

Probability Theory

Kevin P. Murphy: Probabilistic Machine Learning Ch2 & Ch3

Presented by Wending Liu

RSE ML and Econometrics Reading Group

Oct 23 2023

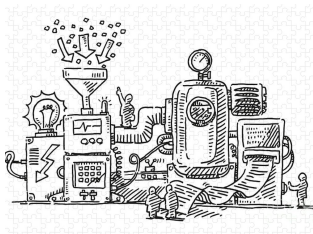
Table of Contents

- 1 What is Probability
- 2 Probability: Univariate Models
- 3 Probability: Multivariate Models

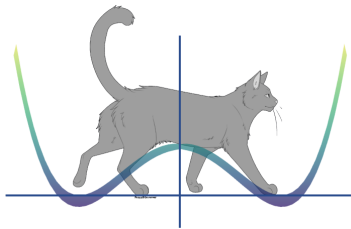
- What do we mean by saying that "the probability of a fair coin will land head is 50%"?
 - ▶ **Frequentist** interpretation: If we flip the coin 1,000,000 times, we will observe the coin land heads about 500,000 times.
 - ▶ **Bayesian** interpretation: We ***believe*** the coin is equally likely to land head or tail on our next toss.
- In the Bayesian view, probability is used to quantify our **uncertainty** about something.
- Bayesian interpretation can be used to model uncertainty about one-off events.
- Bayesian interpretation is used throughout the whole book.

Uncertainty

- But what is **uncertainty** in Bayesian interpretation?
- Epistemic/Model uncertainty
- Aleatoric/Data uncertainty



Complex machine, little human



Schrodinger's cat

Table of Contents

- 1 What is Probability
- 2 Probability: Univariate Models
- 3 Probability: Multivariate Models

Bayes' rule

- Bayes' rule is a formula for computing the probability distribution over possible values of unknown quantity H given some observed data $Y = y$:

$$p(H = h \mid Y = y) = \frac{p(H = h)p(Y = y \mid H = h)}{p(Y = y)}$$

- $p(H = h)$: prior distribution.
- $p(Y = y \mid H = h)$: likelihood.
- $p(Y = y)$: marginal likelihood.
- $p(H = h \mid Y = y)$: posterior distribution.
- posterior \propto prior \times likelihood.

Example: Testing for Covid-19

- Suppose COVID-19 prevalence is 1% now and you take a diagnostic test that has a positive result. What's the probability that you are infected?
- H : the indicator of infection. Y : the indicator of a positive test result.
- prior: $p(H = 1) = 0.01$, $p(H = 0) = 0.99$.
- likelihood: $p(Y = 1|H = 1) = 0.875$. (true positive rate)
- false positive rate: $p(Y = 1|H = 0) = 0.025$.

$$\begin{aligned} & p(H = 1 \mid Y = 1) \\ &= \frac{p(Y = 1 \mid H = 1)p(H = 1)}{p(Y = 1 \mid H = 1)p(H = 1) + p(Y = 1 \mid H = 0)p(H = 0)} \\ &= \frac{\text{TPR} \times \text{prior}}{\text{TPR} \times \text{prior} + \text{FPR} \times (1 - \text{prior})} \\ &= \frac{0.875 \times 0.01}{0.875 \times 0.01 + 0.025 \times 0.99} = 0.261 \end{aligned}$$

Bernoulli and binomial distribution

- $\text{Ber}(y \mid \theta) \triangleq \theta^y (1 - \theta)^{1-y} = \begin{cases} 1 - \theta & \text{if } y = 0 \\ \theta & \text{if } y = 1 \end{cases}$
- $\text{Bin}(s \mid N, \theta) \triangleq \binom{N}{s} \theta^s (1 - \theta)^{N-s}$
- We want to predict a binary variable $y \in \{0, 1\}$ given some inputs $\mathbf{x} \in \mathcal{X}$:

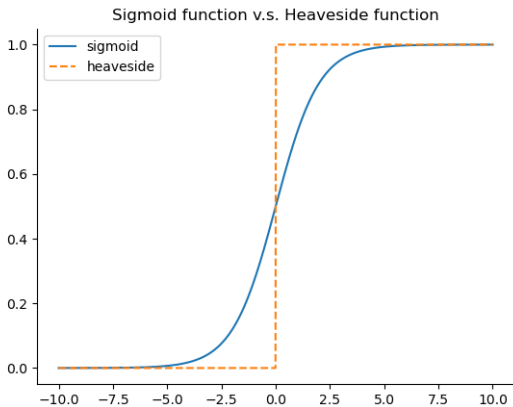
$$p(y \mid \mathbf{x}, \boldsymbol{\theta}) = \text{Ber}(y \mid f(\mathbf{x}; \boldsymbol{\theta}))$$

- To avoid the requirement that $0 \leq f(\mathbf{x}; \boldsymbol{\theta}) \leq 1$, we can let f be an unconstrained function, and use the following model:

$$p(y \mid \mathbf{x}, \boldsymbol{\theta}) = \text{Ber}(y \mid \sigma(f(\mathbf{x}; \boldsymbol{\theta})))$$

- σ is the **sigmoid** function: $\sigma(a) \triangleq \frac{1}{1+e^{-a}}$
- Binary logistic regression: $f(\mathbf{x}; \boldsymbol{\theta}) = \mathbf{w}^T \mathbf{x} + b$

Sigmoid function



Categorical and multinomial distributions

- **one-hot vector:** $\mathbf{y} \in \{0, 1\}^C$, $\sum_{c=1}^C y_c = 1$
- $\text{Cat}(y \mid \boldsymbol{\theta}) \triangleq \prod_{c=1}^C \theta_c^{I(y=c)} = \prod_{c=1}^C \theta_c^{y_c}$
- $\text{Mu}(\mathbf{s} \mid N, \boldsymbol{\theta}) \triangleq \binom{N}{s_1 \dots s_C} \prod_{c=1}^C \theta_c^{s_c}$
- $p(y \mid \mathbf{x}, \boldsymbol{\theta}) = \text{Cat}(y \mid f(\mathbf{x}; \boldsymbol{\theta}))$
- We require that $0 \leq f_c(\mathbf{x}; \boldsymbol{\theta}) \leq 1$ and $\sum_{c=1}^C f_c(\mathbf{x}; \boldsymbol{\theta}) = 1$.
- To avoid this requirement, we pass the output from f into the **softmax function**, also called the **multinomial logit**:

$$\mathcal{S}(\mathbf{a}) \triangleq \left[\frac{e^{a_1}}{\sum_{c'=1}^C e^{a_{c'}}}, \dots, \frac{e^{a_C}}{\sum_{c'=1}^C e^{a_{c'}}} \right]$$

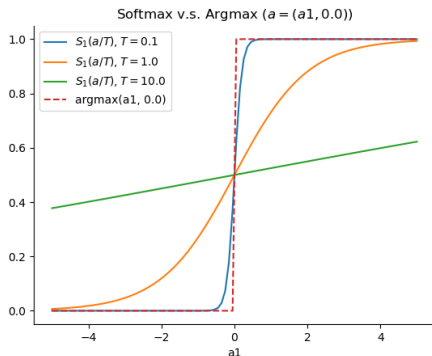
- Multinomial logistic regression: $f_c(\mathbf{x}; \boldsymbol{\theta}) = \mathbf{w}_c^T \mathbf{x} + b_c$:

$$p(y \mid \mathbf{x}; \boldsymbol{\theta}) = \text{Cat}(y \mid \mathcal{S}(\mathbf{W}\mathbf{x} + \mathbf{b}))$$

Softmax function

- Softmax function is related to Boltzmann distribution in physics.
- Let us divide each a_c by a constant T called the temperature. Then as $T \rightarrow 0$, we find:

$$\mathcal{S}(\mathbf{a}/T)_c = \begin{cases} 1.0 & \text{if } c = \operatorname{argmax}_{c'} a_{c'} \\ 0.0 & \text{otherwise} \end{cases}$$



Guassian distribution

- Gaussian pdf: $\mathcal{N}(y | \mu, \sigma^2) \triangleq \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(y-\mu)^2}$
- Gaussian distribution is widely used.
 - ▶ only two parameters.
 - ▶ central limit theorem.
 - ▶ maximum entropy given fixed mean and variance.
- A robust alternative to Gaussian is the **Student's t-distribution**:

$$\mathcal{T}(y | \mu, \sigma^2, \nu) \propto \left[1 + \frac{1}{\nu} \left(\frac{y - \mu}{\sigma} \right)^2 \right]^{-\left(\frac{\nu+1}{2}\right)}$$

where μ is the mean, $\sigma > 0$ is the scale parameter, and $\nu > 0$ is the degree of freedom (a better term would be the **degree of normality**).

- **Laplace** distribution (heavy tail): $\text{Lap}(y | \mu, b) \triangleq \frac{1}{2b} \exp\left(-\frac{|y-\mu|}{b}\right)$

Robustness of t distribution

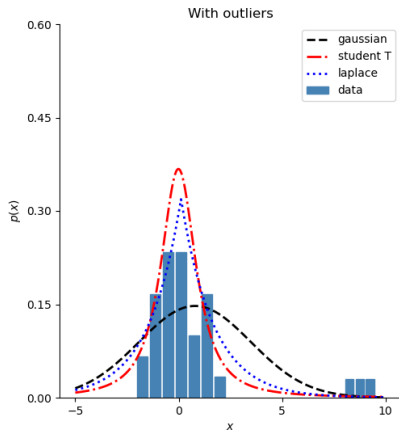
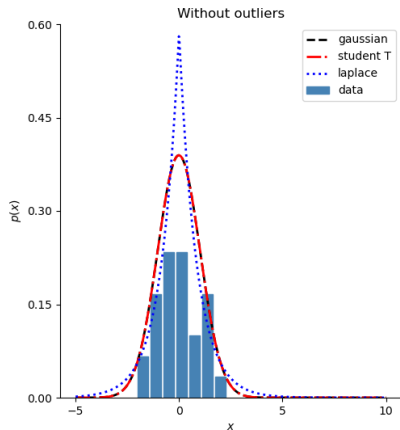
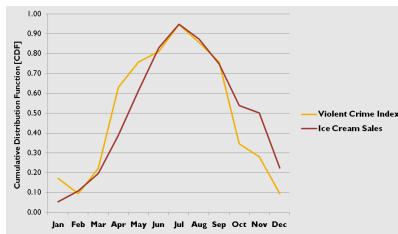


Table of Contents

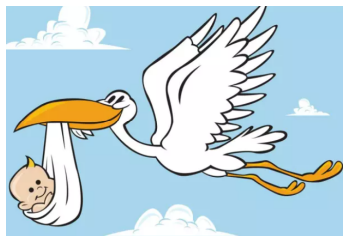
- 1 What is Probability
- 2 Probability: Univariate Models
- 3 Probability: Multivariate Models**

Causality, Independence and Correlation

- We use multivariate models to study the dependence of variables on each other.
- causality \nleftrightarrow dependence \Leftarrow correlation



ice cream makes people angry?



storks deliver babies?

Multivariate Gaussian distribution

MVN

The MVN density is defined by the following:

$$\mathcal{N}(\mathbf{y} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) \triangleq \frac{1}{(2\pi)^{D/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{y} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \boldsymbol{\mu}) \right]$$

where $\boldsymbol{\mu} = \mathbb{E}[\mathbf{y}] \in \mathbb{R}^D$ is the mean vector, and $\boldsymbol{\Sigma} = \text{Cov}[\mathbf{y}]$ is the $D \times D$ covariance matrix, defined as follows:

$$\begin{aligned} \text{Cov}[\mathbf{y}] &\triangleq \mathbb{E} \left[(\mathbf{y} - \mathbb{E}[\mathbf{y}])(\mathbf{y} - \mathbb{E}[\mathbf{y}])^\top \right] \\ &= \begin{pmatrix} \mathbb{V}[Y_1] & \text{Cov}[Y_1, Y_2] & \cdots & \text{Cov}[Y_1, Y_D] \\ \text{Cov}[Y_2, Y_1] & \mathbb{V}[Y_2] & \cdots & \text{Cov}[Y_2, Y_D] \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}[Y_D, Y_1] & \text{Cov}[Y_D, Y_2] & \cdots & \mathbb{V}[Y_D] \end{pmatrix} \end{aligned}$$

- Suppose $\mathbf{y} = (\mathbf{y}_1, \mathbf{y}_2)$ is jointly Gaussian with parameters

$$\boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix}, \boldsymbol{\Lambda} = \boldsymbol{\Sigma}^{-1} = \begin{pmatrix} \boldsymbol{\Lambda}_{11} & \boldsymbol{\Lambda}_{12} \\ \boldsymbol{\Lambda}_{21} & \boldsymbol{\Lambda}_{22} \end{pmatrix}$$

where $\boldsymbol{\Lambda}$ is the **precision matrix**.

Marginal Distributions

$$p(\mathbf{y}_1) = \mathcal{N}(\mathbf{y}_1 \mid \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11})$$

$$p(\mathbf{y}_2) = \mathcal{N}(\mathbf{y}_2 \mid \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_{22})$$

Posterior Conditional Distributions

$$p(\mathbf{y}_1 \mid \mathbf{y}_2) = \mathcal{N}(\mathbf{y}_1 \mid \boldsymbol{\mu}_{1|2}, \boldsymbol{\Sigma}_{1|2})$$

$$\boldsymbol{\mu}_{1|2} = \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{y}_2 - \boldsymbol{\mu}_2)$$

$$= \boldsymbol{\mu}_1 - \boldsymbol{\Lambda}_{11}^{-1}\boldsymbol{\Lambda}_{12}(\mathbf{y}_2 - \boldsymbol{\mu}_2)$$

$$= \boldsymbol{\Sigma}_{1|2}(\boldsymbol{\Lambda}_{11}\boldsymbol{\mu}_1 - \boldsymbol{\Lambda}_{12}(\mathbf{y}_2 - \boldsymbol{\mu}_2))$$

$$\boldsymbol{\Sigma}_{1|2} = \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21} = \boldsymbol{\Lambda}_{11}^{-1}$$

Linear Gaussian systems

Linear Gaussian system

Let $\mathbf{z} \in \mathbb{R}^L$ be an unknown vector of values, and $\mathbf{y} \in \mathbb{R}^D$ be some noisy measurement of \mathbf{z} with the following joint distribution:

$$\begin{aligned}p(\mathbf{z}) &= \mathcal{N}(\mathbf{z} \mid \boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z) \\p(\mathbf{y} \mid \mathbf{z}) &= \mathcal{N}(\mathbf{y} \mid \mathbf{W}\mathbf{z} + \mathbf{b}, \boldsymbol{\Sigma}_y)\end{aligned}$$

where \mathbf{W} is a matrix of size $D \times L$. The corresponding joint distribution, $p(\mathbf{z}, \mathbf{y}) = p(\mathbf{z})p(\mathbf{y} \mid \mathbf{z})$, is a $L + D$ dimensional Gaussian, with mean and covariance given by

$$\begin{aligned}\boldsymbol{\mu} &= \begin{pmatrix} \boldsymbol{\mu}_z \\ \mathbf{W}\boldsymbol{\mu}_z + \mathbf{b} \end{pmatrix} \\ \boldsymbol{\Sigma} &= \begin{pmatrix} \boldsymbol{\Sigma}_z & \boldsymbol{\Sigma}_z \mathbf{W}^\top \\ \mathbf{W}\boldsymbol{\Sigma}_z & \boldsymbol{\Sigma}_y + \mathbf{W}\boldsymbol{\Sigma}_z \mathbf{W}^\top \end{pmatrix}\end{aligned}$$

Bayes rule for Gaussians

$$p(\mathbf{z} | \mathbf{y}) = \mathcal{N}(\mathbf{z} | \boldsymbol{\mu}_{\mathbf{z}|\mathbf{y}}, \boldsymbol{\Sigma}_{\mathbf{z}|\mathbf{y}})$$

$$\boldsymbol{\Sigma}_{\mathbf{z}|\mathbf{y}}^{-1} = \boldsymbol{\Sigma}_{\mathbf{z}}^{-1} + \mathbf{W}^{\top} \boldsymbol{\Sigma}_{\mathbf{y}}^{-1} \mathbf{W}$$

$$\boldsymbol{\mu}_{\mathbf{z}|\mathbf{y}} = \boldsymbol{\Sigma}_{\mathbf{z}|\mathbf{y}} \left[\mathbf{W}^{\top} \boldsymbol{\Sigma}_{\mathbf{y}}^{-1} (\mathbf{y} - \mathbf{b}) + \boldsymbol{\Sigma}_{\mathbf{z}}^{-1} \boldsymbol{\mu}_{\mathbf{z}} \right]$$

- Gaussian prior $p(\mathbf{z})$, combined with the Gaussian likelihood $p(\mathbf{y} | \mathbf{z})$, results in a Gaussian posterior $p(\mathbf{z} | \mathbf{y})$.
- Thus Gaussians are closed under Bayesian conditioning.
- The Gaussian prior is a conjugate prior for the Gaussian likelihood.

Example: Inferring an unknown scalar

- Assume we have one noisy measurement y for an unknown quantity z .
- Prior: $p(z) = \mathcal{N}(z \mid \mu_0, \Sigma_0)$.
- Likelihood: $p(y \mid z) = \mathcal{N}(y \mid z, \Sigma_y)$.
- Posterior:

$$\begin{aligned} p(z \mid y) &= \mathcal{N}(z \mid \mu_1, \Sigma_1) \\ \Sigma_1 &= \left(\frac{1}{\Sigma_0} + \frac{1}{\Sigma_y} \right)^{-1} = \frac{\Sigma_y \Sigma_0}{\Sigma_0 + \Sigma_y} \\ \mu_1 &= \Sigma_1 \left(\frac{\mu_0}{\Sigma_0} + \frac{y}{\Sigma_y} \right) \\ &= y - (y - \mu_0) \frac{\Sigma_y}{\Sigma_y + \Sigma_0} \quad (\text{shrinkage}) \end{aligned}$$

- strong prior \rightarrow large shrinkage.