

Probability Theory Review.

• Probability Space:

(1) triple: (Ω, \mathcal{F}, P)

Ω : outcome space

\mathcal{F} : $\mathcal{F} \subseteq \mathcal{P}(\Omega)$ event space

P : probability measure: $E \in \mathcal{F} \rightarrow P(E) \in [0, 1]$

(2) restriction of \mathcal{F}

*: not all the event space has measure.

① $\Omega \in \mathcal{F}, \emptyset \in \mathcal{F}$

② \mathcal{F} is closed under (countable) unions.

③ \mathcal{F} is closed under complement.

(3) Properties of \mathcal{F} :

① $\forall A \in \mathcal{F}: P(A) \geq 0$

② $P(\Omega) = 1$

③ $\forall A, B \in \mathcal{F}, A \cap B = \emptyset: P(A \cup B) = P(A) + P(B)$

• Random Variables:

(1) not variables \rightarrow functions.

outcome $\rightarrow X \in \mathbb{R}$.

(2) indicator variable

*: the difference between Ω and \mathcal{F} .

◦ Joint distribution & Marginal Distribution.

(1) Joint distribution:

$$P(X=a, Y=b) = \delta.$$

(2) Marginal distribution:

$$P(X) = \sum_{b \in \text{Val}(Y)} P(X, Y=b)$$

◦ Conditional Distributions:

$$(1) P(X=a | Y=b) = \frac{P(X=a, Y=b)}{P(Y=b)}$$

knowing some events are true.

(2) $P(X|Y) \rightarrow$ knowing Y .

◦ Independence:

(1) machine learning:

data \rightarrow independent.

$$P(X) = P(X|Y)$$

$$P(X, Y) = P(X) \cdot P(Y)$$

(2) conditional independence:

$$P(X, Y | Z) = P(X | Z) \cdot P(Y | Z)$$

*: Naive Bayes.

◦ Chain Rule & Bayes Rule

(1) Chain Rule:

$$P(X_1, X_2, \dots, X_n) = P(X_1) \cdot P(X_2 | X_1) \cdot P(X_3 | X_1, X_2) \cdots P(X_n | X_1, X_2, \dots, X_{n-1})$$

*: calculating joint probability.

(2) Bayes Rule:

$$P(X|Y) = \frac{P(Y|X) \cdot P(X)}{P(Y)}$$

not knowing $P(Y)$:

$$P(Y) = \sum_x P(Y|X) P(X) \quad (\text{total probability})$$

◦ Probability Distribution:

(1) discrete & continuous

unified: measure theory.

(2) discrete:

probability mass function.

$$P(X=a) = p_i \quad \sum p_i = 1.$$

(3) continuous:

probability density function.

f : non-negative, integrable.

$$\int f(x) dx = 1.$$

$$P(a \leq X \leq b) = \int_a^b f(x) dx.$$

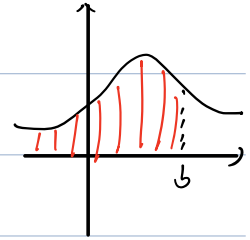
$$P(X=c) = \int_c^c f(x) dx = 0.$$

e.g. uniform distribution over $[a, b]$

$$f(x) = \begin{cases} \frac{1}{b-a} & x \in [a, b] \\ 0 & \text{otherwise.} \end{cases}$$

cumulative distribution function:

$$F(b) = P(X \leq b) = \int_{-\infty}^b f(x) dx$$



(4) Joint distribution:

$$P(a_1 \leq x_1 \leq b_1, a_2 \leq x_2 \leq b_2) = \int_{a_1}^{b_1} \int_{a_2}^{b_2} f(x_1, x_2) dx_1 dx_2.$$

(5) Conditional distribution:

$$f(y|x) = \frac{f(x, y)}{f(x)}$$

$$\text{e.g. } P(a \leq Y \leq b | X=c) = \int_a^b \frac{f(x, y)}{f(x)} dy = \int_a^b \frac{f(c, y)}{f(c)} dy.$$

• Expectations & Variance:

(1) discrete: $E(X) = \sum x P(X=a)$ * : first moment

continuous: $E(X) = \int_{-\infty}^{+\infty} x f(x) dx$

(2) linearity:

$$E(ax + by) = aE(x) + bE(y)$$

(3) Independent:

$$E(XY) = \sum_x \sum_y P(X=a, Y=b) xy$$

$$= \sum_x P(X=a) \cdot x \cdot \sum_y P(Y=b) y$$

$$= E(X) \cdot E(Y) \Rightarrow E(XY) = E(X)E(Y)$$

(4) Variance: * : second moment.

$$\text{Var}(X) = E[(X - E(X))^2] = \sigma^2 \Rightarrow \sigma = \sqrt{\text{Var}(X)}$$

$$\text{Var}(X) = E(X^2) - E^2(X)$$

$$\text{Var}(ax+b) = a^2 \text{Var}(x).$$

(5) Independent:

$$\text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y)$$

Covariance: (closely related)

$$\text{Cov}(X, Y) = E(X - E(X))(Y - E(Y))$$

• Important Distribution:

(1) Bernoulli:

$$P(X=k) = \binom{n}{k} p^k (1-p)^{n-k}$$

*: 2-classification tasks.

logistic regression

(2) Poisson:

fixed arrival rate λ

$$P(X=k) = e^{-\lambda} \frac{\lambda^k}{k!}$$

(3) Gaussian: *: normal distribution.

approximate \rightarrow binomial distribution.

*: noise \rightarrow Gaussian white noise.

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

multi-variate: (μ, Σ)

$\mu \in \mathbb{R}^k$ Σ : covariance matrix $\in \mathbb{R}^{k \times k}$

$$\Sigma_{ij} = \text{Cov}(X_i, X_j).$$

$$f(x) = \frac{1}{\sqrt{2\pi^k |\Sigma|}} \exp\left(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

• Efficient Manipulation:

(1) log trick:

product \rightarrow sum.

*: Likelihood Function:

$$L(\theta) = \prod_{i=1}^m (h_{\theta}(x^{(i)})^{y^{(i)}} [1 - h_{\theta}(x^{(i)})]^{1-y^{(i)}})$$

$$\log(L(\theta)) = \sum_{i=1}^m y^{(i)} \log(h_{\theta}(x^{(i)})) + \sum_{i=1}^m [1 - y^{(i)}] \log(1 - h_{\theta}(x^{(i)}))$$

(2) delayed normalization.

(3) Jensen's Inequality:

f : convex function.

$$f(E(x)) \leq E(f(x))$$

*: bound \rightarrow exact value.