

# **Probability and Information Theory Overview**

Instructor: Xiaoqian Wang  
8/26/2024

Slides prepared based on the Lectures slides of Probability and Information Theory from  
[https://www.deeplearningbook.org/lecture\\_slides.html](https://www.deeplearningbook.org/lecture_slides.html)

# Probability Mass Function

- The domain of  $P$  must be the set of all possible states of  $x$ .
- $\forall x \in X, 0 \leq P(x) \leq 1$ . An impossible event has probability 0 and no state can be less probable than that. Likewise, an event that is guaranteed to happen has probability 1, and no state can have a greater chance of occurring.
- $\sum_{x \in X} P(x) = 1$ . We refer to this property as being **normalized**. Without this property, we could obtain probabilities greater than one by computing the probability of one of many events occurring.

Example: uniform distribution:  $P(X = x_i) = \frac{1}{k}$

# Probability Density Function

- The domain of  $p$  must be the set of all possible states of  $x$ .
- $\forall x \in X, p(x) \geq 0$ . Note that we do not require  $p(x) \leq 1$ .
- $\int p(x)dx = 1$ .

Example: uniform distribution:  $u(x; a, b) = \frac{1}{b-a}$ .

# Computing Marginal Probability with the Sum Rule

$$\forall x \in X, P(X = x) = \sum_y P(X = x, Y = y). \quad (3.3)$$

$$p(x) = \int p(x, y) dy. \quad (3.4)$$

# Conditional Probability

$$P(y = y \mid x = x) = \frac{P(y = y, x = x)}{P(x = x)}. \quad (3.5)$$

# Chain Rule of Probability

$$P(x^{(1)}, \dots, x^{(n)}) = P(x^{(1)}) \prod_{i=2}^n P(x^{(i)} \mid x^{(1)}, \dots, x^{(i-1)}). \quad (3.6)$$

# Bayes' Rule

$$P(x | y) = \frac{P(x)P(y | x)}{P(y)}. \quad (3.42)$$

# Independence

$$\forall x \in \mathbf{x}, y \in \mathbf{y}, p(\mathbf{x} = x, \mathbf{y} = y) = p(\mathbf{x} = x)p(\mathbf{y} = y). \quad (3.7)$$

# Conditional Independence

$$\forall x \in \mathbf{x}, y \in \mathbf{y}, z \in \mathbf{z}, p(\mathbf{x} = x, \mathbf{y} = y \mid \mathbf{z} = z) = p(\mathbf{x} = x \mid \mathbf{z} = z)p(\mathbf{y} = y \mid \mathbf{z} = z). \quad (3.8)$$

# Expectation

$$\mathbb{E}_{x \sim P}[f(x)] = \sum_x P(x)f(x), \quad (3.9)$$

$$\mathbb{E}_{x \sim p}[f(x)] = \int p(x)f(x)dx. \quad (3.10)$$

linearity of expectations:

$$\mathbb{E}_x[\alpha f(x) + \beta g(x)] = \alpha \mathbb{E}_x[f(x)] + \beta \mathbb{E}_x[g(x)], \quad (3.11)$$

# Variance and Covariance

$$\text{Var}(f(x)) = \mathbb{E} \left[ (f(x) - \mathbb{E}[f(x)])^2 \right]. \quad (3.12)$$

$$\text{Cov}(f(x), g(y)) = \mathbb{E} [(f(x) - \mathbb{E}[f(x)]) (g(y) - \mathbb{E}[g(y)])]. \quad (3.13)$$

Covariance matrix:

$$\text{Cov}(\mathbf{x})_{i,j} = \text{Cov}(\mathbf{x}_i, \mathbf{x}_j). \quad (3.14)$$

# Bernoulli Distribution

$$P(x = 1) = \phi \tag{3.16}$$

$$P(x = 0) = 1 - \phi \tag{3.17}$$

$$P(x = x) = \phi^x(1 - \phi)^{1-x} \tag{3.18}$$

$$\mathbb{E}_x[x] = \phi \tag{3.19}$$

$$\text{Var}_x(x) = \phi(1 - \phi) \tag{3.20}$$

# Gaussian Distribution

Parametrized by variance:

$$\mathcal{N}(x; \mu, \sigma^2) = \sqrt{\frac{1}{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right). \quad (3.21)$$

Parametrized by precision:

$$\mathcal{N}(x; \mu, \beta^{-1}) = \sqrt{\frac{\beta}{2\pi}} \exp\left(-\frac{1}{2}\beta(x - \mu)^2\right). \quad (3.22)$$

# Gaussian Distribution

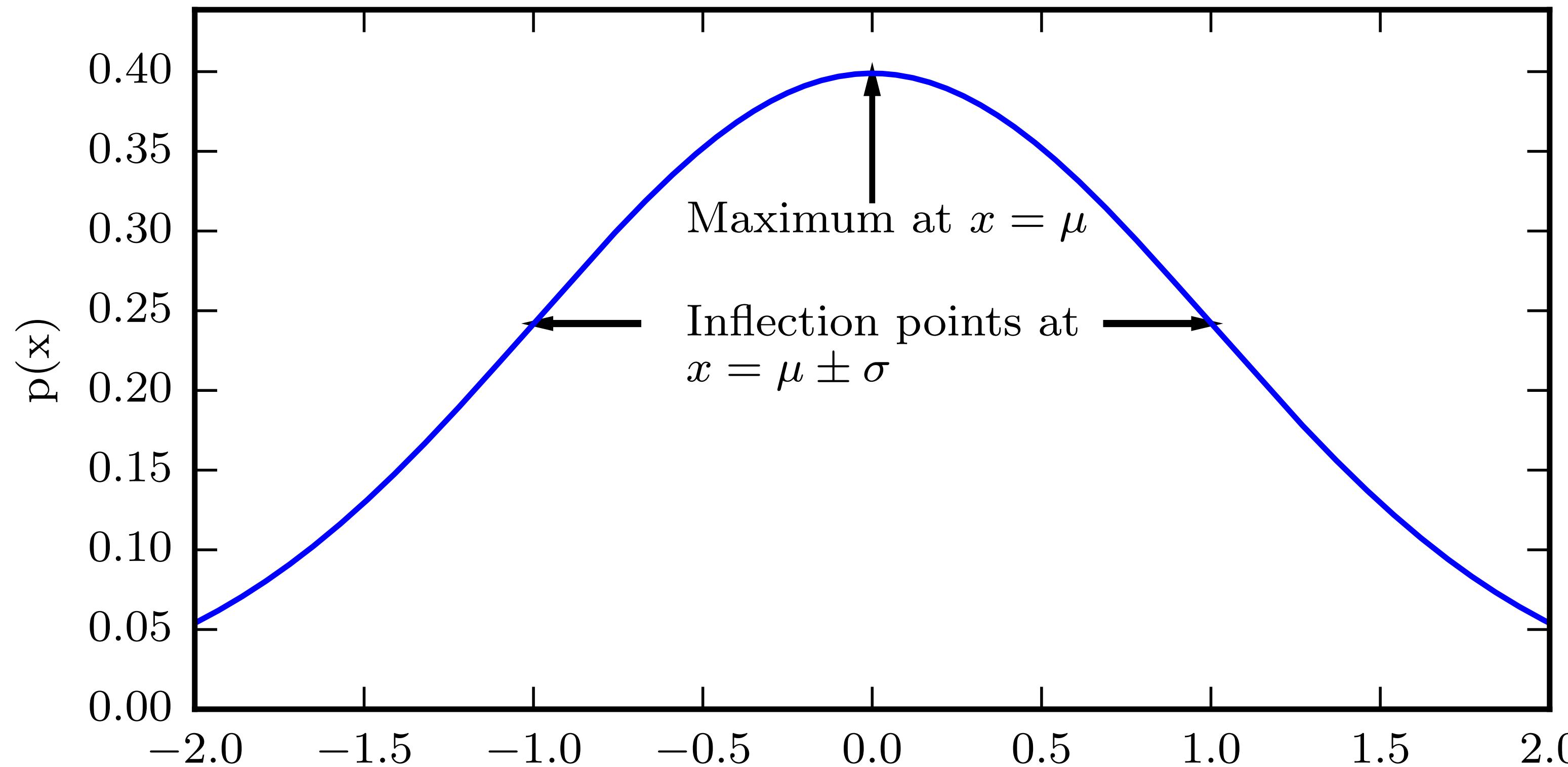


Figure 3.1

# Multivariate Gaussian

Parametrized by covariance matrix:

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sqrt{\frac{1}{(2\pi)^n \det(\boldsymbol{\Sigma})}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right). \quad (3.23)$$

Parametrized by precision matrix:

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\beta}^{-1}) = \sqrt{\frac{\det(\boldsymbol{\beta})}{(2\pi)^n}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\beta} (\mathbf{x} - \boldsymbol{\mu})\right). \quad (3.24)$$

# Mixture Distributions

$$P(\mathbf{x}) = \sum_i P(c = i)P(\mathbf{x} \mid c = i) \quad (3.29)$$

Gaussian mixture  
with three  
components

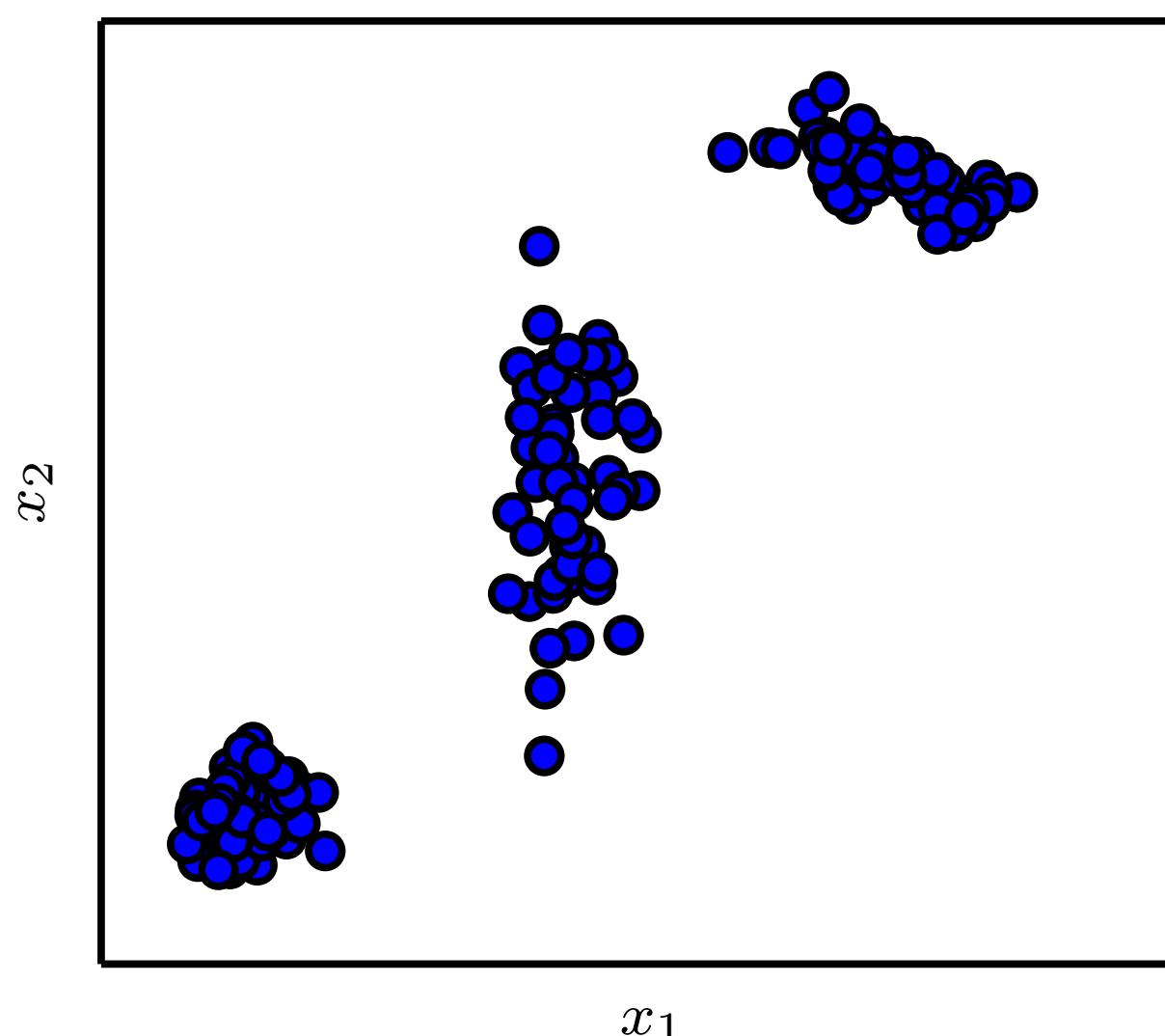


Figure 3.2

(Goodfellow 2016)

# Information Theory

Information:

$$I(x) = -\log P(x). \quad (3.48)$$

Entropy:

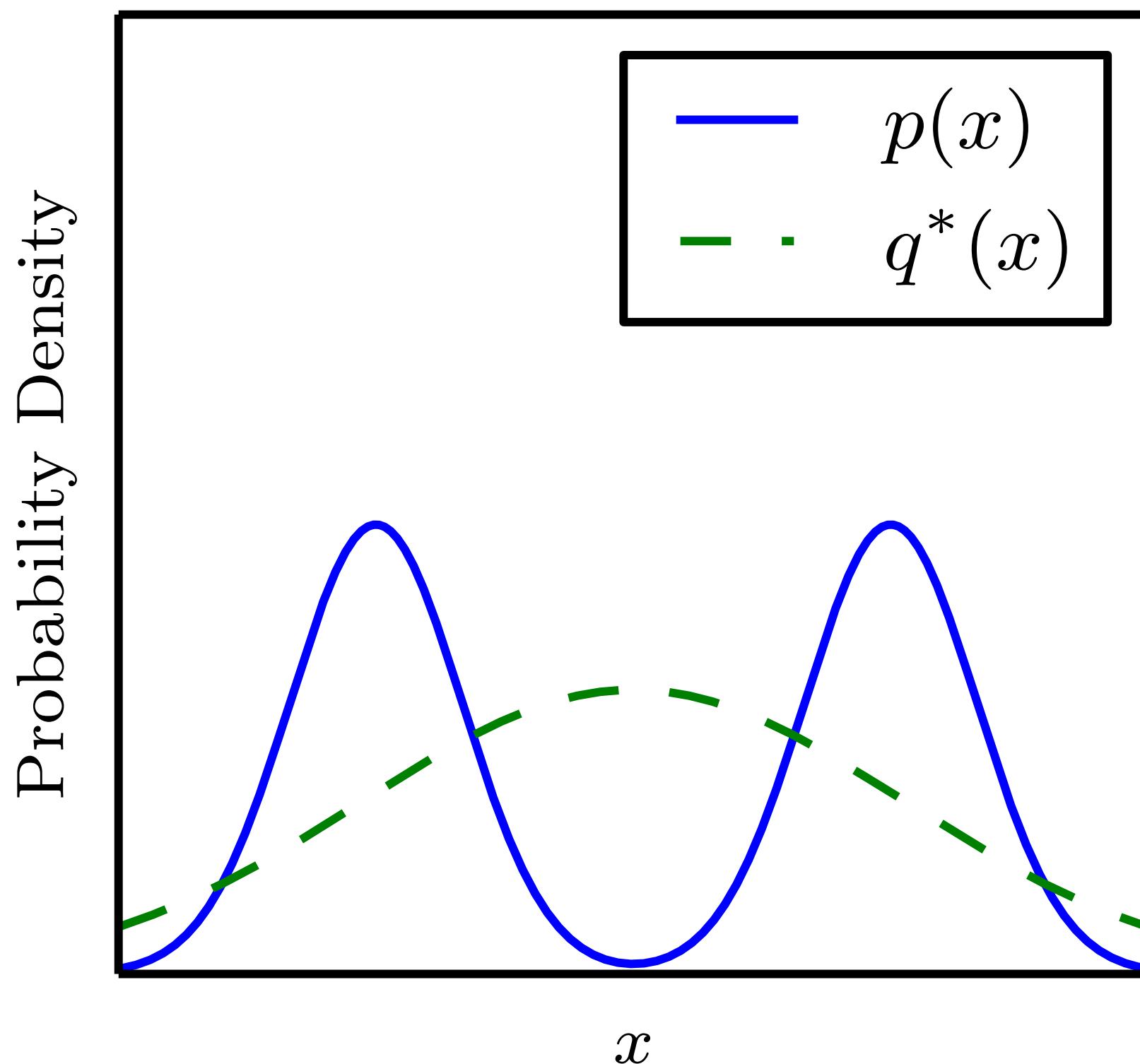
$$H(x) = \mathbb{E}_{x \sim P}[I(x)] = -\mathbb{E}_{x \sim P}[\log P(x)]. \quad (3.49)$$

KL divergence:

$$D_{\text{KL}}(P \| Q) = \mathbb{E}_{x \sim P} \left[ \log \frac{P(x)}{Q(x)} \right] = \mathbb{E}_{x \sim P} [\log P(x) - \log Q(x)]. \quad (3.50)$$

# The KL Divergence is Asymmetric

$$q^* = \operatorname{argmin}_q D_{\text{KL}}(p\|q)$$



$$q^* = \operatorname{argmin}_q D_{\text{KL}}(q\|p)$$

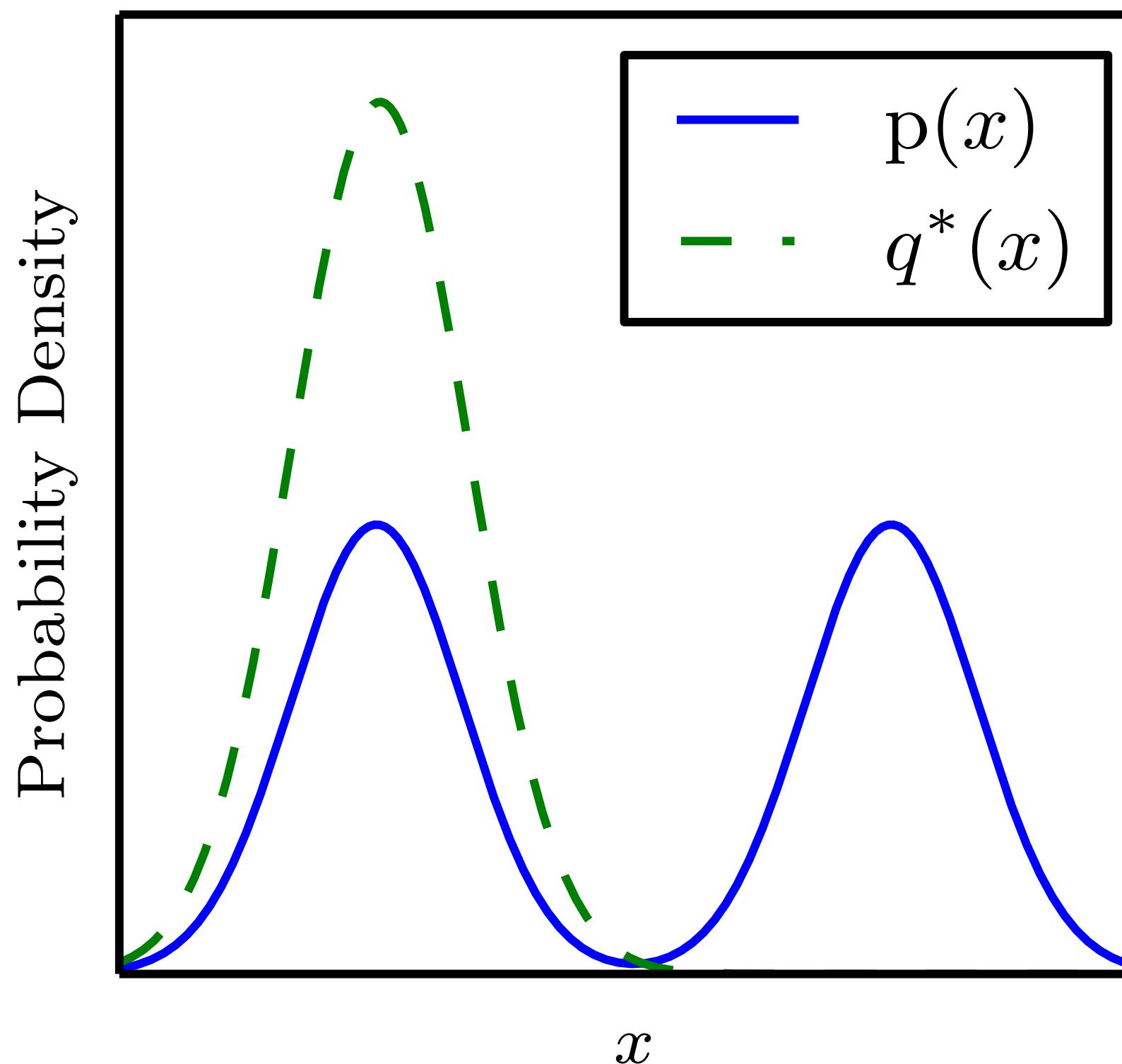


Figure 3.6

# Causality V.S. Correlation

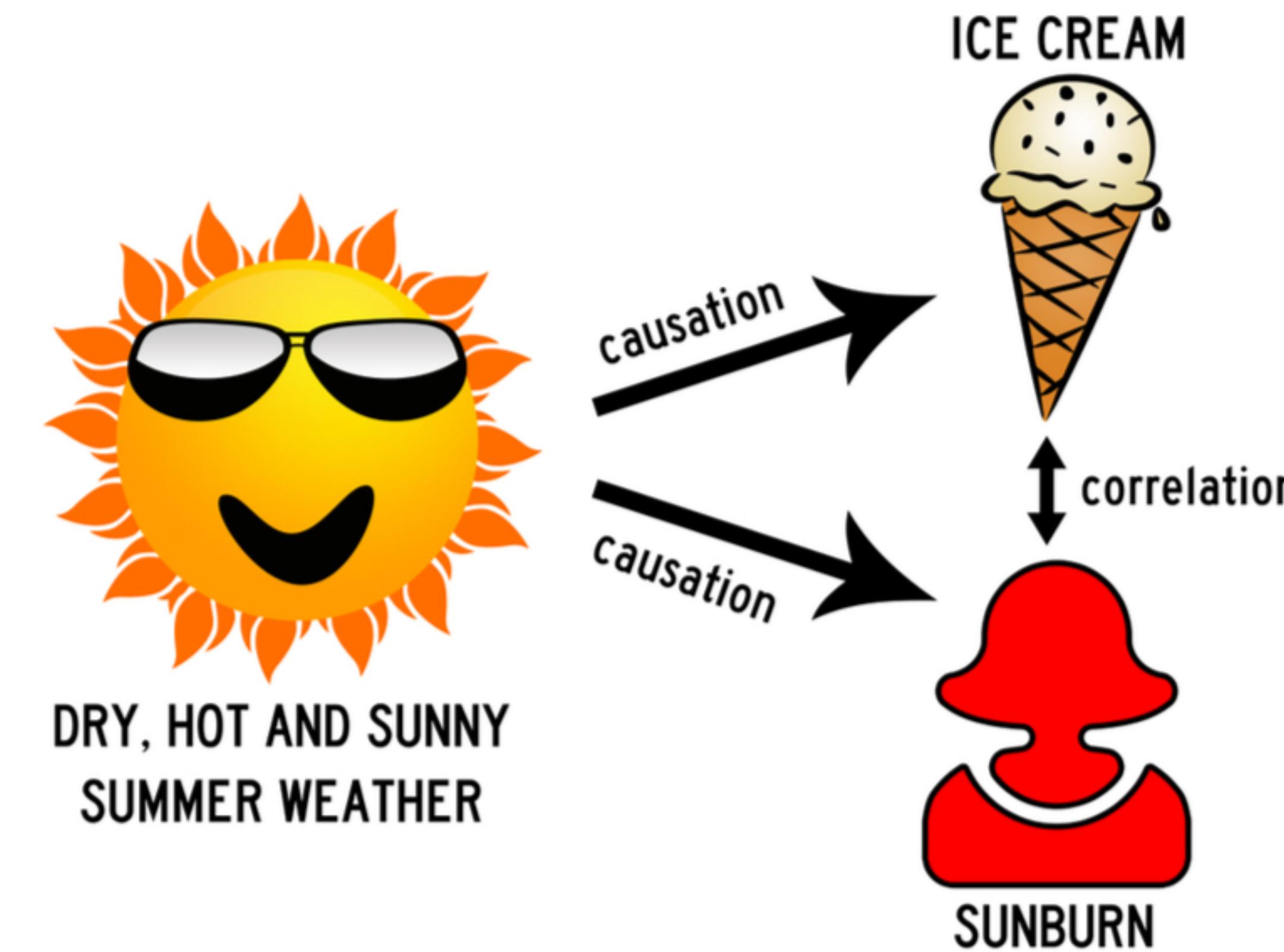


Image from <https://towardsdatascience.com/correlation-is-not-causation-ae05d03c1f53?gi=2fdda0721e2e>