

Robotics

Probability Basics

Marc Toussaint U Stuttgart

Probability Theory

• Why do we need probabilities?

Probability Theory

• Why do we need probabilities?

- Obvious: to express inherent stochasticity of the world (data)

Probability Theory

- Why do we need probabilities?
 - Obvious: to express inherent stochasticity of the world (data)
- But beyond this: (also in a "deterministic world"):
 - lack of knowledge!
 - hidden (latent) variables
 - expressing uncertainty
 - expressing information (and lack of information)
- Probability Theory: an information calculus

Probability: Frequentist and Bayesian

- Frequentist probabilities are defined in the limit of an infinite number of trials *Example:* "The probability of a particular coin landing heads up is 0.43"
- Bayesian (subjective) probabilities quantify degrees of belief *Example:* "The probability of it raining tomorrow is 0.3"
 - Not possible to repeat "tomorrow"

Probabilities & Sets

- Sample Space/domain Ω , e.g. $\Omega = \{1, 2, 3, 4, 5, 6\}$
- **Probability** $P: A \subset \Omega \mapsto [0, 1]$ e.g., $P(\{1\}) = \frac{1}{6}$, $P(\{4\}) = \frac{1}{6}$, $P(\{2, 5\}) = \frac{1}{3}$,
- Axioms: $\forall A, B \subseteq \Omega$
 - Nonnegativity $P(A) \ge 0$
 - Additivity $P(A \cup B) = P(A) + P(B)$ if $A \cap B = \emptyset$
 - Normalization $P(\Omega) = 1$
- Implications

$$0 \le P(A) \le 1$$

$$P(\emptyset) = 0$$

$$A \subseteq B \Rightarrow P(A) \le P(B)$$

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

$$P(\Omega \setminus A) = 1 - P(A)$$

Probabilities & Random Variables

• For a random variable X with discrete domain dom $(X) = \Omega$ we write:

```
 \forall_{x \in \Omega} : 0 \le P(X = x) \le 1  \sum_{x \in \Omega} P(X = x) = 1
```

Example: A dice can take values $\Omega = \{1, .., 6\}$. *X* is the random variable of a dice throw. $P(X=1) \in [0, 1]$ is the probability that *X* takes value 1.

• A bit more formally: a random variable relates a measureable space with a domain (sample space) and thereby introduces a probability measure on the domain ("assigns a probability to each possible value")

Probabilty Distributions

P(X=1) ∈ ℝ denotes a specific probability
 P(X) denotes the probability distribution (function over Ω)

Probabilty Distributions

• $P(X=1) \in \mathbb{R}$ denotes a specific probability P(X) denotes the probability distribution (function over Ω)

Example: A dice can take values $\Omega = \{1, 2, 3, 4, 5, 6\}$. By P(X) we discribe the full distribution over possible values $\{1, ..., 6\}$. These are 6 numbers that sum to one, usually stored in a *table*, e.g.: $[\frac{1}{6}, \frac{1}{6}, \frac{1}{6$

- In implementations we typically represent distributions over discrete random variables as tables (arrays) of numbers
- Notation for summing over a RV: In equation we often need to sum over RVs. We then write

 $\sum_{X} P(X) \cdots$ as shorthand for the explicit notation $\sum_{x \in \text{dom}(X)} P(X=x) \cdots$

Joint distributions

Assume we have *two* random variables X and Y

$$P(X = x, Y = y)$$



y

Definitions:

Joint: P(X, Y)Marginal: $P(X) = \sum_{Y} P(X, Y)$ Conditional: $P(X|Y) = \frac{P(X,Y)}{P(Y)}$

The conditional is normalized: $\forall_Y : \sum_X P(X|Y) = 1$

- X is *independent* of Y iff: P(X|Y) = P(X)(table thinking: all columns of P(X|Y) are equal)
- The same for *n* random variables $X_{1:n}$ (stored as a rank *n* tensor) Joint: $P(x_{1:n})$, Marginal: $P(X_1) = \sum_{X_{2:n}} P(X_{1:n})$, Conditional: $P(X_1|X_{2:n}) = \frac{P(X_{1:n})}{P(X_{2:n})}$

Joint distributions

joint: P(X, Y)marginal: $P(X) = \sum_{Y} P(X, Y)$ conditional: $P(X|Y) = \frac{P(X,Y)}{P(Y)}$

• Implications of these definitions: *Product rule:* P(X,Y) = P(X|Y) P(Y) = P(Y|X) P(X)

Bayes' Theorem
$$P(X|Y) = \frac{P(Y|X) P(X)}{P(Y)}$$

• The same for n variables, e.g., (X, Y, Z):

$$P(X_{1:n}) = \prod_{i=1}^{n} P(X_i | X_{i+1:n})$$
$$P(X_1 | X_{2:n}) = \frac{P(X_2 | X_1, X_{3:n}) P(X_1 | X_{3:n})}{P(X_2 | X_{3:n})}$$

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

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

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

Bayes' Theorem

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

$$\mathsf{posterior} = rac{\mathsf{likelihood} \cdot \mathsf{prior}}{\mathsf{normalization}}$$

Distributions over continuous domain

• Let *X* be a continuous RV. The **probability density function (pdf)** $p(x) \in [0, \infty)$ defines the probability

$$P(a \le X \le b) = \int_a^b p(x) \ dx \ \in [0, 1]$$

The (cumulative) probability distribution

 $F(x) = P(X \le x) = \int_{-\infty}^{x} dx \ p(x) \in [0, 1]$ is the cumulative integral with $\lim_{x\to\infty} F(x) = 1$.

(In discrete domain: probability distribution and probability mass function $P(X) \in [0, 1]$ are used synonymously.)

• Two basic examples:

 $\begin{array}{ll} \mbox{Gaussian:} & \mathcal{N}(x \,|\, a, A) = \frac{1}{|2\pi A|^{1/2}} \,\, e^{-\frac{1}{2}(xa)^\top \,A^{-1} \,\, (xa)} \\ \mbox{Dirac or } \delta \mbox{ ("point particle")} & \delta(x) = 0 \mbox{ except at } x = 0, \, \int \delta(x) \,\, dx = 1 \\ \delta(x) = \frac{\partial}{\partial x} H(x) \mbox{ where } H(x) = [x \geq 0] = \mbox{Heavyside step function.} \end{array}$

Gaussian distribution

• 1-dim:
$$\mathcal{N}(x \mid \mu, \sigma^2) = \frac{1}{|2\pi\sigma^2|^{1/2}} e^{-\frac{1}{2}(x-\mu)^2/\sigma^2}$$

$$N(x|\mu,\sigma^2)$$

• *n*-dim:
$$\mathcal{N}(x \mid a, A) = \frac{1}{|2\pi A|^{1/2}} e^{-\frac{1}{2}(x-a)^{\top} A^{-1} (x-a)}$$



Useful identities:

Mre identities: see "Gaussian identities" http://userpage.fu-berlin.de/~mtoussai/notes/gaussians.pdf

Particle Approximation of a Distribution

- We approximate a distribution p(x) over a continuous domain \mathbb{R}^n .
- A particle distribution q(x) is a weighed set of N particles $\{(x^i,w^i)\}_{i=1}^N$
 - each particle has a location $x^i \in \mathbb{R}^n$ and a weight $w^i \in \mathbb{R}$
 - weights are normalized $\sum_i w^i = 1$

$$q(x) := \sum_{i=1}^{N} w^{i} \delta(x - x^{i})$$

where $\delta(x - x^i)$ is the δ -distribution.

Particle Approximation of a Distribution



Particle Approximation of a Distribution

• For *q*(*x*) to approximate a given *p*(*x*) we want to choose particles and weights such that for any (smooth) *f*:

$$\lim_{N \to \infty} \langle f(x) \rangle_q = \lim_{N \to \infty} \sum_{i=1}^N w^i f(x^i) = \int_x f(x) p(x) dx = \langle f(x) \rangle_p$$

• How to do this? See An Introduction to MCMC for Machine Learning www.cs.ubc.ca/~nando/papers/mlintro.pdf

Some continuous distributions

Gaussian Dirac or δ Student's t (=Gaussian for $\nu \rightarrow \infty$, otherwise heavy tails)

$$\mathcal{N}(x \mid a, A) = \frac{1}{|2\pi A|^{1/2}} e^{-\frac{1}{2}(xa)^{\top} A^{-1} (xa)}$$
$$\delta(x) = \frac{\partial}{\partial x} H(x)$$
$$p(x; \nu) \propto [1 + \frac{x^2}{\nu}]^{-\frac{\nu+1}{2}}$$

Exponential (distribution over single event time)

Laplace ("double exponential")

Chi-squared

Gamma

$$p(x;\lambda) = [x \ge 0] \ \lambda e^{-\lambda x}$$

$$p(x; \mu, b) = \frac{1}{2b} e^{-|x-\mu|/b}$$

$$\begin{split} p(x;k) &\propto [x \geq 0] \; x^{k/2-1} e^{-x/2} \\ p(x;k,\theta) &\propto [x \geq 0] \; x^{k-1} e^{-x/\theta} \end{split}$$