

Given Ω is chosen, a probability law on Ω is a mapping P that assigns a number to every event (i.e. to every subset of Ω) such that:

- 1) $P(A) \geq 0$ for every event A
- 2) $P(\Omega) = 1$ (normalization)
- 3) Additivity Rules
 - \Rightarrow If $A \cap B = \emptyset$ (A, B are events) then $P(A \cup B) = P(A) + P(B)$
 - \Rightarrow If A_1, A_2, A_3, \dots is a countable sequence of mutually disjoint events (i.e. $A_i \cap A_j = \emptyset \forall i, j$), then $P(A_1 \cup A_2 \cup \dots) = \sum_{i=1}^{\infty} P(A_i)$

So, given an event $A \subset \Omega$, $P(A)$ is a model for "the likelihood that the outcome of the uncertain experiment is in A "

"Event A occurs" means "outcome of experiment is in A "

Conditional Independence: Ω, P ; say events A and B are conditionally independent given (event) C when

$$P(A \cap B | C) = P(A|C)P(B|C)$$

$$P(A|B \cap C) = P(A|C); \quad P(B|A \cap C) = P(B|C)$$

Knowledge of B gives no functional info about probability of A on top of knowledge of C . To see this just play w/ formulas

$$P(A|B \cap C) = \frac{P(A \cap B \cap C)}{P(C)} = \frac{P(A|C)P(B|C)}{P(C)} = P(A|C)$$

Given a discrete r.v. X w/ $P_X(x)$ prob, define the expected value (or expectation)

$$E(X) = \sum_{x \in X} x P_X(x)$$

Next, multiple discrete rvs and joint pmfs, etc. Given Ω, P and two discrete rvs X, Y defined on Ω , define the joint pmf of X and Y as

$$P_{XY}(x,y) = P\left(\begin{array}{c} X=x \\ Y=y \end{array}\right) = P(X=x \cap Y=y)$$

Note: for any set V of possible value pairs for XY , we have

$$\sum_{(x,y) \in V} P_{XY}(x,y) = P(\text{event that } (X,Y) \in V)$$

Since X, Y are discrete rvs, they have their own pmfs P_X, P_Y . These are determined as follows from the joint pmf $P_{XY}(x,y)$:

$$① \quad P_X(x) = \sum_{y \in Y} P_{XY}(x,y) \quad \forall x$$

$$② \quad P_Y(y) = \sum_{x \in X} P_{XY}(x,y) \quad \forall y$$

Why are these true? by Total Probability Rule

$$P_X(x) \cdot P(Y=x) = P_X(x) = \sum_y P_{XY}(x,y) = P_{XY}(x,y) \forall x$$

These definitions generalize in an obvious way to > 2 rvs defined on same Ω, P .

KEY THING: - joint pmf determines the marginals
- marginals do NOT determine the joint

Recall the expected-value rules: Given $X, P_X(x), Y=g(X)$, have

$$E(X) = \sum_x g(x) P_X(x) \quad \leftarrow \text{Zeros be causes as from computing } P_X(x)$$

Similarly, given X, Y w/ joint pmf $P_{XY}(x,y)$ and some real-valued function $z = g(X, Y)$, we have

$$E(z) = \sum_{x,y} g(x,y) P_{XY}(x,y) \quad \leftarrow \text{don't need } P_Z(z) \quad \text{to get } E(z)$$

Next, Conditional Stuff

Given Ω, P and a discrete r.v. Y defined on Ω , and an event $A \subset \Omega$ w/ $P(A) > 0$, and a possible value x for X , the conditional pmf of X given A is defined as

$$P_{XA}(x) = \frac{P\left(\begin{array}{c} X=x \\ X \in A \end{array}\right)}{P(A)} = "P(B|A)" \text{ where } B \text{ is the event } \{x\}$$

Observe that for any A w/ $P(A) > 0$, $P_{XA}(x)$ as x ranges over X value space defines a pmf - i.e. $P_{XA}(x) \geq 0 \forall x$ and $\sum_x P_{XA}(x) = 1$

Here's a fact that's similar to (and follows directly from) the Total Probability Thm: If events A_1, A_2, \dots, A_n partition Ω , and $P(A_k) > 0$ for $1 \leq k \leq n$, then for any discrete r.v. X on Ω ,

$$P(x) = \sum_{k=1}^n P_{XA_k}(x) P(A_k)$$

More often, encounter conditional pmf of X given some other rv Y (defined on same Ω, P): given X, Y defined on Ω, P , conditional pmf of X given Y is defined for all x and for all y with $P(Y=y) > 0$: $P_{XY}(x,y) = P_{XY}(x|y)$ as

$$P_{XY}(x,y) = \frac{P_{XY}(x,y)}{P_Y(y)} \quad \leftarrow \text{Same as } P_{XA}(x) \text{ where } A = \{Y=y\}$$

Standard Notation

Note that for any y w/ $P_Y(y) > 0$, $P_{XY}(x|y)$ as x ranges over X values defines a pmf

$$\text{i.e. } P_{XY}(x|y) \geq 0 \text{ and } \sum_x P_{XY}(x|y) = 1$$

Given X, P_X , and $Y=g(X)$,

$$E(Y) = \sum_x g(x) P_X(x)$$

Conditional Variance

Given X, Y conditional variance of X given Y is the random variable

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

A recipe similar to the "g-thing" for computing $\text{Var}(X|Y)$

* Given g , compute $\text{Var}(X|Y=g) = E[(X - E(X|Y=g))^2 | Y=g]$

- Do this by finding conditional pmf $P_{X|Y=g}$ or

pmf $f_{X|Y=g}$ and then computing variance of it

* This yields a function of $y = g(Y)$ - plug Y in for y that yields

$$\text{Var}(X|Y) = g(Y)$$

How to compute? In general

- Find conditional pdf/pmf $f_{X|Y=g}$ / $P_{X|Y=g}$

- Mean of that is $E(X|Y=g)$

- Variance of that is

$$\text{Var}(X - E(X|Y)) | Y=y = \begin{cases} \int_x (x - E(X|Y))^2 f_{X|Y}(x) dx & \text{OR} \\ \sum_x (x - E(X|Y))^2 P_{X|Y}(x) \end{cases}$$

Law of Total Variance: (a sometimes useful identity)

$$\text{Var}(X) = E[\text{Var}(X|Y)] + \text{Var}[E(X|Y)]$$

Total probability rule: If A_1, \dots, A_n partition Ω , then

$$P(X) = \sum_{i=1}^n P(A_i) P_X(x|A_i)$$

Given two rvs on same Ω, P - say X and Y - define

$$P_{XY}(x,y) = \frac{P\left(\begin{array}{c} X=x \\ Y=y \end{array}\right)}{P(X=x)} = \frac{P(X=x \cap Y=y)}{P(X=x)} = \frac{P(X=x)P(Y=y|X=x)}{P(X=x)} = \frac{P(Y=y|X=x)}{P(X=x)}$$

For fixed y , $P_{X|Y=y}$ defines a pmf over x -values - i.e.

$$P_{X|Y=y}(x) > 0 \quad \forall x \quad \text{and} \quad \sum_x P_{X|Y=y}(x) = 1$$

" $P_{X|Y=y}(x)$ = conditional pmf of X given $Y=y$ "

Like the conditional "event-centered story", have a product rule of sorts

$$P_{XY}(x,y) = P_X(x)P_{Y|X=x}(y) \quad \forall x, y$$

OR

$$P_{XY}(x,y) = P_Y(y)P_{X|Y=y}(x) \quad \forall x, y$$

This expresses joint in terms of marginals + conditional(s).

Also have a total-probability rule of sorts:

$$P_X(x) = \sum_y P_{XY}(x,y) P_{Y|X=x}(y)$$

$$P_Y(y) = \sum_x P_{XY}(x,y) P_{X|Y=y}(x)$$

Moment Generating Function (MGF)

Given X , continuous or discrete, define MGF of X as

$$M_X(s) = E(e^{sX}) = s \text{ is a variable!}$$

$M_X(s)$ is a function of s - need to be careful about its domain of definition

$$IE(X^k) = \frac{d^k}{ds^k} M_X(s) \Big|_{s=0}$$

Let

$$M_n = \frac{X_1 + \dots + X_n}{n}, \quad X_k \text{ iid}$$

Since

$$IE(M_n) = \mu \neq n; \quad \text{Var}(M_n) = \frac{\sigma^2}{n} \neq \sigma^2$$

Chebyshev Inequality:

$$P(|X - \mu| \geq c) \leq \frac{\text{Var}(X)}{c^2} \neq c$$

From this, it follows that

$$P(|M_n - \mu| \geq \varepsilon) \leq \frac{\sigma^2}{n \varepsilon^2} \neq \varepsilon^2$$

Consequence

$$P(|M_n - \mu| \geq \varepsilon) \rightarrow 0 \text{ as } n \rightarrow \infty \neq \varepsilon^2 \rightarrow 0$$

WLLN

Markov Inequality: If X is a nonnegative-valued rv, then for every $c > 0$,

$$P\{|X| \geq c\} \leq \frac{E(X)}{c}$$

Think of these both as quantitative bounds on tail probabilities

Chubyshev

Markov



Another consequence of Chubyshev:

$$P\{|X - \mu| \geq k\sigma\} \leq \frac{1}{k^2} \quad \text{"probability of big std deviation away from mean is 1/k^2"}$$

The kind of convergence taking place in WLLN converges in probability

Definition: Given a sequence of rvs Y_n , $n \in \mathbb{N}$, and a number a , say Y_n converges in probability to a as $n \rightarrow \infty$,

$$\lim P(|Y_n - a| > \varepsilon) = 0$$

WLLN just says " $Y_n \rightarrow \mu$ in probability as $n \rightarrow \infty$ "

Last thing: convergence w/ probability 2 of a sequence Y_1, Y_2, \dots of random variables.

Consider the sequence $\{Y_n : n > 0\}$.

Given some random variable Y ,

$$\lim_{n \rightarrow \infty} Y_n = Y$$

is an event need to refer back to Ω, P , etc.

Say $Y_n \rightarrow Y$ with probability 1 (w.p. 1)

$$Y_n \xrightarrow{w.p. 1} Y$$

$Y_n \xrightarrow{w.p. 1} Y$ as almost surely

when this event has probability 1.

$$X, Y \text{ independent} \Leftrightarrow F_{X,Y}(x,y) = F_X(x)F_Y(y) \quad \forall x, y$$

\uparrow implies (take $\partial/\partial x, \partial/\partial y$)

$$X, Y \text{ independent} \Leftrightarrow f_{X,Y}(x,y) = f_X(x)f_Y(y) \quad \forall x, y$$

Comment: When X, Y independent, we have

$$- E(XY) = E(X)E(Y)$$

$$- E(g(X)h(Y)) = E(g(X))E(h(Y))$$

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

\uparrow implies (take $\partial/\partial x, \partial/\partial y$)

Given Ω , P , and $X: \Omega \rightarrow \mathbb{R}$ a rv. Say X is a continuous rv wrt there exists a function $f_X(x)$ - called the probability density function (pdf) of X - such that "any" $\forall V \in \mathbb{R}$,

$$P(\{X \in V\}) = \int_V f_X(x) dx \quad \text{if } f_X(x) \text{ has to be reasonable enough for integrals to make sense}$$

$\Rightarrow f_X(x) \geq 0 \nabla x$ (need this to ensure $P(\{x \in V\}) \geq 0 \nabla V \in \mathbb{R}$)

$$\lim_{n \rightarrow \infty} \int_{-\infty}^{\infty} f_X(x) dx = 1 \rightarrow \int_{-\infty}^{\infty} f_X(x) dx = P(X \in (-\infty, \infty)) = 1$$

- Given $x \in \mathbb{R}$, $f_X(x)$ is NOT $P(\text{some event})$ - in particular, $f_X(x) \neq P(\{X=x\})$

Turns out $P(\{X=x\})=0 \nabla x \in \mathbb{R}$ when X is a continuous random variable

Expected Value

The expected value of a continuous rv X w/ pdf $f_X(x)$:

$$E[X] = \int_{-\infty}^{+\infty} x f_X(x) dx \quad \text{Caution: NOT always defined - integral might fail to exist}$$

Expected Value Rule

Given X w/ pdf $f_X(x)$ and $Y = g(X)$, we have

$$E[Y] = \int_{-\infty}^{+\infty} g(x) f_X(x) dx \quad \text{enables } E[Y] \text{ computation w/o finding } f_Y(y) \text{ or } F_Y(y)$$

Variance

Variance of continuous rv:

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

By expected value rule, we have

$$\text{Var}(X) = \int_{-\infty}^{+\infty} (X - E[X])^2 f_X(x) dx$$

At ∞ , as before,

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

Next, define - for ANY rv X (discrete or continuous) the cumulative distribution function (cdf) by

$$F_X(x) = P(\{X \leq x\}) \nabla x \in \mathbb{R}$$

Observation: If X is a continuous rv w/ pdf $f_X(x)$, then since

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

we have

$$F_X(x) = \int_{-\infty}^x f_X(t) dt \quad \text{and} \quad f_X(x) = \frac{d}{dx} F_X(x)$$

Discrete version: If X is a discrete rv w/ pmf $p_X(x)$

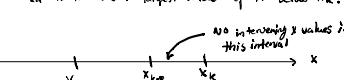
we have

$$F_X(x) = \sum_{x_k: x_k \leq x} p_X(x_k) \quad \text{set of all possible } X\text{-values that don't exceed } x$$

Can invert this formula to get $p_X(x)$ in terms of $F_X(x)$:

$$p_X(x_k) = F_X(x_k) - F_X(x_{k-1})$$

where x_{k+1} is the "NEXT largest value" of X below x_k .



General Properties of CDFs

$$\lim_{x \rightarrow -\infty} F_X(x) = 0 \quad \text{and} \quad \lim_{x \rightarrow \infty} F_X(x) = 1$$

(2) When X is a continuous rv, $F_X(x)$ is continuous in x and differentiable "almost everywhere" (comes in finite) correspond to jumps in $f_X(x)$

(3) X is a discrete iff $f_X(x)$ is a piecewise constant.

(4) $F_X(x)$ is monotonically increasing in x .

$$x_1 < x_2 \Rightarrow F_X(x_1) \leq F_X(x_2)$$

Say X, Y rvs defined in same Ω, P are jointly continuous w/ joint pdf $f_{X,Y}(x,y)$ when

$$P(\{(X,Y) \in V\}) = \iint_V f_{X,Y}(x,y) dx dy \quad \forall V \subset \mathbb{R}^2$$

Special case of a $V: [a_1, b_1] \times [a_2, b_2]$

\Rightarrow

Then

$$P(\{(X,Y) \in V\}) = \int_{a_1}^{b_1} dx \int_{a_2}^{b_2} dy (f_{X,Y}(x,y))$$

Again, have marginals

$$f_X(x) = \int_{-\infty}^{+\infty} f_{X,Y}(x,y) dy \quad ; \quad f_Y(y) = \int_{-\infty}^{+\infty} f_{X,Y}(x,y) dx$$

An official way to get this: Get $F_X(x)$ first then take $\frac{d}{dx} F_X(x)$

$$F_X(x) = P(\{X \leq x\}) = P(\{(X,Y) \in (-\infty, x] \times (-\infty, \infty)\})$$

$$= \int_{-\infty}^x dt \int_{-\infty}^{+\infty} dy (f_{X,Y}(t,y))$$

then

$$\frac{d}{dx} F_X(x) = \int_{-\infty}^{+\infty} dy (f_{X,Y}(x,y))$$

Could also derive marginal formulas as follows:

$$\forall V \subset \mathbb{R}, P(\{X \in V\}) = P(\{(X,Y) \in V \times (-\infty, \infty)\}) = \int_V \int_{-\infty}^{+\infty} f_{X,Y}(x,y) dx dy = \int_V \int_{-\infty}^{+\infty} f_{X,Y}(x,y) dy dx \quad \text{Must be } f_{X,Y}(x,y) \text{ integrate over } V \text{ to get } P(\{X \in V\})$$

Other stuff

$$- \int_{-\infty}^{+\infty} dx \int_{-\infty}^{+\infty} dy (f_{X,Y}(x,y)) = 1$$

$$- \text{Joint CDF: } F_{X,Y}(x,y) = P(\{(X,Y) \in (-\infty, x] \times (-\infty, y]\})$$

$$- f_{X,Y}(x,y) = \frac{\partial^2}{\partial x \partial y} F_{X,Y}(x,y)$$

Conditional Stuff For Continuous Random Variables

Given a continuous rv X on Ω, P and some event $A \subset \Omega$, the conditional pdf of X given A "defined" as follows:

For any $V \subset \mathbb{R}$, we have

$$P(\{X \in V\} | A) = \int_V f_{X|A}(x) dx$$

In general, no decent formula for $f_{X|A}(x)$ in terms of $f_X(x)$.

One way to compute it:

- First get conditional cdf of x given A

$$F_{X|A}(x) = P(\{X \leq x\} | A)$$

- Then take $\frac{d}{dx}$ to get $f_{X|A}(x)$

However, if A is an event of the form $\{X \in V\}$, and $P(A) > 0$, we have

$$f_{X|A}(x) = \begin{cases} \frac{f_X(x)}{P(\{X \in V\})}, & \text{when } x \in V \\ 0, & \text{otherwise} \end{cases}$$

How does this arise?

$$P(\{X \in V\} | A) = \frac{P(\{X \in V \cap A\})}{P(A)} = \frac{\int_V f_X(x) dx}{P(A)}$$

Total Probability Theorem in context of $f_{X|A}(x)$:

If X is a continuous rv and A_1, \dots, A_n are events of positive probability that partition Ω , then

$$f_X(x) = \sum_{k=1}^n f_{X|A_k}(x) P(A_k)$$

To see this: go via cdfs.

$$F_{X|A_k}(x) = \frac{P(\{X \leq x\} \cap A_k)}{P(A_k)}$$

$$\frac{d}{dx} F_{X|A_k}(x) = f_{X|A_k}(x)$$

By Total Probability Theorem,

$$f_X(x) = P(\{X \leq x\}) = \sum_{k=1}^n F_{X|A_k}(x) P(A_k) \xrightarrow{\frac{d}{dx}} \sum_{k=1}^n f_{X|A_k}(x) P(A_k) = f_X(x)$$

Comment: this holds when A_k aren't of the special form $\{X \in V\}$!

Bottom line: conditional pdf of X given $Y=y$ is

$$f_{X|Y}(x|y) = \frac{f_{X,Y}(x,y)}{f_Y(y)} \quad \text{What you integrate over for any } x \in V \text{ to get } P(\{X \in V | Y=y\})$$

e.g.

$$f_{X|Y}(x|y) = f_{X,Y}(x,y) f_Y(y)$$

Integrate over x or y to get

$$f_X(x) = \int_{-\infty}^{+\infty} f_{X|Y}(x|y) dy \quad \text{or} \quad f_Y(y) = \int_{-\infty}^{+\infty} f_{X|Y}(x|y) f_X(x) dx$$

\Rightarrow X is a continuous rv whose density $f_X(x)$ is concentrated on a single interval $a < x < b$. $a = -\infty$ and/or $b = +\infty$ allowed.

$\Rightarrow Y=g(X)$, $\text{g strictly monotonic and differentiable}$, implying that $f_Y(y)$ is concentrated on $(g(a), g(b))$ OR $(g(b), g(a))$

- Let h be the inverse function of g - defined only on $(g(a), g(b))$

or $(g(b), g(a))$ - h is also strictly monotonic and differentiable on its domain of definition

$$\text{Then } f_Y(y) = \begin{cases} \frac{dh(y)}{dy} f_X(h(y)) & , y \in (g(a), g(b)) \text{ OR } y \in (g(b), g(a)) \\ 0 & , \text{else} \end{cases}$$

Covariance

Given any X, Y rvs (discrete, continuous, whatever) defined on same probability space, the covariance of X and Y defined as

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

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

Terminology: When $\text{Cov}(X, Y) = 0$; say X and Y are uncorrelated.

Fact: If X, Y independent, then X, Y uncorrelated

Terminology: Given X, Y , the correlation coefficient of X and Y is

$$\rho = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}$$

Central Limit Theorem Converging in the sense that for every z ,

Recall that if X_k iid w/ common μ, σ^2

$$F_{Z_n}(z) \leftarrow \Phi(z)$$

$$M_n = \frac{X_1 + \dots + X_n}{n} \Rightarrow E(M_n) = \mu \text{ and } \text{Var}(M_n) = \frac{\sigma^2}{n}$$

Form Z_n by renormalizing so $E(Z_n)=0$ and $\text{Var}(Z_n)=1$.

$$Z_n = \frac{X_1 + \dots + X_n - n\mu}{\sqrt{n}\sigma} = \frac{\sqrt{n}(M_n - \mu)}{\sigma} \quad \text{In this context, } Z_n \text{ converges to, as } n \rightarrow \infty, \text{ a Gaussian with mean } \mu=0 \text{ and variance } \sigma^2=1.$$

Conditional Expectation Revisited

Terminology: $E(X|Y)$ = conditional expectation of X given Y

Question: What is $E(E(X|Y))$?

Fact: Law of iterated expectations

$$E(X) = E(E(X|Y))$$

Idea: $E(E(X|Y)) = g(Y)$ for some function g .

Thus

$$E(E(X|Y)) = E(g(Y)) \quad \text{use expected value rule to get this}$$

Law of Iterated Expectations: $E(E(X|Y)) = E(X)$

For any function $h(Y)$,

$$\bullet E(h(Y)|Y) = h(Y)$$

$$\bullet E(h(Y)X|Y) = h(Y)E(X|Y)$$

Can think of $E(X|Y)$ as an estimator of X given Y .

In what sense does it "act like an estimator?"

• $E(Y|Y) = E(Y)$ by law of iterated expectations

• The estimation error $X - E(X|Y)$ is uncorrelated w/ the estimate $E(X|Y)$ - in fact, $X - E(X|Y)$ is uncorrelated with Y - More generally, w/ any function $h(Y)$

Fact (Major): $E(X|Y)$ is the function of Y that minimizes $E((X-h(Y))^2)$

$$f_{X|Y}(x|y) = \frac{f_{X,Y}(y|x)}{f_Y(y)} = \frac{f_{X,Y}(y|x)f_X(x)}{\int_{-\infty}^{+\infty} f_{X,Y}(y|x)f_X(x) dx} \quad \left\{ \begin{array}{l} \text{Continuous Bayes's Rule} \\ \text{Rule} \end{array} \right.$$