

0.1 Law of Total Expectation

Suppose X and Y are random variables over the same probability space (i.e., the events or outcomes happens simultaneously, e.g., gender=girl, boy, and haircolor=black, brown, other happens at the same time), then the Law is

$$\mathbf{E}(X) = \mathbf{E}_Y(\mathbf{E}_{X|Y}(X|Y)) \quad (1)$$

0.1.1 Proof in the discrete case

$$\begin{aligned} \mathbf{E}(X) &= \mathbf{E}_Y\left(\sum_x xP(X = x|Y)\right) \quad \text{weighted sum of the random outcomes of } x \\ &= \sum_y \left(\sum_x xP(X = x|Y = y)\right)P(Y = y) \\ &= \sum_x x\left(\sum_y P(X = x|Y = y)P(Y = y)\right) \\ &= \sum_x xP(X = x) \end{aligned} \quad (2)$$

0.2 Law of Total Variance (conditional variance formula)

$$\mathbf{Var}(Y) = \mathbf{E}(\mathbf{Var}(Y|X)) + \mathbf{Var}(\mathbf{E}(Y|X)) \quad (3)$$

0.2.1 Proof

$$\mathbf{Var}(Y) = \mathbf{E}(Y^2) - \mathbf{E}^2(Y) \quad (4)$$

0.3 Kernel Density Estimation

Draw n samples randomly from some unknown distribution f , the samples (x_1, \dots, x_n) are thus iid. The shape of f is estimated by using the n samples

$$\hat{f}_h(x) = \frac{1}{n} \sum_{i=1}^n K_h(x - x_i) = \frac{1}{hn} \sum_{i=1}^n K\left(\frac{x - x_i}{h}\right), \quad (5)$$

where K is a nonnegative kernel that integrates to one and has mean zero, $K_h(x) = \frac{1}{h}K(\frac{x}{h})$, and h is a positive smoothing parameter called the bandwidth. Notice that $\hat{f}_h(x)$ becomes a deterministic function once the n samples are determined, but changes as new n samples are drawn. If K is a standard normal kernel with $z = (x - x_i)/h \sim \mathcal{N}(0, 1)$, x is treated as a random variable with mean x_i , and standard deviation h .

Finding the best bandwidth

The mean integrated squared error (mean of the integrated squared error over multiple batches of n samples, where each batch gives an estimation of $\hat{f}_h(x)$)

$$\text{MISE}(h) = \mathbf{E} \int (\hat{f}_h(x) - f(x))^2 dx. \quad (6)$$

0.4 Subgrid-scale Parametrization with CMC (Conditional Markov Chain)

The Markov chain is a stochastic process/sequence indexed by t with the Markov property. The multi-dimensional state X has N_x outcomes indexed by i and j as $X^i(t)$ and $X^j(t + dt)$, with a corresponding N_b outcomes subgrid-scale parameter $B^n(t)$ and $B^m(t + dt)$. The i, j, n and m are the indices for the

uniformly separated intervals of the domain of X and B . Suppose $i = \{1, \dots, N_x\}$, and $n = \{1, \dots, N_B\}$. At fixed i and j , the stochastic transition matrix of size $N_B \times N_B$ is

$$\mathbf{P}^{ij} = P(B_j^m \mid B_i^n, X^i, X^j), \quad (7)$$

therefore there are N_x^2 numbers of matrices. Given the initial stochastic row vector of size N_x with elements summed to one, we can multiply the transition matrix to get the next stochastic vector. The i th index of the vector indicates the probability of ending up at the i th state.

Example:

$$[0.3 \ 0.7] \begin{pmatrix} 0.2 & 0.8 \\ 0.3 & 0.7 \end{pmatrix} \quad (8)$$

The probability of ending up at the 1st state equals $0.3 \times 0.2 + 0.7 \times 0.3$, which is the transition probability of 1 to 1, and 2 to 1 added.

Instead of using the transition matrix to calculate the final probability vector, we can use it to sample the states to obtain realizations of stochastic processes.

Example:

Suppose one given the triplet (i, n, j) , the m th state of B is sampled according to the n th row stochastic vector of the transition matrix \mathbf{P}^{ij} . If m is sampled at $m = 2$, then the next triplet $(j, 2, k)$, where $X^k(t+dt)$ is determined from the dynamic system, samples at the 2nd row of the transition matrix \mathbf{P}^{jk} .

0.5 Hidden Markov Model

0.6 Correlation

Definition. The linear relationship between two vectors/variables x and y

$$\text{corr}(x, y) = \frac{\mathbf{E}((x - \bar{x})(y - \bar{y}))}{\sigma(x)\sigma(y)}. \quad (9)$$

Derivation. The geometric intuition is by treating the rhs as the dot product of two unit vectors

$$\frac{x - \bar{x}}{|x - \bar{x}|} \cdot \frac{y - \bar{y}}{|y - \bar{y}|} = \cos \theta, \quad (10)$$

which is just the cosine of the angle between the two vectors x' and y' . This is seen by dividing both side by the unit vector squared $x' \cdot x' = 1$,

$$\frac{x' \cdot y'}{x' \cdot x'} = \frac{\cos \theta}{x' \cdot x'} = \cos \theta \quad (11)$$

and the lhs indicates the cosine is just a linear regression coefficient, i.e., $y' = \cos \theta x'$.