

Estimators

Key distribution used:

$$\overline{P_n} = \frac{1}{n} \sum_{i=1}^n P_i \text{, where } P_1 = P_2 \dots = P_n$$

(NOTE: $P_n = \sum_{i=1}^n P_i$)

→ This is the sample mean distribution.
It models the distribution form of sample means from n independent samples drawn from a common distribution P .

Properties of sample mean dist. :

- 1) $\mu(\overline{P_n}) = \mu(P)$
- 2) $\text{Var}(\overline{P_n}) = \frac{\text{Var}(P)}{n}$

The purpose of sample mean is to estimate the actual distribution mean.

i.e. sample mean is an estimator of the distribution's mean.

To generalise...

i.e. a map for each n value

An estimator ~~is~~ is a collection of maps

$$T^n : X^n \rightarrow \mathbb{R}$$

implementing some operation tuples (size n) of samples

The purpose of an estimator is to estimate some real value associated with probability measures on X (ex. mean, variance, etc.)

In statistics, we study the disto. of estimators applied to products the product measure

$$\underbrace{P \times P \times \dots \times P}_{n \text{ times}} \text{ on } X^n$$

(i.e. n independently & identically distributed samples)

2 properties of estimators:

1) Unbiasedness: on average, they estimate correctly, i.e. if an estimator T^n is meant to estimate parameter $\phi(P)$ of a disto. & prob. measure P , then T^n is unbiased if

$$\mu(T^n * \underbrace{(P \otimes \dots \otimes P)}_{n \text{ times}}) = \phi(P)$$

↳ disto. mean

2) Consistency: as no. of samples grows, estimator converges to true value of the parameters it is meant to estimate, i.e. T^n is consistent if

$$\lim_{n \rightarrow \infty} T^n * \underbrace{(P \otimes \dots \otimes P)}_{n \text{ times}} ([\phi(P) - \varepsilon, \phi(P) + \varepsilon]) = 1$$

where the expression within the limit is the prob. mass of the interval $[\phi(P) - \varepsilon, \phi(P) + \varepsilon]$ w.r.t. the estimator's disto. & where ε is an arbitrarily small value.

Some results

1) Sample mean is unbiased

2) Biased sample variance i.e.

$$\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \text{ for sample } (x_1, x_2, \dots, x_n)$$

is biased

Limit Laws

1) Law of Large numbers

1.1) Weak law

If P is a prob. measure with mean $\mu(P)$, then for any $\epsilon > 0$,

$$\lim_{n \rightarrow \infty} \overline{P}_n([\mu(P) - \epsilon, \mu(P) + \epsilon]) = 1$$

i.e. sample mean is a consistent estimator of true (disto.) mean.

NOTE: WLLN presupposes that $\mu(P)$ actually exists, it cannot be applied else.

1.2) Strong law

If P is a disto. with mean $\mu(P)$, then

$$P^\infty(\{(x_1, x_2, \dots) \mid \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n x_i = \mu(P)\}) = 1$$

where

$$P^\infty = P \otimes P \otimes \dots \otimes P$$

(x_1, x_2, \dots, x_n) = samples, where $n \rightarrow \infty$

i.e. the probability mass of an infinite set of samples whose finite means converge to the true mean is 1 (provided the samples are independently & identically disto-ed)

All ~~IID~~ IID ~~infinite~~ infinite samples have this property.

NOTE: SLLN also presuppose that $\mu(P)$ actually exists.

2) Central limit theorem

If P is a prob. measure with mean $\mu(P)$ & variance $\text{Var}(P)$, then

$$\lim_{n \rightarrow \infty} \sqrt{n} (\overline{P}_n - \mu(P))((- \infty, t]) \\ = N(0, \text{Var}(P))((- \infty, t])$$

In other terms

$$\sqrt{n}(\bar{P}_n - \mu(P)) \sim N(0, \text{Var}(P))$$

as $n \rightarrow \infty$

Represents the prob. dist. of the deviation of sample means from true mean, multiplied by \sqrt{n} .

We can rearrange the terms to get:

a) $\frac{\sqrt{n}(\bar{P}_n - \mu(P))}{\sqrt{\text{Var}(P)}} \approx N(0, 1)$ as $n \rightarrow \infty$

b) $\bar{P}_n \approx N(\mu(P), \frac{\text{Var}(P)}{n})$ as $n \rightarrow \infty$

Hypothesis testing

Basics

For given data X , we define an estimator to serve a required purpose (ex. if we want to test for disto. mean, the estimator is sample mean). This estimator or some function of it is used as the test statistic. Then, we either identify or approximate the estimator's probability disto. Using this, we can see how likely or unlikely is the particular value we obtain for the estimator for the given data X .

We formulate such a test by formulating the appropriate hypotheses..

Null hypothesis (H_0): The assumption we want to test. called "null" because it is the assumption of no difference btw. assumed & true value. (of a parameter) \rightarrow (or estimated)

\hookrightarrow We reject H_0 if our observation is deemed sufficiently unlikely.

Defining "sufficiently unlikely"

We define a parameter α , the "confidence interval level". This is a proportion the probability mass of some proportion (indicated by α itself, naturally) of the most likely / high-probability / high-density values of the disto. of the test statistic.

ex. $\alpha = 90\%$ \Rightarrow the prob. mass of 90% of the most likely values of the disto.

$\hookrightarrow \alpha = 90\%$ \Rightarrow if an observed value lies outside this range of values, we will consider it sufficiently unlikely

Usually, α is chosen as 95%, 99%, 99.9% etc.

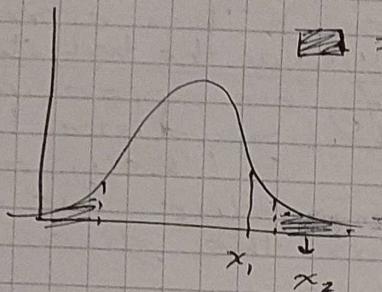
Critical region

The least likely values of the disto. that together have a mass of $1-\alpha$ (called the significance level) \hookleftarrow from the critical region. i.e. they are the $1-\alpha$ least likely proportion

of least likely values in the disto.

If an observation lies in the critical region, it is considered sufficiently unlikely.

Visual example:



■ \Rightarrow critical region

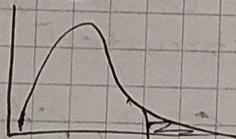
$x_2 \in$ (crit. region)
 \Rightarrow reject H_0

$x_1 \notin$ (crit. region)
 \Rightarrow cannot reject H_0 .

NOTE: Depending on the type of distribution, the $1 - \alpha$ proportion of least likely (i.e. most extreme) values may either:

i) lie only on one side of the disto.

ex.

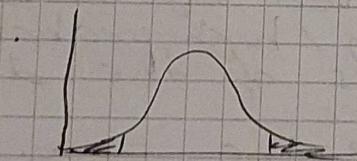


\rightarrow requires one-sided test

- ↳ most extreme values (at the current proportion) lie only to the right
- ↳ (at a higher proportion, we may begin to include from left also).

ii) lie on either side of the disto.

ex.



\rightarrow requires two-sided test

↳ is the case for any symmetric disto.

p-value (to critical region)

+ logically

An alternative[†] (Δ practically equivalent) method to test extremeness of an observation.

mass

p-value = probability[†] (under the assumption that H_0 is true)

↳ (i.e. assuming a disto. for test statistic)

of ~~as~~ all values at least as extreme (i.e. unlikely) as the observed value.

NOTE :

| For ~~one~~-sided test, if p-value < $\frac{1-\alpha}{2}$, reject H_0 .

| For ~~two~~-sided test, if p-value < $1-\alpha$, reject H_0 .

) This is if we calculate p-value as $\min(CDF(t), 1 - CDF(t))$ (where t is observed value)

) If we instead calculate p-value as $2 \min(CDF(t), 1 - CDF(t))$ for two-sided test only, then if p-value < $1-\alpha$, reject H_0 .

WHY? Consider: what does p-value mean? What set of values would be considered "at least as extreme" as the observed value in two-sided test?
Think about it...

Testing errors

Type 1: Rejecting H_0 even if it holds.

By definition, the probability of Type 1 error is $1-\alpha$ (i.e. the significance level).

Type 2: Failing to reject H_0 even if it does not hold

Prob. of committing type 2 is called β . & is harder to derive.

NOTE: We can never validly say we "accept" H_0 , only that we "fail to reject" it.

WHY? Consider: what kind of evidence does statistical testing give? Can it give positive evidence for something, or only a lack of negative evidence against something?

NOTE: H_0 either explicitly or implicitly hypothesis that a test statistic follows a certain distn. If.

Exact test: if exact dist. of test statistic is known

Approximate test: if we can only approx. the dist. of the test statistic.

Approximate test

Preliminary concepts: useful distributions
(helpful in approximating test dists.)

1) $N(0, 1)$ → standard normal dist.

2) $\chi^2(k)$ → Chi squared dist. with k degrees of freedom (df)

NOTE: (extra info)

$$\sum_{i=1}^n (N(0, 1) - \bar{N}(0, 1)_n)^2 = \chi^2(n-1)$$

3) $t(k)$ → Student's t-distr. with k df

NOTE: (useful later)

$$t(k) = \frac{N(0, 1)}{\sqrt{\frac{\chi^2(k)}{k}}}$$

Preliminary concepts: key facts

1) By the CLT (central limit theorem):

$$\frac{\sqrt{n}(\bar{P}_n - \mu(P))}{\sqrt{\text{Var}(P)}} \approx N(0, 1) \quad [\text{as } n \rightarrow \infty]$$

$$2) \frac{\text{SVar}_n(P)}{\text{Var}(P)} \approx \frac{\chi^2(k)}{k} \quad \left| \begin{array}{l} \text{where } k = \frac{2n}{c - \frac{n-3}{n-1}} \\ (c \text{ is some constant}) \end{array} \right.$$

$\text{SVar}_n(P) = \frac{\text{dist. of sample variance of a sample of size } n \text{ drawn from } P}{\text{sample variance of a sample of size } n \text{ drawn from } P}$ where $c = \text{Kurtosis}(P)$

NOTE: $\text{SVar}_n(P)$ is a dist., but $\text{Var}(P)$ is a single value!

To present more clearly: Here,
 $k = \frac{2n}{c - \frac{n-3}{n-1}}$, where $c = \text{Kurtosis}(P)$

3) (1) & (2) \Rightarrow

$$\frac{\sqrt{n}(\bar{P}_n - \mu(P))}{\sqrt{\text{Var}_n(\bar{P})}} \approx \frac{N(0, 1)}{\sqrt{\chi^2(k)}} \quad \text{IMPORTANT here}$$

$$\frac{\sqrt{n}(\bar{P}_n - \mu(P))}{\sqrt{\text{Var}_n(\bar{P})}} \approx t(k)$$

$$\Rightarrow \frac{\sqrt{n}(\bar{P}_n - \mu(P))}{\sqrt{\text{Var}_n(\bar{P})}} \approx t(k)$$

Notice that $\text{Var}(\bar{P})$ is cancelled out; this is useful if $\text{Var}(\bar{P})$ is unknown but we want to test for mean.

NOTE: Estimating df i.e. k .

In general, we may estimate

$df = n$ (sample size)
with reasonable accuracy, if
Kurtosis (P) is unknown.

Uses of the above:

- Z-test: using test statistic in (1)
→ testing mean with known variance
- χ^2 -test: using test statistic in (2)
→ testing sample variance
- t-test: using test statistic in (3)
→ testing mean with unknown variance.

Goodness of fit

(comes under hypothesis testing)

Preliminary concept: Distance b/w. distrs

Formally, a distance function (defined on $X \times Y$ where X & Y are 2 sets) should fulfill the following:

- $d(x, x) = \text{dist}(x, x) = 0$
- $\text{dist}(x, y) = \text{dist}(y, x)$
- $\text{dist}(x, z) \leq \text{dist}(x, y) + \text{dist}(y, z)$

Distance (as general concept) serves as a measure of similarity or dissimilarity: higher dist. \Rightarrow lower similarity
lower " " \Rightarrow higher "

OUR AIM: Measure i.e. quantify the similarity (or dissimilarity) b/w. 2 distrs.

GENERAL IDEA TO MEASURE DISTANCE BTW DISTR.S:

Two distributions are close if they compute similar integrals.

(NOTE: One of the main things we do with prob. distrs is integrate them i.e. obtain expectations.)

- We can compute the integrals along with different functions that would focus on different aspects of the dists. \rightarrow (unclear, need to clarify!)
- There are "test functions" to clarify!

Suppose distrs P_1 & P_2 have densities f_1 & f_2 , then P_1 & P_2 are close if $\int g(x) f_1(x) dx$ is close to $\int g(x) f_2(x) dx$ for all test functions g .

Generalising with "test functions"

Measuring distances b/w. dists. using integrals

Suppose dists. P_1 & P_2 have densities f_1 & f_2 . Then, distance may be calculated in one of the following ways:

- Max. value of $\left| \int g(x) f_1(x) dx - \int g(x) f_2(x) dx \right|$ over all test functions, i.e.

$$\text{dist}(P_1, P_2) = \sup_g \left| \int g(x) f_1(x) dx - \int g(x) f_2(x) dx \right|$$

- Weighted avg. of

$$\left[\int g(x) f_1(x) dx - \int g(x) f_2(x) dx \right]^2$$

over all test functions.

Preliminary concept: Empirical PDF, CDF & probability measure

Empirical PDF:

Suppose we have

- n samples x_1, x_2, \dots, x_n from dist. P

Then \rightarrow Any sort density-histogram we build from these is the empirical density dist.

For this, choose a sequence of k bins $y_1 < y_2 < \dots < y_k$

(a good choice could be $k = \sqrt{n}$)

Then we define the empirical probability density function as

$$f_{\text{emp}}(x) = \sum_{i=1}^{k-1} \mathbb{1}_{[y_i, y_{i+1}]}(x) \cdot \left(\frac{\sum_{j=1}^n \mathbb{1}_{[y_i, y_{i+1}]}(x_j)}{n} \right)$$

NOTE: $\mathbb{1}_{[y_i, y_{i+1}]}(x) = \begin{cases} 1 & \text{if } x \in [y_i, y_{i+1}], \\ 0 & \text{if } x \notin [y_i, y_{i+1}] \end{cases}$

Also, $\sum_{j=1}^n \mathbb{1}_{[y_i, y_{i+1}]}(x_j) = \text{# of all samples belonging to bin } [y_i, y_{i+1}]$

(Proposition)

+ NOTE: As $n \rightarrow \infty$, we expect empirical PDF to converge to actual PDF of P .

CONSIDER: Empirical PDF's define prob. measures. Why?

Empirical CDF

As before, suppose we have

- n samples x_1, x_2, \dots, x_n from dist σ . P
- The unique empirical CDF for w.r.t. these samples is

$$F_n(x) = \frac{1}{n} \sum_{i=1}^{*n} 1_{[x_i, \infty)}(x)$$

CONSIDER: what does empirical probability cumulative prob. mean?
It means you add up the empirical probabilities proportional to frequencies
i.e. the proportion of the

Essentially equals the proportion of times value the no. of times value samples at least less than or equal to x appeared in the sample data set.

+ NOTE: As $n \rightarrow \infty$, we expect

F_n to converge to the actual CDF of P .

Empirical measure

Again, suppose we have

- n samples x_1, x_2, \dots, x_n from dist σ . P .

Also, define δ_x $\delta_x(A) = \begin{cases} 1 & \text{if } x \in A \\ 0 & \text{else} \end{cases}$

For a single value, $\delta_x(a) = \begin{cases} 1 & \text{if } x = a \\ 0 & \text{else} \end{cases}$

- Then, the unique empirical prob. measure corresponding to the samples is

$$P_n = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$$

$$\text{Hence, } P_n(A) = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}(A)$$

Clearly, we see that

$$P_n(A) = \frac{\text{No. of samples that are present in } A}{n}$$

from x_1, x_2, \dots, x_n

+ NOTE: As $n \rightarrow \infty$, $P_n \rightarrow P$

+ NOTE: It is easy to show that empirical \Leftrightarrow CDF can be derived from P_n .

Main topic: Goodness of fit tests

Discrete case

→ Test To evaluate distance between theoretical & empirical dists. :

- Test function used: $1_t(x) = \begin{cases} 1 & \text{if } x=t \\ 0 & \text{else} \end{cases}$

- Distance method used:
Weighted avg. of

$$\left(\int g(x) f_1(x) dx - \int g(x) f_2(x) dx \right)^2$$

NOTE: In discrete case, the "integral" is essentially summation.

- Consider $\{z_1, z_2, \dots, z_k\}$ as the exhaustive support of the theoretical dists., P .

- Consider a sample x_1, x_2, \dots, x_n drawn from P

- Let f be the PMF of P . Then, we can obtain the probabilities of each point of the support as

$$\sum_{i=1}^k 1_{z_k}(z_i) f(z_i) = P(z_k) \dots \quad (1)$$

- From x_1, x_2, \dots, x_n , we build the empirical measure
(turn to next page)

$$P_n = \frac{1}{n} \sum_{i=1}^n \delta_{x_i} \rightarrow \text{also the PMF empirical}$$

Hence, we have defined

- The theoretical disto.
- The empirical disto.
- The test function

"Integral" of test function w.r.t. P is in fact $\sum_{i=1}^k 1_{z_k}(z_i) f(z_i) = P(z_k)$

Further, "integral" of test function w.r.t. P_n is in fact

$$\sum_{i=1}^k 1_{z_k}(z_i) P_n(z_i) = P_n(z_k)$$

Hence we have the "integrals" needed to calculate distance btw. P_n & P .

But weighting scheme?

We choose a scheme s.t. the disto. of the distance can be approximated using central limit theorem (CLT).

Hence, we define

$$\text{dist}(P, P_n) = \sum_{i=1}^k \frac{1}{P(z_i)} (P(z_i) - P_n(z_i))^2$$

NOTE: This is for discrete case

or difference of "integrals" weighting scheme

Results on the above computed distance measure for discrete case:

1) As $n \rightarrow \infty$, $\text{dist}(P, P_n) \rightarrow 0$

(is derived from the next theorem)

2) Theorem (Pearson):

$$n \cdot \text{dist}(P, P_n) \rightarrow \chi^2(k-1) \text{ as}$$

$$n \rightarrow \infty$$

NOTE: k = size of support of P

Hence, we can approximate this sample distances' distribution!

The goodness of fit test tests how well the dist^r. P agrees with the empirical dist^rs. of the drawn samples i.e. P_n .

Thus Thanks to the previous results, we can approximate how the distance b/w. theoretical & empirical dist^rs are distributed in general & thus, we can evaluate the plausibility of observed cases (i.e. observed empirical dist^rs).

H_0 : The samples $x_1, x_2 \dots x_n$ are drawn from P

$$\boxed{\text{Under } H_0, n \cdot \text{dist}(P, P_n) \approx \chi^2(k-1)}$$

H_0 is rejected if $n \cdot \text{dist}(P, P_n)$ is too large i.e. too unlikely.

Since this test is based on Pearson's theorem, it is called

Pearson's χ^2 -test

NOTE: χ^2 -test \Rightarrow one-sided test

(right-tailed) $\Rightarrow p\text{-value} = \chi^2([\text{ndist}(P, P_n), \infty))$

Goodness of fit, next topic ...)

Continuous case

To evaluate distance b/w. ~~the~~ theoretical & empirical dist^rs:

- Test function used: $1_{(-\infty, t]}$
(rather, the class of test functions used)

$$g_t(x) = 1_{(-\infty, t]}(x)$$

- Distance method used:

Since we need to deal with a class of test functions, \Rightarrow we use max. value of

$$\left| \int g_t(x) f_1(x) - \int g_t(x) f_2(x) dx \right|$$

(maximized using t) dx

Now, note that

$$\int g_t(x) f(x) dx = \int_{-\infty}^t f(x) dx = F(t)$$

↓
density func.
of P
(theoretical dist.)

↓
theoretical
CDF

Similarly, note that

$$\int g_t(x) f_{emp}(x) dx = \int_{-\infty}^t f_{emp}(x) dx = F_{emp}(t)$$

↓
(empirical PDF)

Hence, we get the distance b/w. theoretical & empirical dists as

$$dist(P, P_n) = \sup_t |F(t) - F_{emp}(t)|$$

[Oops, F_{emp} is notated as F_n !]

$$dist(P, P_n) = \sup_t |F(t) - F_n(t)|$$

Results for distances b/w. dists in continuous case:

1) Gilivenko-Cantelli theorem:

$$dist(P, P_n) \rightarrow 0 \text{ as } n \rightarrow \infty$$

2) Kolmogorov theorem:

$$\sqrt{n} dist(P, P_n) \rightarrow K \text{ as } n \rightarrow \infty$$

(K here is the Kolmogorov dist.)

Just like for discrete case, now we can do hypothesis testing (goodness of fit testing) for continuous dists. too.

This test is the Kolmogorov-Smirnov (KS) test, wherein we

$$KS\text{-statistic} = \sqrt{n} dist(P, P_n)$$
$$p\text{-value} = K([\sqrt{n} dist(P, P_n), \infty))$$

& also a one-sided right tailed test.

NOTE: H_0 is the same as before (except for distribution used).

Independence, dependence & linear dependence

Preliminary concept: Joint distributions
(or distrs over product spaces)

A product space is the cross product of 2 or more sets, i.e. $X_1 \times X_2 \times \dots \times X_n$. \downarrow (Cartesian)

A distribution on a product space is often called a joint distribution.

NOTE: We have seen product measures,
ex. $P_1 \otimes P_2 (X_1 \times X_2)$
which are a kind of joint dists.

Partial Independence

Product measure captures the idea of independence. For example

$$P_1 \otimes P_2 (A \times B) = P_1(A) P_2(B)$$

= Prob. that the random experiment modelled by P_1 returns outcome A

& Prob. that the random experiment modelled by P_2 returns outcome B

One outcome does not affect the other. A order of outcomes / of observing the outcomes does not change the joint probability above.

i.e. no information flow btw. the experiments.

KEY POINT TO REITERATE: All product measures are joint dists., but not all joint dists. are product measures.

Hence, $P(A \times B) \neq P_1(A) P_2(B)$
for any P_1 & P_2 .

In general, we cannot assume independence in a joint dist.

Perfect dependence

Consider the "copy" function

$$\text{copy}(x) = (x, x), \forall x \in X$$

$$\text{i.e. } \text{copy} : X \rightarrow X \times X, x \mapsto (x, x)$$

We can use this to create a perfectly dependent joint dist.

$$P(\{(x, x) \mid x \in X\})$$

SIDE NOTE: Alternatively, the above joint dist. is representable as

$$\text{copy}_* P(A) = P(\{x \mid (x, x) \in A\})$$

Marginals (key concept)

Any product space $X_1 \times X_2 \dots X_n$

comes with n special maps called projections:

$$\pi_1 : X_1 \times X_2 \dots X_n \rightarrow X_1, (x_1, x_2, \dots, x_n) \mapsto x_1$$

$$\pi_2 : X_1 \times X_2 \dots X_n \rightarrow X_2, (x_1, x_2, \dots, x_n) \mapsto x_2$$

$$\pi_1 : X_1 \times X_2 \dots X_n \rightarrow X_1, (x_1, x_2, \dots, x_n) \mapsto x_1$$

$$\pi_2 : X_1 \times X_2 \dots X_n \rightarrow X_2, (x_1, x_2, \dots, x_n) \mapsto x_2$$

etc. upto π_n

Basically - projection i.e. π_i maps any given tuple from the product space to the i th coordinate of the tuple i.e. $\pi_i((x_1, x_2, \dots, x_n)) = x_i$

(It's basically a way to pick a particular coordinate from a tuple of the product space.)

The marginals of a joint dist. P on $X_1 \times X_2 \dots \times X_n$ are the pushforward of the projections i.e. (turn to next page)

i.e. marginals are measures on particular axes of the product space, obtained via the ~~joint dist.~~'s measure defined for the product space as a whole.

Hence, if P is a measure on the whole product space $X_1 \times X_2 \dots X_n$, then,

$\pi_1 * P$ is a measure on X_1

$\pi_2 * P$ is a measure on X_2

etc. upto $\pi_n * P$

Hence, if $A \subseteq X_i$, then

$$\pi_i * P(A) = P(X_1 \times \dots \times$$

$$= P(X_1 \times X_2 \dots X_{i-1} \times \underline{A} \times X_{i+1} \times \dots \times X_n)$$

Basically, it's the joint dist. where you keep ~~all but~~ one coordinate as constant. (all others)

Illustrative example:

Consider product space ~~{0, 1}~~ $\times \{0, 1\}$,

~~A $\times A$, where $A = X \times X$, where~~

$X = \{0, 1\}$. Consider the following joint dist.:

$$P(0,0) = \frac{1}{4}, P(0,1) = \frac{1}{4}, P(1,0) = \frac{1}{2}$$

$$\text{and } P(1,1) = 0.$$

Then, $\pi_1 * P$ is a measure defined on $\{0, 1\}$, s.t.

$$\pi_1 * P(0) = P(0 \times X)$$

$$= P(0 \times \{0, 1\})$$

$$= P(\{(0,0), (0,1)\})$$

$$= P(0,0) + P(0,1)$$

$$= \frac{1}{4} + \frac{1}{4} = \frac{1}{2}$$

(Essentially, we have defined the prob. of observing 0 in the 1st coordinate)

$$\text{Similarly, } \pi_1 * P(1) = P(1 \times X) \\ = P(1 \times \{0, 1\}) \\ = \frac{1}{2} + 0 = \frac{1}{2}$$

— Similarly we can define $\pi_2 * P$ —

Continuous example...

Let P be a ~~product~~ joint measure defined on $\mathbb{R} \times \mathbb{R}$, with its density being $f: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$

Now, consider $\pi_1 * P$ defined for a specific interval $[a, b] \subseteq \mathbb{R}$

$$\pi_1 * P \pi_1 * P([a, b])$$

$$= \int_{-\infty}^{\infty} \int_a^b f(x, y) dx dy$$

i.e., x is integrated over range $[a, b]$
(1st coordinate)
while y is integrated over range
 $(-\infty, \infty)$ (no bounds).

Some results on marginals:

1) Theorem: Let P on $\mathbb{R} \times \mathbb{R}$ have density $f(x, y)$. Then,

$$\mu_x = \mu(\pi_1 * P) = \iint x f(x, y) dx dy$$

mean of 1st marginal's dist.

2) Similarly, we have

$$\mu_y = \mu(\pi_2 * P) = \iint y f(x, y) dx dy$$

Results: Marginals for independent & perfectly dependent joint distns

For independent measure $P_1 \otimes P_2$ on $X \times Y$,

$$\pi_1 * (P_1 \otimes P_2)(A) = P_1 \otimes P_2(A \times Y)$$

$$= P_1(A) P_2(Y)$$

$\therefore = P_1(A)$
i.e. total prob. of A .

Similarly,

$$\pi_2 * (\text{IP}_1 \otimes \text{IP}_2)(B) = \text{IP}_2(B)$$

Hence, marginals of a product measure are just its components.

For perfectly dependent disto-s, copy* IP defined on $X \times X$,

$$\pi_{L_1} * (\text{copy}_* \text{IP})(A) = \text{IP}(A)$$

$$\pi_{L_2} * (\text{copy}_* \text{IP})(A) = \text{IP}(A)$$

(the 1st & 2nd coordinates are perfectly dependent, so the total prob. for each is equal)

Linear dependence

Preliminary: covariance, correlation

Theoretical covariance:

Discrete case:

$$\text{Cov}(\text{IP}) = \sum_x \sum_y (x - \mu_x)(y - \mu_y) f(x, y)$$

NOTE: IP is the joint disto., μ_x & μ_y are the means of $\pi_{L_1} * \text{IP}$ & $\pi_{L_2} * \text{IP}$ respectively (i.e. 1st & 2nd marginals) & f is the density of IP.

Continuous case

$$\text{Cov}(\text{IP}) = \iint (x - \mu_x)(y - \mu_y) f(x, y) dx dy$$

+ special cases:

Cov for independent joint disto.:

$$\text{Cov}(\text{IP}) =$$

Key result:

$$\text{Cov}(\text{IP}) = \iint xy f(x, y) dx dy - \mu_x \mu_y$$

+ Corollary: $\text{Cov}(\text{IP}_1 \otimes \text{IP}_2) = 0$

independent joint disto.

+ Special Corollary: $\text{Cov}(\text{copy}_* \text{IP}) = \text{Var}(\text{IP})$
 perfectly dependent joint dists.

Sample covariance

$$\text{Scov} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

where we have n samples
 $(x_1, y_1), (x_2, y_2) \dots (x_n, y_n)$

& Scov is the sample cov. for them.

NOTE: For $(x_1, x_1), (x_2, x_2) \dots (x_n, x_n)$,
 clearly, $\text{Scov} = \text{Unbiased sample variance}$.

+ POINTS:

$\text{Scov}_*(\underbrace{\text{IP}(x) \dots \text{IP}}_{n \text{ times}})$ (where IP is a joint dist., & $\text{IP}(x) \dots \text{IP}$ is the dist. of n -sized tuple of samples each drawn from IP)
 gives the dist. of sample covs when samples are drawn from IP .

+ NOTE: Sample cov. is unbiased estimator of cov.

Correlation

It is essentially "standardised" covariance...

$$\text{Corr}(\text{P}) = \frac{\text{Cov}(\text{P})}{\sigma_x \sigma_y}, \text{ where}$$

$$\sigma_x^2 = \text{Var}(\pi_1 * \text{IP}) \quad (\text{variance of 1st marginal})$$

$$\sigma_y^2 = \text{Var}(\pi_2 * \text{IP}) \quad (" \quad " \text{ 2nd } ")$$

NOTE:

$$\sigma_x^2 = \sum \sum (x - \mu_x)^2 f(x, y) . \quad \left. \begin{array}{l} \text{discrete} \\ \text{case} \end{array} \right\}$$

$$\sigma_y^2 = \iint (y - \mu_y)^2 f(x, y) dx dy \quad \left. \begin{array}{l} \text{continuous} \\ \text{case} \end{array} \right\}$$

Results on correlation:

- 1) $-1 \leq \text{Corr}(P) \leq 1$
- 2) $\text{Corr}(P_1 \otimes P_2) = 0$
- 3) $\text{Corr}(\text{copy}_* P) = 1$

NOTE: No dependence \Rightarrow no correlation,
but no correlation does not imply
no dependence (ex. dependence may be
non-linear)

Sample correlation

$$\text{SCorr} = \frac{\text{SCov}}{\sqrt{\text{SVar}_x \cdot \text{SVar}_y}}$$

sample variance of 1st coordinate (x 's)
 sample variance of 2nd coordinate (y 's)

$$= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}}$$

+ POINTS: $\underbrace{n}_{n \text{ times}}$

$\text{SCorr} * (\overline{IP} \otimes \dots \otimes \overline{IP})$ (where IP is joint dist.)
 is dist. of sample correlation for
 n -sized samples drawn from IP .

Linear dependence

Consider the function

$$\text{lin}_{a,b}(x) = (x, ax+b)$$

(it is a generalisation of the previous copy function; copy is basically $\text{lin}_{0,1}$)

Now, consider a distribution P with support X .

Then, $\text{lin}_{a,b} * P$ defines a joint disto. with support

$$\{(x, ax+b) \mid x \in X\}$$

~~$\text{lin}_{a,b} * P$~~ $\text{lin}_{a,b} * P$ is a joint disto. displaying linear dependence, wherein one coordinate in a pair is a linear function of the other.

Results on linear dependence:

1) Theorem: $\text{Cov}(\text{lin}_{a,b} * P) = a \text{Var}(P)$

2) Theorem: $\text{Corr}(\text{lin}_{a,b} * P)$

$$= \begin{cases} 1 & \text{if } a > 0 \\ -1 & \text{if } a < 0 \end{cases}$$