

CSCI 3022 Intro to Data Science

Continuous Random Variables

Announcement and Reminders:

- ▶ HW3 Due tonight; HW4 posted tomorrow.
- ▶ Past Exams and solutions posted to Canvas!
- ▶ Exam is untimed, takehome. Designed to take around 2-3 hours, but you'll have most of a week. Due Mar 8, **likely** posted Mar 2. Think of it as an "all pen-and-paper" homework. Submission via Gradescope:
 1. You can TeX it yourself (template will be provided)
 2. You can take the given pdf and annotate on it
 3. You can print the pdf and write on it, then picture/scan your solutions
 4. You can work the exam on regular paper and picture/scan your work
- ▶ Tested content through lecture Mar 1 (Expectation & Variance)

Last Time...: the blocks of discrete probability

1. Bernoulli: *one* binary outcome experiment.
2. Binomial: binary outcome experiment success *count* in n tries.
3. Geometric: Total trials *until a success* of a binary outcome experiment.
4. Negative Binomial: Trials until r binary outcome experiment *successes*.
5. Poisson: *counting* outcomes with a fixed rate λ .

Last Time...: the blocks of discrete probability

1. Bernoulli: *one* binary outcome experiment.

$$f(x) = p^x(1-p)^{1-x}$$

2. Binomial: binary outcome experiment success *count* in n tries.

$$f(x) = \binom{n}{x} p^x (1-p)^{(n-x)}$$

3. Geometric: Total trials *until a success* of a binary outcome experiment.

$$f(x) = (1-p)^{x-1} p$$

4. Negative Binomial: Trials until r binary outcome experiment *successes*.

$$f(x) = \binom{x-1}{r-1} p^r (1-p)^{(x-r)}$$

5. Poisson: *counting* outcomes with a fixed rate λ .

$$f(x) = \frac{e^{-\lambda} \lambda^x}{x!}$$

Last Time...: the blocks of discrete probability

The underlying pieces of discrete RVs:

1. The random variable X takes inputs/events in the (discrete) sample space Ω and maps them to a (discrete) finite or infinite set of probability values a_1, a_2, a_3, \dots .
2. We find probabilities in the probability mass function or probability density function

$$f(x) = P(X = x).$$

3. We can find cumulative probabilities or probability on ranges of outcomes in the cumulative density function

$$F(x) = P(X \leq x) = \sum_{X \leq x} f(x).$$

Continuous RVs

Many real-life random processes must be modeled by random variables that can take on continuous (non-discrete) values. Some example:

1. Peoples' heights: $X \in$
2. Final grades in a class: $X \in$
3. Time between people checking out at a store : $t \in$

Continuous RVs

Many real-life random processes must be modeled by random variables that can take on continuous (non-discrete) values. Some example:

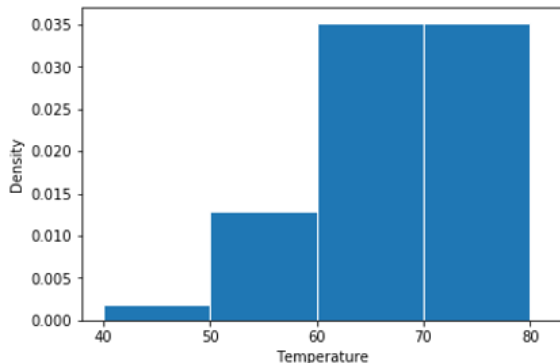
1. Peoples' heights: $X \in \{[0, 7.5ft]\}$
2. Final grades in a class: $X \in \{[0, 100]\}$
3. Time between people checking out at a store : $t \in \{[0, \infty]\}$

More precision!

Suppose your friend asks you what the high temperature will be today. They want to know the probability it will be between 70F and 80F, so they can wear shorts.

Probability:

Add up the share of outcomes between 70F and 80F!

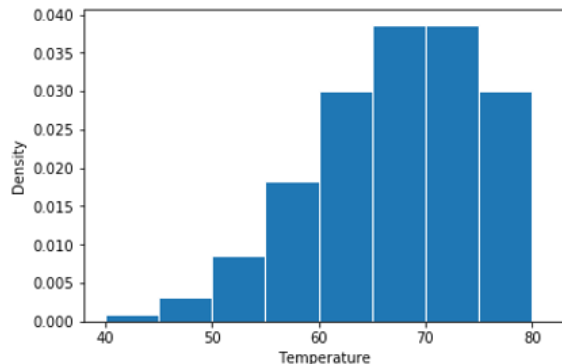


More precision!

Suppose your friend asks you what the high temperature will be today. They want to know the probability it will be between 70F and 80F, so they can wear shorts.

Probability:

Add up the share of outcomes between 70F and 80F!

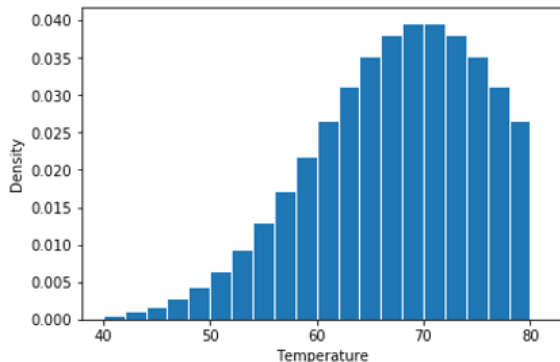


More precision!

Suppose your friend asks you what the high temperature will be today. They want to know the probability it will be between 70F and 80F, so they can wear shorts.

Probability:

Add up the share of outcomes between 70F and 80F!

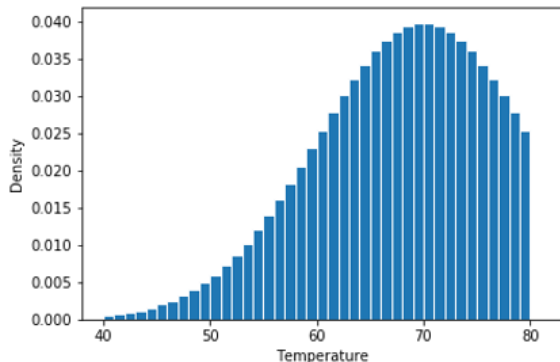


More precision!

Suppose your friend asks you what the high temperature will be today. They want to know the probability it will be between 70F and 80F, so they can wear shorts.

Probability:

Add up the share of outcomes between 70F and 80F!

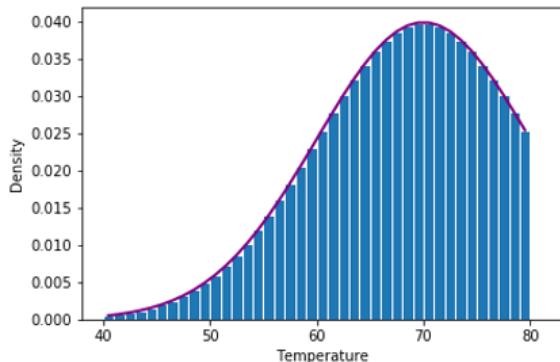


More precision!

Suppose your friend asks you what the high temperature will be today. They want to know the probability it will be between 70F and 80F, so they can wear shorts.

Probability:

Add up the share of outcomes between 70F and 80F!

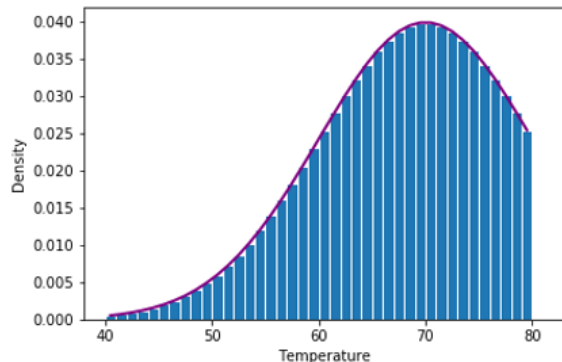


More precision!

Suppose your friend asks you what the high temperature will be today. They want to know the probability it will be between 70F and 80F, so they can wear shorts.

Probability:

Integrate up the share of outcomes between 70F and 80F!



Continuous Distributions

Example:

Consider the reference line connecting the valve stem on a tire to the center point.

Let X be the angle measured clockwise to the location of an imperfection. The pdf for X is:

$$f(x) = \begin{cases} \frac{1}{360} & 0 \leq X < 360 \\ 0 & \text{else} \end{cases}$$

Continuous Distributions

Example, cont'd:

$$f(x) = \begin{cases} \frac{1}{360} & 0 \leq X < 360 \\ 0 & \text{else} \end{cases}$$

Graphically, the pdf of X is:

Continuous Distributions

Example, cont'd: How can we show that:

1. the total area of the pdf of x is 1?
2. How do we calculate $P(90 \leq X \leq 180)$?
3. What is the probability that the angle of occurrence is within 90 of the reference line?
(The reference line is at 0 degrees.)

Continuous Distributions

Example, cont'd: How can we show that:

1. the total area of the pdf of x is 1?

$$\int_0^{360} f(x) dx = 1?$$

2. How do we calculate $P(90 \leq X \leq 180)$?

$$\int_{90}^{180} f(x) dx = \dots?$$

3. What is the probability that the angle of occurrence is within 90 of the reference line?
(The reference line is at 0 degrees.)

$$P(X < 90 \text{ OR } X > 270) = \int_0^{90} f(x) dx + \int_{270}^{360} f(x) dx = \dots?$$

Uniform Distribution

The previous problem was an example of the uniform distribution.

Definition: *Uniform Distribution*

A continuous rv X is said to have a *uniform distribution* on the interval $[a, b]$ if the pdf of X is:

NOTATION: We write _____ to indicate that X is a uniform rv with lower bound a and upper bound b .

Uniform Distribution

The previous problem was an example of the uniform distribution.

Definition: *Uniform Distribution*

A continuous rv X is said to have a *uniform distribution* on the interval $[a, b]$ if the pdf of X is:

$$f(x) = \frac{1}{b-a}; \quad x \in [a, b]$$

NOTATION: We write $X \sim U(a, b)$ to indicate that X is a uniform rv with lower bound a and upper bound b .

Exponential Distribution

The family of exponential distributions provides probability models that are very widely used in engineering and science disciplines to describe time-to-event data.

It can be thought of as a continuous analogue to the Poisson distribution, but instead of events-per-time, it measure time-per-events.

Examples:

Exponential Distribution

Definition: *Exponential Distribution*

A continuous rv X is said to have an *exponential distribution* with rate parameter λ if the pdf of X is:

NOTATION: We write _____ to indicate that X is an exponential rv with rate λ .

Exponential Distribution

Definition: *Exponential Distribution*

A continuous rv X is said to have an *exponential distribution* with rate parameter λ if the pdf of X is:

$$f(x) = \lambda e^{-\lambda x}; \quad x \geq 0$$

NOTATION: We write $X \sim \exp(\lambda)$ to indicate that X is an exponential rv with rate λ .

Exponential Distribution

Example:

Suppose a light bulb's lifetime is exponentially distributed with parameter $\lambda = 1/1000$.

1. What are the units for λ ?
2. What is the probability that the lifetime of the light bulb lasts less than 400 hours?
3. What is the probability that the lifetime of the light bulb lasts more than 5 hours?

Exponential Distribution

Example:

Suppose a light bulb's lifetime is exponentially distributed with parameter $\lambda = 1/1000$.

1. What are the units for λ ?

Same as Poisson: outcomes per time; so maybe burnouts per hour?

2. What is the probability that the lifetime of the light bulb lasts less than 400 hours?

$$P(X < 400) = \int_0^{400} \lambda e^{-\lambda x} = -e^{-\lambda x} \Big|_0^{400} = 1 - e^{-400/1000} = 1 - e^{-2/5}$$

3. What is the probability that the lifetime of the light bulb lasts more than 5 hours?

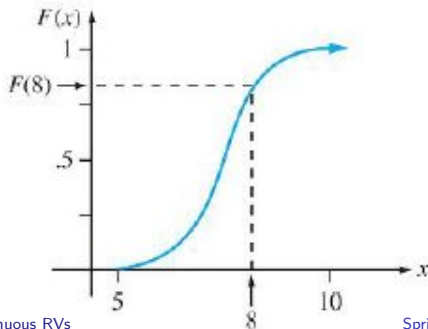
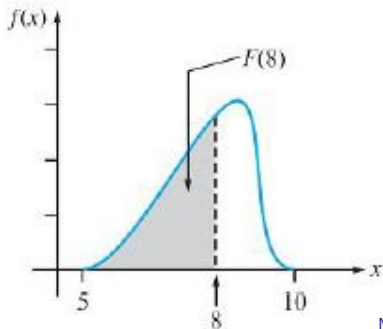
$$P(X > 5) = \int_5^{\infty} \lambda e^{-\lambda x} = -e^{-\lambda x} \Big|_5^{\infty} = 0 - (-e^{-5/1000}) = e^{-1/200} \approx 1$$

Cumulative Density Function

Definition: *Cumulative Density Function*

The *cumulative distribution function* (cdf) is denoted with $F(x)$. For a continuous r.v. X with pdf $f(x)$, $F(x)$ is defined for every real number x by:

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



Continuous CDFs

Example:

The distribution of the amount of gravel (in tons) sold by a particular construction supply company in a given week is a continuous rv X with pdf

$$f(x) = \begin{cases} \frac{3}{2}(1 - x^2) & 0 \leq X < 1 \\ 0 & \text{else} \end{cases}$$

1. What is the cdf of sales for any x ?
2. Find the probability that X is less than .25?
3. X is greater than .75?
4. $P(.25 < X < .75)$?

Continuous CDFs

Example:

The distribution of the amount of gravel (in tons) sold by a particular construction supply company in a given week is a continuous rv X with pdf

$$f(x) = \begin{cases} \frac{3}{2}(1 - x^2) & 0 \leq X < 1 \\ 0 & \text{else} \end{cases}$$

1. What is the cdf of sales for any x ?

$$F(x) = P(X \leq x) = \int_0^x \frac{3}{2}(1 - t^2) dt$$

$$F(x) = \frac{3x}{2} - \frac{x^3}{2}$$

2. Find the probability that X is less than .25? $F(.25)$

3. X is greater than .75? $1 - F(.75)$

4. $P(.25 < X < .75)$? $F(.75) - F(.25)$

Continuous CDFs

Wait, we've seen this before...

Recall: *The Fundamental Theorem of Calculus.*

Suppose F is an anti-derivative of f . Then:

1.

$$\frac{d}{dx} \int_a^x f(t) dt = f(x);$$

a.k.a.

$$\frac{d}{dx} F(x) = f(x);$$

2.

$$\int_a^b f(x) dx = F(B) - F(A).$$

Percentiles of a Distribution

Definition: The median \tilde{x} of a continuous distribution is the 50th percentile or .5 quantile of the distribution.

How can we express this in terms of $f(x), F(x)$?

Notation:

Visually:

Percentiles of a Distribution

Definition: The median \tilde{x} of a continuous distribution is the 50th percentile or .5 quantile of the distribution.

How can we express this in terms of $f(x)$, $F(x)$?

Notation:

\tilde{x} satisfies $F(\tilde{x}) = .5$, or

Visually:

$$.5 = \int_{-\infty}^{\tilde{x}} f(x) dx$$

Pops and Samples

Today marks the start of a large jump in how we approach data science problems:

1. We know about *sample statistics* like \bar{X} , s_X .
2. We've defined some *processes* that gives rise to distributions like the binomial, exponential, etc.
3. **Now:** we start bridging the gap! Given data and sample statistics, how do we estimate or infer properties of the underlying reality process? (parameters like p , λ).

To do this, we need an understanding of centrality and dispersion of a process or density function might be.

Mean/Expected Value

Example:

Consider a university having 15,000 students and let X equal the number of courses for which a randomly selected student is registered.

The pdf of X is given to you as follows:

x	1	2	3	4	5	6	7
$f(x) = P(X = x)$.01,	.03	.13	.25	.39	.17	.02

Students pay more money when enrolled in more courses, and so the university wants to know what the *average* number of courses students take per semester.

Mean/Expected Value

Definition: *Expected Value:*

For a discrete random variable X with pdf $f(x)$, the *expected* value or *mean* value of X is denoted as $E(X)$ and is calculated as:

Mean/Expected Value

Definition: *Expected Value:*

For a discrete random variable X with pdf $f(x)$, the *expected* value or *mean* value of X is denoted as $E(X)$ and is calculated as:

$$E[X] = \sum_{x \in \Omega} x \cdot P(X = x)$$

Mean/Expected Value

Example:, cont'd:

The pdf of X is given to you as follows:

x	1	2	3	4	5	6	7
$f(x) = P(X = x)$.01,	.03	.13	.25	.39	.17	.02

What is $E[X]$?

Mean/Expected Value

Example:, cont'd:

The pdf of X is given to you as follows:

x	1	2	3	4	5	6	7
$f(x) = P(X = x)$.01,	.03	.13	.25	.39	.17	.02

What is $E[X]$?

$$E[X] = \sum_{x \in \Omega} x \cdot P(X = x) = 1 \cdot .01 + 2 \cdot .03 + 3 \cdot .13 + 4 \cdot .25 + 5 \cdot .39 + 6 \cdot .17 + 7 \cdot .02$$

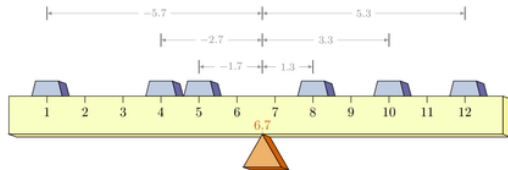
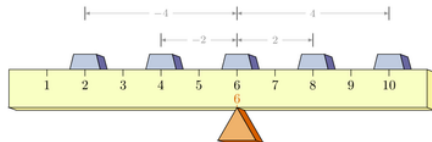
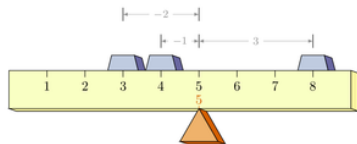
$$E[X] = 4.57$$

Interpreting Expected Value: Relative Frequency

One way to interpret expected value of a discrete distribution (especially on a finite support) is the sample mean if we managed to observe observations that *exactly* mirror the probability mass function.

In the preceding example, the pmf was given at 7 values of X with a precision up to 1%. In this case, if we had exactly 100 students and their proportions *observed* exactly mirrored the probabilities given in the example, the sample mean would be identical to the population mean.

Interpreting Expected Value



- ▶ The "center of mass" of a set of point masses
- ▶ Each mass exerts an " $r \times f$ " force on the balancing point.
- ▶ Same idea holds in continuous space: we're looking for a centroid of an object.

<http://www.texample.net/media/tikz/examples/TEX/balance.tex>

Daily Recap

Today we learned

1. Continuous pdfs!

Moving forward:

- nb day Friday!

Next time in lecture:

- More on “average values” of pdfs, plus calculating variances from them!