

# **Probabilistic Modeling**

## Today

- Discrete random variables
- Continuous random variables
- P.d.f.'s and c.d.f.'s
- Mean and variance
- Dependence and independence; joint and marginal probabilities

# What is/why probabilistic modeling?

## What is a random variable?

- Something that has not happened yet.
    - Does a tossed coin come up heads or tails?
    - Does the cancer recur or not?
  - Something you do not know ...
    - Did a tossed coin come up heads or tails?
    - Is X a transcription factor for gene Y?
    - How does the protein fold?
- ... because you have not/cannot observe it directly or compute it definitively from what you have observed.

## Discrete random variables

## Examples

A discrete r.v.  $X$  takes values from a discrete set  $\Omega_X$ .

- $X$  = result of a coin toss;  $\Omega_X = \{\text{Head, Tail}\}$ .
- $X$  = roll of a die;  $\Omega_X = \{1, 2, 3, 4, 5, 6\}$ .
- $X$  = nucleotide a position 1, chromosome 1, in a particular person;  $\Omega_X = \{A, C, G, T\}$ .
- $X$  = amino acid 12 in a particular person's hemoglobin;  $\Omega_X = \{A, R, N, D, C, Q, E, G, H, I, L, K, M, F, P, S, T, W, Y, V\}$ .
- $X$  = copy number of gene Z in a particular person;  $\Omega_X = \{0, 1, 2, 3, \dots\}$ .

## Probabilities

- For a discrete r.v.  $X$ , each value  $x \in \Omega_X$  has a probability of occurring, denoted variously by

$$\begin{array}{ccc} \text{Prob}(X = x) & \text{Prob}_X(x) & \text{Prob}(x) \\ \text{Pr}(X = x) & \text{Pr}_X(x) & \text{Pr}(x) \\ \mathbf{P}(X = x) & \mathbf{P}_X(x) & \mathbf{P}(x) \end{array}$$

- $0 \leq \mathbf{P}(x) \leq 1$
- $\sum_{x \in \Omega_X} \mathbf{P}(x) = 1$
- $\mathbf{P}(X)$  denotes the *probability distribution function* for r.v.  $X$ . It can be thought of as a table.

$x$	A	C	G	T
$\mathbf{P}(x)$	0	0.2	0.7	0.1

## Cumulative distribution functions

- If  $X$  takes values from an ordered set  $\Omega_X$  (such as integers) then the *cumulative distribution function* is

$$\text{c.d.f.}(x) = P(X \leq x) = \sum_{x' \leq x} P(x')$$

- For example, if  $X$  is the roll of a die, then:

$x$	1	2	3	4	5	6
$P(x)$	1/6	1/6	1/6	1/6	1/6	1/6
c.d.f.( $x$ )	1/6	2/6	3/6	4/6	5/6	1



## Mean and variance

- If  $\Omega_X$  is a set of numbers, then the expected value of  $X$  is

$$\mathbf{E}(X) = \sum_{x \in \Omega_X} x\mathbf{P}(x)$$

- The variance of  $X$  is

$$\begin{aligned}\mathbf{Var}(X) &= \mathbf{E}(X^2) - (\mathbf{E}(X))^2 \\ &= \left( \sum_x x^2 \mathbf{P}(x) \right) - \left( \sum_x x \mathbf{P}(x) \right)^2 \\ &\geq 0\end{aligned}$$

- Example: If  $X$  is a die roll, then the mean value is 3.5 and the standard deviation is approximately 3.4157.

## Continuous random variables

## Examples

A continuous r.v.  $X$  takes real values.

- $X$  = expression level for a gene as reported by a microarray.
- $X$  = time until a patient's cancer recurs.
- $X$  = size of a tumor.
- $X$  = mass of a peptide as reported by mass-spec.
- $X$  = binding energy between a TF and DNA. (?)
- $X$  = fraction of time a TF is bound to DNA.

## Cumulative distribution functions

- Any continuous r.v.  $X$  has a *cumulative distribution function*

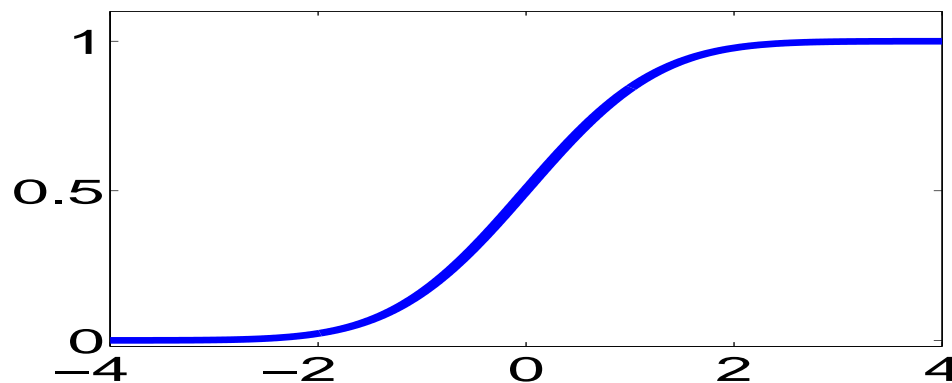
$$\text{c.d.f.}(x) = \text{P}(X \leq x)$$

- $\text{c.d.f.}(x)$  is a non-decreasing function;  $\text{c.d.f.}(x) \leq \text{c.d.f.}(x')$  whenever  $x \leq x'$ .

- $\lim_{x \rightarrow -\infty} \text{c.d.f.}(x) = 0$ .

- $\lim_{x \rightarrow +\infty} \text{c.d.f.}(x) = 1$ .

- Example: The c.d.f. of a mean-zero, variance-one Gaussian r.v.:



## Probability density functions

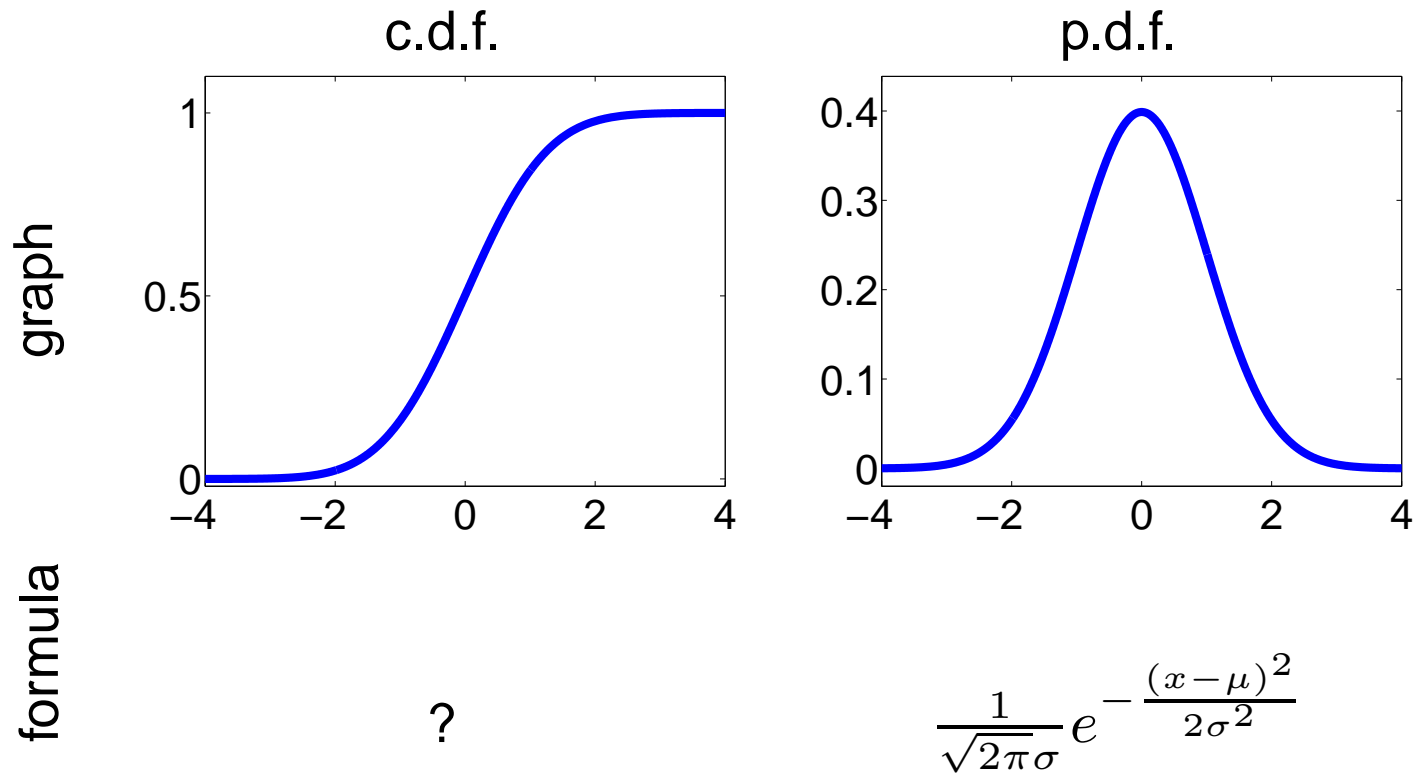
- If c.d.f.  $(x)$  is continuous and differentiable (at least, in most places) then its derivative is the *probability density function*, analogous to the probability *distribution* function of a discrete r.v.

$$\frac{d}{dx} \text{c.d.f.}(x) = \text{p.d.f.}(x) = P(x)$$

- $P(x)$  is the “probability”, or more properly, likelihood that  $X$  takes value  $x$ .
- $0 \leq P(x) < \infty$ . Observe that  $P(x) > 1$  is allowed, unlike for discrete r.v.'s.
- $\int_x P(x) dx = 1$ , similar to discrete r.v.'s.

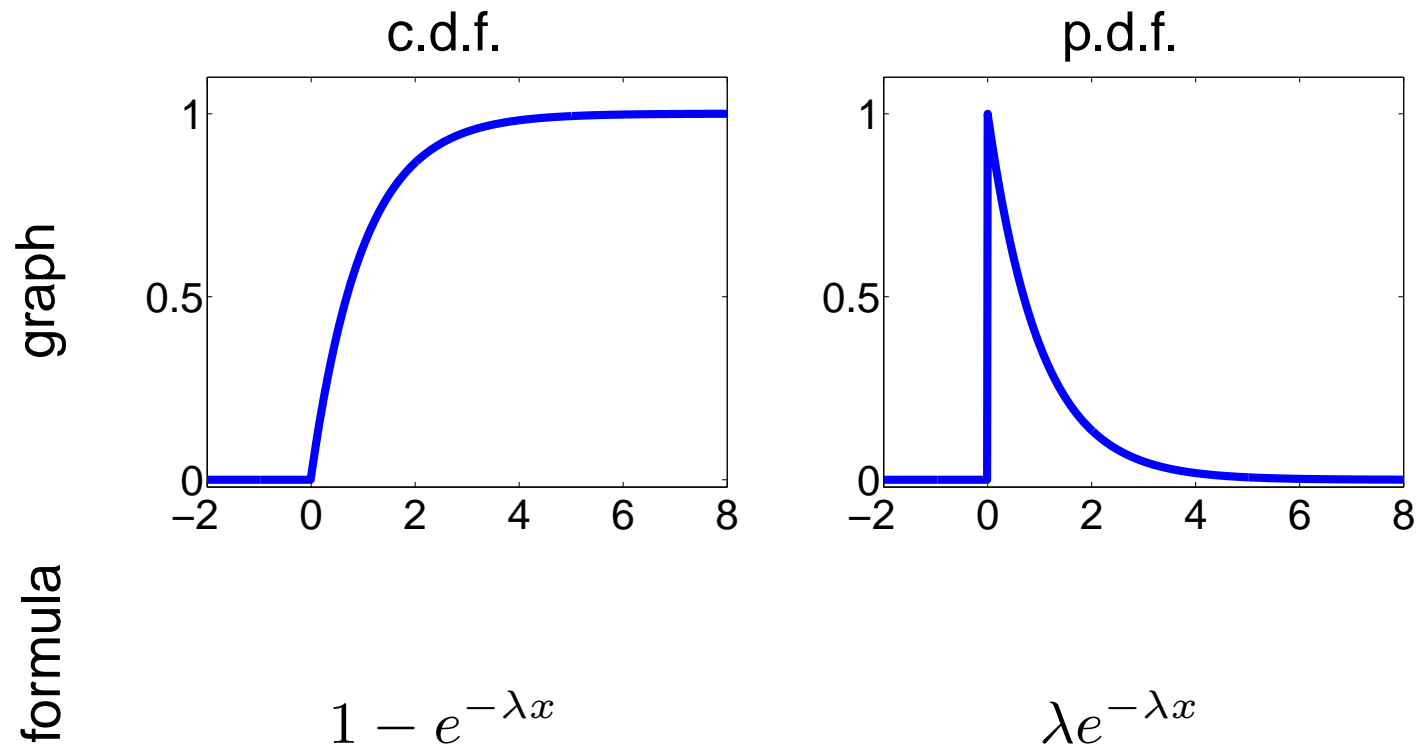
## Gaussian random variables

$X \sim N(\mu, \sigma)$  has mean  $\mu$  and standard deviation  $\sigma$ .



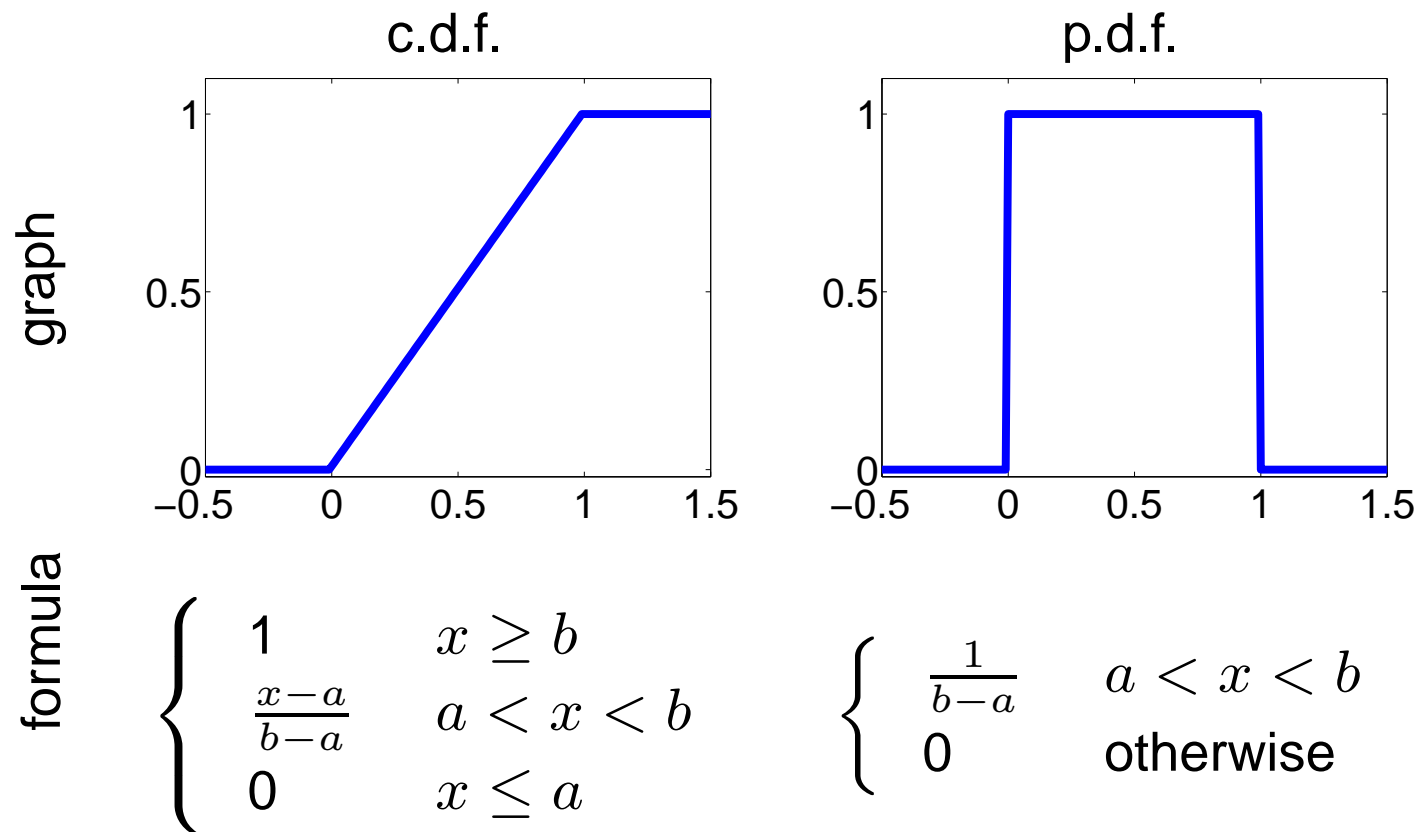
## Exponential random variables

$X \sim \text{Exp}(\lambda)$  has mean  $1/\lambda$  and standard deviation  $1/\lambda$ .



## Uniform random variables

$X \sim U(a, b)$  has mean  $\frac{a+b}{2}$  and standard deviation  $\frac{(b-a)}{\sqrt{12}}$ .



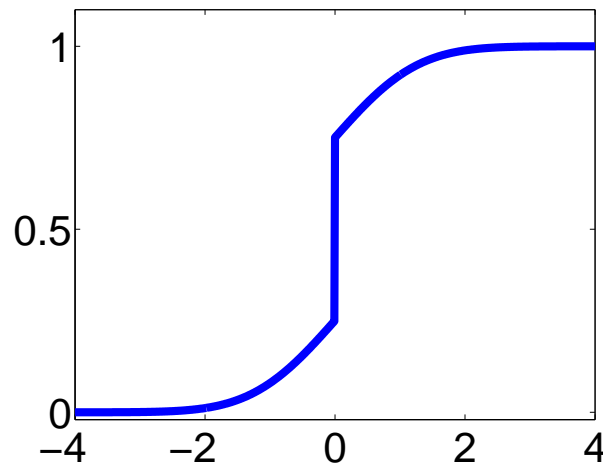


## A continuous r.v. with no p.d.f.

- Suppose  $X$  equal to zero with probability  $\frac{1}{2}$  and otherwise is distributed according to  $N(0, 1)$ .

- Then the c.d.f. is
$$\text{c.d.f.}(x) = \begin{cases} \frac{1}{2}f(x) & x < 0 \\ \frac{1}{2}f(x) + \frac{1}{2} & x \geq 0 \end{cases}$$

where  $f(x)$  denotes the c.d.f. of a  $N(0, 1)$  r.v.



- There is no p.d.f. because of the discrete jump in the c.d.f.

## Mean and variance

- We will almost always restrict attention to continuous r.v.'s with p.d.f.'s.
- Then, the expected value is defined as

$$\mathbf{E}(X) = \int_x x\mathbf{P}(x)dx$$

- Variance is

$$\begin{aligned}\mathbf{Var}(X) &= \mathbf{E}(X^2) - (\mathbf{E}(X))^2 \\ &= \int_x x^2\mathbf{P}(x)dx - \left(\int_x x\mathbf{P}(x)dx\right)^2 \\ &\geq 0\end{aligned}$$

## [In]dependent random variables

## Example

- Let  $X_1 = \text{true}$  iff a rolled die comes out even.
- Let  $X_2 = \text{true}$  iff the same rolled die comes out odd.

$$P(X_1 = \text{true}) = P(X_1 = \text{false}) = \frac{1}{2}$$

$$P(X_2 = \text{true}) = P(X_2 = \text{false}) = \frac{1}{2}$$

- What is the probability  $P(X_1 = \text{true} \text{ and } X_2 = \text{true})$ ?

## Example

- Let  $X_1 = \text{true}$  iff a rolled die comes out even.
- Let  $X_2 = \text{true}$  iff the same rolled die comes out odd.

$$P(X_1 = \text{true}) = P(X_1 = \text{false}) = \frac{1}{2}$$

$$P(X_2 = \text{true}) = P(X_2 = \text{false}) = \frac{1}{2}$$

- What is the probability  $P(X_1 = \text{true} \text{ and } X_2 = \text{true})$ ?
  - We know it is zero.
  - But there is no way of knowing just from  $P(X_1)$  and  $P(X_2)$ .
- ⇒ There are several ways we can specify the relationships between variables. They all come down to specifying *joint probability distributions/densities*.

## Joint probabilities

- When considering r.v.'s  $X_1, X_2, \dots, X_m$ , the joint probability function specifies the probability of every combination of values.

$$P(X_1 = x_1 \text{ and } X_2 = x_2 \text{ and } \dots \text{ and } X_m = x_m)$$

- When the r.v.'s are discrete, the joint probability can be viewed as a table.

	even=true	odd=true
odd=true	0	1/2
odd=false	1/2	0

	die=1	2	3	4	5	6
even=true	0	1/6	0	1/6	0	1/6
even=false	1/6	0	1/6	0	1/6	0

## Marginal probabilities

Given r.v.'s  $X_1, X_2, \dots, X_m$  with joint probability  $P(x_1, x_2, \dots, x_m)$ .

- The marginal probability of a r.v.  $X_i$  is

$$P(X_i = x_i) = \sum_{x_1, x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_m} P(x_1, x_2, \dots, x_m)$$

- That is, you get the marginal probability by summing (or integrating) over all possible values of the other r.v.'s.

	die=1	2	3	4	5	6	P(even)
even=true	0	1/6	0	1/6	0	1/6	1/2
even=false	1/6	0	1/6	0	1/6	0	1/2
P(die)	1/6	1/6	1/6	1/6	1/6	1/6	

- Similarly for the marginal probability of a subset of the r.v.'s.

## Independent r.v.'s

- Two r.v.'s  $X$  and  $Y$  are independent if and only if

$$P(X = x \text{ and } Y = y) = P(X = x)P(Y = y)$$

for all  $x$  and  $y$ .

- This is often abbreviated as  $P(X, Y) = P(X)P(Y)$ .



## Conditional probability

- For two r.v.'s  $X$  and  $Y$ ,  $P(X = x|Y = y)$  denote the probability that  $X = x$  *given that*  $Y = y$ .
  - $P(\text{die}=1|\text{odd} = \text{true}) = 1/3$ .
  - $P(\text{die}=1|\text{odd} = \text{false}) = 0$ .
- The conditional probability can be defined (and computed) as

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

as long as  $P(y) > 0$ .

- This is sometimes used as

$$P(x) = \sum_y P(x, y) = \sum_y P(x|y)P(y)$$

## Conditional probability (2)

Conditional probabilities are interesting because we often observe something and want to infer something/make a guess about something unobserved but related.

- $P(\text{cancer recurs} | \text{tumor measurements})$
- $P(\text{TF binds} | \text{TF and DNA properties})$
- $P(\text{Gene expressed} > 1.3 | \text{TF concentrations})$

## Bayes' Rule

(Or possibly Bayes's Rule.)

- Bayes' Rule:  $P(x|y) = \frac{P(y|x)P(x)}{P(y)}$ .
- E.g., suppose we know based on past data collected:

$P(\text{tumor measurements}|\text{cancer})$

$P(\text{tumor measurements}|\text{not cancer})$

$P(\text{cancer})$

$P(\text{not cancer})$

$$P(\text{cancer}|\text{tumor meas.}) = \frac{P(\text{tumor meas.}|\text{cancer})P(\text{cancer})}{P(\text{tumor meas.})}$$

$$= \frac{P(\text{tumor meas.}|\text{cancer})}{P(\text{tumor meas.}|\text{cancer})P(\text{cancer}) + P(\text{tumor meas.}|\text{not cancer})P(\text{not cancer})}$$