu)(x - \mu)^\top] \right) \end{aligned}$$

Rewriting the expectation,

$$\begin{aligned} \mathbf{E}_{x \sim \mathcal{N}(\mu', \Sigma')} [(x - \mu)(x - \mu)^\top] &= \mathbf{E}_{x \sim \mathcal{N}(\mu', \Sigma')} [(x - \mu' + \Delta)(x - \mu' + \Delta)^\top] \\ &= \mathbf{E}_{x \sim \mathcal{N}(\mu', \Sigma')} [(x - \mu')(x - \mu')^\top + (x - \mu')\Delta^\top + \Delta(x - \mu')^\top + \Delta\Delta^\top] \\ &= \Sigma' + \Delta\Delta^\top \end{aligned}$$

Therefore we have

$$\begin{aligned} H(\mathcal{N}(\mu', \Sigma'), \mathcal{N}(\mu, \Sigma)) &= \frac{1}{2} \text{tr}(\Sigma^{-1} \Sigma' + \Sigma^{-1} \Delta\Delta^\top) + \frac{1}{2} \log((2\pi)^d \det(\Sigma)) \\ &= \frac{1}{2} \text{tr}(\Sigma^{-1} \Sigma') + \frac{1}{2} \Delta^\top \Sigma^{-1} \Delta + \frac{1}{2} \log((2\pi)^d \det(\Sigma)) \end{aligned}$$

□

**Corollary D.7** (Of Lemma D.6).

$$H(\mathcal{N}(\mu, \Sigma)) = \frac{1}{2} \log((2\pi e)^d \det(\Sigma))$$

**Corollary D.8** (Of Lemma D.6 and Corollary D.7). Define  $\Delta := \mu' - \mu$ . Then

$$D_{\text{KL}}(\mathcal{N}(\mu', \Sigma') || \mathcal{N}(\mu, \Sigma)) = \frac{1}{2} \Delta^\top \Sigma^{-1} \Delta + \frac{1}{2} \text{tr}(\Sigma^{-1} \Sigma' - I_{d \times d}) + \frac{1}{2} \log \frac{\det(\Sigma)}{\det(\Sigma')}$$

**Lemma D.9.** Let  $A \in \mathbb{R}^{d \times d}$ . The **main-diagonal block matrix** of  $A$  at index  $k \in \{1 \dots d\}$  with size  $n$  is a matrix  $B(k, n) \in \mathbb{R}^{n \times n}$  with entries  $B_{i,j}(k, n) = A_{k+i-1, k+j-1}$  for  $i, j \in \{1 \dots n\}$ . If  $A \succ 0$ , then  $B(k, n) \succ 0$  for all valid  $k, n$ .

*Proof.* Suppose  $u^\top B(k, n) u \leq 0$  for some nonzero  $u \in \mathbb{R}^n$ . Define  $v \in \mathbb{R}^d$  where  $v_{k+i-1} = u_i$  for  $i = 1 \dots n$  and other entries are zero. Then  $v$  is nonzero and  $v^\top A v = u^\top B(k, n) u \leq 0$ , contradicting the premise that  $A \succ 0$ . □

**Lemma D.10.** Let  $X \in \mathbb{R}^d$  and  $Y \in \mathbb{R}^{d'}$  be jointly normal with parameters  $(\mu, \Sigma)$ . Assume that  $\Sigma_Y - \Sigma_{YX} \Sigma_X^{-1} \Sigma_{XY}$  is invertible. Then for any  $z = (x, y) \in \mathbb{R}^{d+d'}$ ,

$$\begin{aligned} \frac{1}{(\sqrt{2\pi})^{d+d'} \sqrt{\det(\Sigma)}} \exp\left(-\frac{1}{2}(z - \mu)^\top \Sigma^{-1} (z - \mu)\right) &= \frac{1}{(\sqrt{2\pi})^d \sqrt{\det(\Sigma_X)}} \exp\left(-\frac{1}{2}(x - \mu_X)^\top \Sigma_X^{-1} (x - \mu_X)\right) \\ &\quad \times \frac{1}{(\sqrt{2\pi})^{d'} \sqrt{\det(\Omega)}} \exp\left(-\frac{1}{2}(y - \phi(x))^\top \Omega^{-1} (y - \phi(x))\right) \end{aligned} \quad (37)$$

where  $\Omega \in \mathbb{R}^{d' \times d'}$  and  $\phi(x) \in \mathbb{R}^{d'}$  are defined as

$$\Omega := \Sigma_Y - \Sigma_{YX} \Sigma_X^{-1} \Sigma_{XY} \quad (38)$$

$$\phi(x) := \mu_Y + \Sigma_{YX} \Sigma_X^{-1} (x - \mu_X) \quad \forall x \in \mathbb{R}^d \quad (39)$$

*Proof.* By [block matrix inversion](#) and abbreviating  $O = \Sigma_X^{-1}\Sigma_{XY}$ ,

$$\Sigma^{-1} = \begin{bmatrix} \Sigma_X & \Sigma_{XY} \\ \Sigma_{YX} & \Sigma_Y \end{bmatrix}^{-1} = \begin{bmatrix} \Sigma_X^{-1} + O\Omega^{-1}O^\top & -O\Omega^{-1} \\ -\Omega^{-1}O^\top & \Omega^{-1} \end{bmatrix}$$

Abbreviating  $u = x - \mu_X$  and  $v = y - \mu_Y$ ,

$$\begin{aligned} (z - \mu)^\top \Sigma^{-1}(z - \mu) &= u^\top (\Sigma_X^{-1} + O\Omega^{-1}O^\top) u - u^\top O\Omega^{-1}v - v^\top \Omega^{-1}O^\top u + v^\top \Omega^{-1}v \\ &= u^\top \Sigma_X^{-1}u + u^\top O\Omega^{-1}O^\top u - 2u^\top O\Omega^{-1}v + v^\top \Omega^{-1}v \\ &= u^\top \Sigma_X^{-1}u + (v - O^\top u)^\top \Omega^{-1}(v - O^\top u) \\ &= (x - \mu_X)^\top \Sigma_X^{-1}(x - \mu_X) + (y - \phi(x))^\top \Omega^{-1}(y - \phi(x)) \end{aligned}$$

where we use the fact that  $\Omega$  is symmetric. By the [determinant identity of a block matrix](#), we have  $\det(\Sigma) = \det(\Sigma_X\Omega) = \det(\Sigma_X)\det(\Omega)$ . Applying these identities to the LHS of (37) yields the RHS.  $\square$

**Lemma D.11.** Let  $X \in \mathbb{R}^d$  and  $Y \in \mathbb{R}^{d'}$  be jointly normal with parameters  $(\mu, \Sigma)$ . Assume that  $\Sigma_Y - \Sigma_{YX}\Sigma_X^{-1}\Sigma_{XY}$  is invertible. Then

$$H(Y|X) = \frac{1}{2} \log \left( (2\pi e)^{d'} \det(\Sigma_Y - \Sigma_{YX}\Sigma_X^{-1}\Sigma_{XY}) \right) \quad (40)$$

$$I(X, Y) = \frac{1}{2} \log \left( \frac{\det(\Sigma_X)\det(\Sigma_Y)}{\det(\Sigma)} \right) \quad (41)$$

*Proof.* By Lemma D.10,  $Y|X = x$  is distributed as  $\mathcal{N}(\phi(x), \Omega)$  for any  $x \in \mathbb{R}^d$  where  $\phi(x) := \mu_Y + \Sigma_{YX}\Sigma_X^{-1}(x - \mu_X)$  and  $\Omega := \Sigma_Y - \Sigma_{YX}\Sigma_X^{-1}\Sigma_{XY}$ . Thus

$$H(Y|X = x) = \mathbf{E}[-\log \Pr(Y|X = x)] = \frac{1}{2} \mathbf{E}[(Y - \phi(x))^\top \Omega^{-1}(Y - \phi(x))] + \frac{1}{2} \log((2\pi)^{d'} \det(\Omega))$$

Using trace similarly as in the proof of Lemma D.6, we can verify

$$\mathbf{E}[(Y - \phi(x))^\top \Omega^{-1}(Y - \phi(x))] = \Omega^{-1}(\Sigma_Y - \Sigma_{YX}\Sigma_X^{-1}\Sigma_{XY})(x - \mu_X)(x - \mu_X)^\top \Sigma_X^{-1}\Sigma_{XY}$$

Taking the expectation over  $x$  yields  $I_{d' \times d}$ . This shows (40). To show (41), we have

$$\begin{aligned} I(X, Y) &= H(Y) - H(Y|X) \\ &= \frac{1}{2} \log \left( (2\pi e)^{d'} \det(\Sigma_Y) \right) - \frac{1}{2} \log \left( (2\pi e)^{d'} \det(\Omega) \right) \\ &= \frac{1}{2} \log \left( \frac{\det(\Sigma_Y)}{\det(\Omega)} \right) \\ &= \frac{1}{2} \log \left( \frac{\det(\Sigma_X)\det(\Sigma_Y)}{\det(\Sigma)} \right) \end{aligned}$$

where for the last equality we use the fact that  $\det(\Sigma) = \det(\Sigma_X\Omega) = \det(\Sigma_X)\det(\Omega)$ .  $\square$

**Lemma D.12.** The following statements about  $X \in \mathbb{R}^d$  are equivalent.

1.  $X \sim \mathcal{N}(\mu, \Sigma)$ .
2.  $M_X(t) = \exp(t^\top \mu + \frac{1}{2}t^\top \Sigma t)$  for all  $t \in \mathbb{R}^d$ .
3.  $X = \Sigma^{1/2}Z + \mu$  where  $Z \sim \mathcal{N}(0_d, I_{d \times d})$ .
4.  $Y = a^\top X$  has the density  $\mathcal{N}(a^\top \mu, a^\top \Sigma a)$  for all nonzero  $a \in \mathbb{R}^d$ .

*Proof.* Lemma C.1 gives  $1 \equiv 2$ . To show  $2 \equiv 3$  we note that by (34)

$$M_{\Sigma^{1/2}Z+\mu}(t) = \exp(t^\top \mu) M_Z(\Sigma^{1/2}t) = \exp\left(t^\top \mu + \frac{1}{2}t^\top \Sigma t\right) = M_X(t)$$

We have  $1 \Rightarrow 4$  since the density of  $Y$  is  $\mathcal{N}(a^\top \mu, a^\top \Sigma a)$  by Lemma C.2. To show  $4 \Rightarrow 2$ , pick any nonzero  $a \in \mathbb{R}^d$ . For all  $t \in \mathbb{R}$

$$M_X(ta) = M_{a^\top X}(t) = \exp\left(ta^\top \mu + \frac{1}{2}t^2 a^\top \Sigma a\right)$$

where the first equality uses (34) and the second equality uses Lemma C.1. Setting  $t = 1$  gives  $M_X(a) = \exp(a^\top \mu + \frac{1}{2}a^\top \Sigma a)$ . Additionally,  $M_X(0_d) = 1 = \exp(0_d^\top \mu + \frac{1}{2}0_d^\top \Sigma 0_d)$ . Thus  $M_X(t) = \exp(t^\top \mu + \frac{1}{2}t^\top \Sigma t)$  for all  $t \in \mathbb{R}^d$ .  $\square$

## E Individually Normal But Not Jointly Normal

This is an [example from Wikipedia](#). Let  $X \sim \mathcal{N}(0, 1)$  and, independently,  $\epsilon \sim R$  where  $R$  denotes the Rademacher distribution. Let  $Y = \epsilon X$ . By the symmetry of the distribution of  $X$ , we have  $Y \sim \mathcal{N}(0, 1)$ . More formally,

$$\begin{aligned} \Pr(Y \leq x) &= \Pr(\epsilon = 1) \Pr(X \leq x) + \Pr(\epsilon = -1) \Pr(X \geq -x) \\ &= \Pr(\epsilon = 1) \Pr(X \leq x) + \Pr(\epsilon = -1) \Pr(-X \leq x) \\ &= \frac{1}{2} \Pr(X \leq x) + \frac{1}{2} \Pr(X \leq x) \\ &= \Pr(X \leq x) \end{aligned}$$

Let  $Z = X + Y$ . Then  $Z = 0$  with probability  $\frac{1}{2}$  and  $Z = 2X$  with probability  $\frac{1}{2}$ , so

$$\Pr(Z = z) = \frac{1}{2} \left( \mathbb{1}_{[z=0]} + \mathcal{N}(0, 1)\left(\frac{z}{2}\right) \right) \quad (42)$$

which is not a normal distribution. Then by definition 4,  $(X, Y) \in \mathbb{R}^2$  is not normally distributed. Thus  $X$  and  $Y$  are not jointly normal, even though they are individually normal.

**Mutual information.**  $X$  and  $Y$  are uncorrelated. More formally,

$$\text{Cov}(X, Y) = \mathbf{E}[XY] - \mathbf{E}[X] \mathbf{E}[Y] = \mathbf{E}[\epsilon X^2] = \mathbf{E}[\epsilon] \mathbf{E}[X^2] = 0$$

Thus  $\text{cor}(X, Y) = 0$ . But  $X$  and  $Y$  are not independent. Specifically,  $\Pr(Y = x | X = x) = \frac{1}{2}$  is not equal to  $\Pr(Y = x) = \mathcal{N}(0, 1)(x)$  for any  $x \in \mathbb{R}$ . This illustrates the limitation of linear correlation. On the other hand, the mutual information between  $X$  and  $Y$  is positive:

$$I(X, Y) = H(X) - H(X|Y) = H(X) - \log(2) = \log \sqrt{\frac{\pi e}{2}} \approx 0.73$$