

# Notes on probabilities

## Math

$$0 + 1 + 2 + \dots + n = \frac{n(n+1)}{2}$$

$$0^2 + 1^2 + 2^2 + \dots + n^2 = \frac{1}{6}n(n+1)(2n+1)$$

Sequence: a function such that sequence(S::Set{T},i::N) -> s\_i::T (define an order)

Series: the sum of the elements of a sequence (the cumsum over a given order)

$$\text{Given } |x| < 1 \rightarrow \sum_{i=0}^{\infty} x^i = \frac{1}{1-x};$$

$$\int e^{ax} dx = \frac{1}{a} e^{ax};$$

$$\int u(x)v'(x)dx = u(x)v(x) - \int u'(x)v(x)dx$$

$$\int_{-\infty}^{+\infty} e^{-x^2/2} dx = \sqrt{2\pi}$$

$$\int 1/x dx = \ln(x) + c$$

$$\int \ln(ax) dx = x \ln(ax) - x;$$

$$\ln(x) + \ln(y) = \ln(xy)$$

$$\log_b(a^c) = c \log_b(a) \quad \log_{b^2} x = \frac{\log_{b^1} x}{\log_{b^1} b^2}$$

$$\text{Circumference: } (x-a)^2 + (y-b)^2 = r^2; \sin(0) = \sin(\pi) = 0$$

$$\binom{a+b+c}{a,b,c} = \frac{(a+b+c)!}{a!b!c!}$$

$$x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

$$\lim_{n \rightarrow \infty} \left(1 + \frac{\lambda}{n}\right)^n = e^\lambda$$

- Hessian: second derivatives matrix; Gradient: vector of first derivatives; Jacobian:  $I \times J$  matrix of first derivative of the  $i$  equation for the  $j$  variable
- $x^T H x \leq 0 \quad \forall x \in R^d$  or  $H$  is diagonal and all elements negative  $\leftrightarrow H$  negative semidefinite, all eigenvalues non-positive
- Positive definite:  $|D_i| > 0 \quad \forall i \leq n$  with  $D_i$  the  $i$ -th leading principal minor of the Hessian
- Negative definite:  $(-1)^i |D_i| > 0 \quad \forall i \leq n$
- Vector products:
  - Inner ("dot") product:  $X \in R^d \cdot Y \in R^d \rightarrow OUT = ||X|| * ||Y|| * \cos(\theta) \in R^1$
  - Hadarnard ("elementwise") product:  $X \in R^d \odot Y \in R^d \rightarrow OUT \in R^d$
  - Outer product:  $X \in R^d \otimes Y \in R^d \rightarrow OUT \in R^{(d^2)}$
  - Cross product:  $X \in R^3 \times Y \in R^3 \rightarrow OUT = ||X|| * ||Y|| * \sin(\theta) * n \in R^3$  (where  $n$  is the unit vector perpendicular to the plane containing  $X$  and  $Y$ )
- $\frac{\partial}{\partial x} \int_a^x f(t) dt = f(x)$ ;  $\frac{\partial}{\partial x} \int_x^a f(t) dt = -f(x)$
- "positive semidefinite matrix" := square matrix such that  $x^T A x \geq 0 \quad \forall x \in R^d$ 
  - In particular are spd matrices all diagonal matrices with all non-negative entries and those that can be decomposed as  $A = P^T D P$  with  $D$  a diagonal matrix with only non-negative entries and  $P$  invertible
- "positive semi-def square root": a matrix  $A^{\frac{1}{2}}$  such that  $A^{\frac{1}{2}} A^{\frac{1}{2}} = A$  with  $A$  being spd
  - the roots themselves are positive semi-def
  - for any positive (semi-)definite matrix, the positive (semi-)definite square root is unique.
- "ortogonal" matrix:  $M^T = M^{-1}$
- $(AB)^T = B^T A^T$
- $(A^{-1})^T = (A^T)^{-1}$
- $AIB = AB$  with  $I$  the identity matrix
- $ABsC = sABC = ABCs = \dots$  with  $s$  a scalar value
- $trace(x^T x) = trace(xx^T)$
- $E[trace(\cdot)] = trace(E[\cdot])S$

## Norm

- $||v||^2 = v^T v = trace(vv^T) = \sum_i v_i^2$
- $||v||_t = (\sum_i v_i^t)^{1/t}$

## Vector space

- Projection of vector  $a$  on vector  $b$ :  $c = \frac{a \cdot b}{||b||} * \frac{b}{||b||}$ .
- Distance of a point  $x$  from a plane identified by  $\theta$  and its offset  $\theta_0$ :  $||\vec{d}|| = \frac{\vec{x} \cdot \vec{\theta} + \theta_0}{||\vec{\theta}||}$
- Orthogonal projection of a point  $x$  on plane identified by  $\theta$  and its offset  $\theta_0$ :  $\vec{x}_p = \vec{x} - \frac{\vec{x} \cdot \vec{\theta} + \theta_0}{||\vec{\theta}||} * \frac{\vec{\theta}}{||\vec{\theta}||}$

Vector independence:

A set of  $J$  vectors  $v_j$  are linear dependent i.f.f. there exist a vector  $c$  not all zeros such that  $\sum_{j=0}^J c_j v_j = 0$  (note that each individual  $c_j$  is a scalar while  $v_j$  is a vector).

If a partition of the set of vectors is linearly dependent the whole set is said to be linear dependent.

Any set of  $J$  vectors of  $D$  elements with  $J > D$  is linearly dependent.

$$\text{rank}(AB) \leq \min(\text{rank}(A), \text{rank}(B))$$

## Models and axioms

$$(\cup_n A_n)^c = \cap_n A_n^c \quad (\cap_n A_n)^c = \cup_n A_n^c$$

$$P(A_1 \cup A_2 \cup A_3 \cup \dots) = P(A_1) + P(A_2 \cap A_1^c) + P(A_3 \cap A_1^c \cap A_2^c) + \dots$$

$$P((A \cap B^c) \cup (A^c \cap B)) = P(A) + P(B) - 2 \cdot P(A \cap B)$$

$$P(A_1 \cap A_2 \cap \dots \cap A_n) \geq P(A_1) + P(A_2) + \dots + P(A_n) - (n-1)$$

## Conditioning and independence

**Partition** of a space: array of mutually exclusive ("disjoint") sets whose members are exhaustive ("complementary") of the space.

**Joint:**  $P(A \text{ and } B) = P(A \cap B) = P(A, B)$

→ note that joint PMF/PDF are multidimensional aka multivariate ( $x$  is a vector)

**Marginal** (unconditional):  $P(A)$

→ for PMF (PDF): we sum (integrate) over all or some dimensions to "remove" them and move from the joint toward the marginal

**Conditional:**  $P(A|B) := P(A, B)/P(B)$

→ Valid also for PMF and PDF with respect to an event

Union:  $P(A \text{ or } B) = P(A \cup B) = P(A) + P(B) - P(A \cap B)$

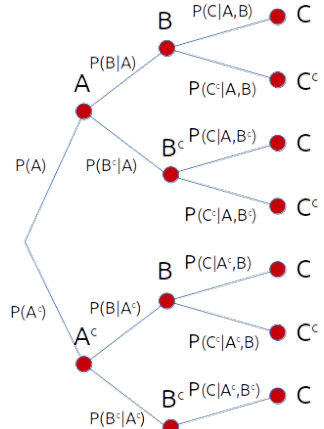
Note the *memoryless* of geometric/exponential:  $Pr(X > t + s | X > t) = Pr(X > s)$ . This is the *remaining* time, not the total time, so it is NOT the independence concept.

## Multiplication rule:

- $P(A_1 \text{ and } A_2) = P(A_1 \cap A_2) = P(A_1) * P(A_2|A_1) = P(A_2) * P(A_1|A_2)$
- $P(A_1 \cap A_2 \cap \dots A_n) = P(A_1) * \prod_{i=2}^n P(A_i | A_1 \cap \dots A_{i-1})$  → also for PMF, PDF

## Total probability/expectation theorem:

- given  $A$  being a partition:  $P(B) = \sum_i P(A) * P(B|A_i)$  → also for PMF, PDF, CDF and expectations





**Bayes' rule:** given  $A$  a partition  $P(A_i|B) = \frac{P(A_i, B)}{P(B)} = \frac{P(A_i)P(B|A_i)}{\sum_j P(A_j)P(B|A_j)}$  where the first relation is by definition and the second one is for the Multiplication rule on the nominator and the total prob. theorem on the denominator  
→ also for PMF, PDF

**Independence:** A, B indep. iff  $P(A \text{ and } B) \equiv P(A \cap B) = P(A) * P(B)$  eq.  $P(A|B) = P(A)$ , equiv.  $P(B|A) = P(B)$   
(a) Indep is symmetric. (b) A collection of event is indep if *every* collection of distinct indices of such collection is indep. (oth. could be pairwise indep.)  
→ also for PMF, PDF, CDF and expectations (but for all  $x, y!$ )

Union rule (De Morgan's law again):

$$P(A_1 \text{ or } A_2 \text{ or } A_3 \text{ or } \dots) = P(A_1 \cup A_2 \cup A_3 \cup \dots) = P(A_1) + P(A_2 \cap A_1^C) + P(A_3 \cap A_1^C \cap A_2^C) + \dots$$

$$P(A_1 \cup A_2 \cup A_3 \cup \dots) = 1 - P(A_1^C \cap A_2^C \cap A_3^C \dots)$$

### Counting

Ways to *order*  $n$  elements ("permutations"):  $n!$

Ways to *partition*  $n$  elements:

(a) in 2 subsets: **(a.1)** defining  $n_1: \binom{n}{n_1}$ ; **(a.2)** Without defining  $n_1: 2^n = \sum_{i=0}^n \binom{n}{i}$ ; (b) in  $K$  subsets: **(b.1)** Defining the  $k_1, k_2, \dots, k_K$  elements of each subset:  $\binom{n}{k_1, k_2, \dots, k_K}$ ; **(b.2)** without defining the number of elements of each subset:  $\left\{ \frac{n}{K} \right\}$  (Sterling); (c) without specifying the number of subsets  $K$ : Bell numbers

Note that the partitioning problem with the  $k$ s all 1 is the problem of ordering a unique set considering each position a "slot".

Ways to sample  $k$  elements from a  $n$  elements bin: **(a)** with replacement: **(a.1)** order matters:  $n^k$ ; **(a.2)** order doesn't matter:  $\binom{n+k-1}{k}$ ; **(b)** without replacement: **(b.1)** order matters:  $k! * \binom{n}{k}$ ; **(b.2)** order doesn't matter:  $\binom{n}{k}$ .

Probability to sample in  $n$  elements of a given type from a bin of  $s$  elements of that type out of total  $k$  elements: **(a)** with

replacement: Binomial( $x; n, s/k$ ); **(b)** without replacement: Hypergeometric( $x; s, k-s, n$ ), i.e.  $\frac{\binom{s}{x} \binom{k-s}{n-x}}{\binom{k}{n}}$  (this reduces to  $\frac{s^x}{k^x}$  for the

probabilities to have *all*  $n$  elements sampled of the given type)

### Distributions

#### Random variable

→ Associate a numerical value to every possible outcome

→ "Discrete" refers to finite or countable infinite values of  $X$ , not necessarily integers

→ "Mixed": those rv that for some ranges are continuous but for some other values have mass concentrated on that values\

- $p_X(x)$ : PMF: Probability Mass Function (discrete)  $P(X = x) = P(\{\omega \in \Omega : X(\omega) = x\}) = p_X(x)$  ("such that")
- $f_X(x)$ : PDF: Probability Density Function (continuous)  $P(a \leq X \leq b) = \int_a^b f_X(x) dx$  (prob per "unit length" - or area, they give the rate at which probability accumulates in the vicinity of a point.) PDF can be discontinue.
- $F_X(x)$ : CDF: Cumulative density function (discrete, continuous or mixed)  $P(X \leq x) = F_X(x)$
- Quantile**( $f$ ) =  $CDF^{-1}(f)$  BUT by convention the  $q_\alpha$  quantile indicates **quantile**( $1 - \alpha$ ) not the quantile of  $\alpha$ .

$\alpha$	0.025	0.05	0.1
$q_\alpha$	1.96	1.645	1.282

◦ "95% CI":  $\mu \pm 1.96 \text{ st.dev}$

$$\bullet \sum_{i=-\infty}^x p_X(i) = F_X(x); \int_{-\infty}^x p_X(i) di = F_X(x); p_X(i) = \frac{dF(x)}{di}$$

→ "Random vector" a multivariate random variable in  $R^k$ . The PDF of the random vector is the joint of all its individual components.

- Gaussian vector** All elements and any linear combination of them is gaussian distributed (e.g. are independent)

#### Discrete distributions:

- Discrete Uniform**: Complete ignorance
- Bernoulli**: Single binary trial

- Binomial**: Number of successes in independent binary trials
- Categorical**: Individual categorical trial
- Multinomial**: Number of successes of the various categories in independent multinomial trials
- Geometric**: Number of independent binary trials until (and including) the first success (discrete time to first success)
- Hypergeometric**: Number of successes sampling without replacement from a bin with given initial number of items representing successes
- Multivariate hypergeometric**: Number of elements sampled in the various categories from a bin without replacement
- Poisson**: Number of independent arrivals in a given period given their average rate per that period length (or, alternatively, rate per period multiplied by number of periods)
- Pascal**: Number of independent binary trials until (and including) the  $n$ -th success (discrete time to  $n$ -th success).

Name	Parameters	Support	PMF	Expectations	Variance	C
<b>D. Unif</b>	$a, b \in \mathbb{Z}$ with $b \geq a$	$x \in \{a, a+1, \dots, b\}$	$\frac{1}{b-a+1}$	$\frac{a+b}{2}$	$\frac{(b-a)(b-a+1)}{12}$	$\frac{x-a+1}{b-a+1}$
<b>Bern</b>	$p \in [0, 1]$	$x \in \{0, 1\}$	$p^x (1-p)^{1-x}$	$p$	$p(1-p)$	$\sum_{i=0}^x p \binom{x}{i} (1-p)^{x-i}$
<b>Bin</b>	$p \in [0, 1], n \in \mathbb{N}^+$	$x \in \{0, \dots, n\}$	$\binom{n}{x} p^x (1-p)^{n-x}$	$np$	$np(1-p)$	$\sum_{i=0}^x \binom{n}{i} p^i (1-p)^{n-i}$
<b>Cat</b>	$p_1, p_2, \dots, p_K$ with $p_k \in [0, 1]$ and $\sum_{k=1}^K p_k = 1$	$x \in \{1, 2, \dots, K\}$	$\prod_{k=1}^K p_k^{x_k}$			
<b>Multin</b>	$n, p_1, p_2, \dots, p_K$ with $p_k \in [0, 1], \sum_{k=1}^K p_k = 1$ and $n \in \mathbb{N}^+$	$x \in \mathbb{N}_0^K$	$\binom{n}{x_1, x_2, \dots, x_K} \prod_{k=1}^K p_k^{x_k}$			
<b>Geom</b>	$p \in [0, 1]$	$x \in \mathbb{N}^+$	$(1-p)^{x-1} p$	$\frac{1}{p}$	$\frac{1-p}{p^2}$	$1 - (1-p)^x$
<b>Hyperg</b>	$n_s, n_f, n \in \mathbb{N}_0$	$x \in \mathbb{N}_0$ with $x \leq n_s$	$\frac{\binom{n_s}{x} \binom{n_f}{n-x}}{\binom{n_s+n_f}{n}}$	$n \frac{n_s}{n_s+n_f}$	$n \frac{n_s}{n_s+n_f} \frac{n_f}{n_s+n_f} \frac{n_s+n_f+n}{n_s+n_f+1}$	
<b>Multiv hyperg</b>	$n_1, n_2, \dots, n_K, n$ with $n \in \mathbb{N}_+, n_i \in \mathbb{N}_0$	$x \in \mathbb{N}_0^K$ with $x_i \leq n_i \forall i, \sum_{i=1}^K x_i = n$	$\frac{\prod_{i=1}^K \binom{n_i}{x_i}}{\binom{\sum_{i=1}^K n_i}{n}}$	$n \frac{n_i}{\sum_{i=1}^K n_i}$	$n \frac{\sum_{j=1}^K n_j - n}{\sum_{j=1}^K n_j - 1} \frac{n_i}{\sum_{j=1}^K n_j} \left(1 - \frac{n_i}{\sum_{j=1}^K n_j}\right)$	
<b>Pois</b>	$\lambda \in \mathbb{N}^+$	$x \in \mathbb{N}_0$	$\frac{\lambda^x e^{-\lambda}}{x!}$	$\lambda$	$\lambda$	
<b>Pasc</b>	$n \in \mathbb{N}^+, p \in [0, 1]$	$x \in [n, n+1, \dots, \infty)$	$\binom{x-1}{n-1} p^n (1-p)^{x-n}$	$\frac{n}{p}$	$\frac{n(1-p)}{p^2}$	

#### Continuous distributions:

- Uniform** Complete ignorance, pick at random, all equally likely outcomes
- Exponential** Waiting time to first event whose rate is  $\lambda$  (continuous time to first success)
- Laplace** Difference between two iid exponential r.v.
- Normal** The asymptotic distribution of a sample means
- Erlang** Time of the  $n$ -th arrival
- Cauchy** The ratio of two independent zero-means normal r.v.
- Chi-squared** The sum of the squared of iid standard normal r.v.
- T distribution** The distribution of a sample means
- F distribution**: The ratio of the ratio of two indep  $X^2$  r.v. with their relative parameter

- **Beta distribution** The Beta distribution
- **Gamma distribution** Generalisation of the exponential, Erlang and chi-square distributions

Name	Parameters	Support	PMF	Expectations	Variance	CDF
<b>Unif</b>	$a, b \in \mathbb{R}$ with $b \geq a$	$x \in [a, b]$	$\frac{1}{b-a}$	$\frac{a+b}{2}$	$\frac{(b-a)^2}{12}$	$\frac{x-a}{b-a}$
<b>Expo</b>	$\lambda \in \mathbb{R}^+$	$x \in \mathbb{R}^+$	$\lambda e^{-\lambda x}$	$\frac{1}{\lambda}$	$\frac{1}{\lambda^2}$	$1 - e^{-\lambda x}$
<b>Laplace</b>	$\mu \in \mathbb{R}$ (location), $b \in \mathbb{R}^+$ (scale)	$x \in \mathbb{R}$	$\frac{1}{2b} e^{-\frac{ x-\mu }{b}}$	$\mu$	$2b^2$	
<b>Normal</b>	$\mu \in \mathbb{R}$ , $\sigma^2 \in \mathbb{R}^+$	$x \in \mathbb{R}$	$\frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$	$\mu$	$\sigma^2$	
<b>Multiv. Normal</b>	$\mu \in \mathbb{R}^d$ , $\Sigma \in \mathbb{R}^{d \times d}$	$x \in \mathbb{R}^d$	$\frac{1}{\sqrt{(2\pi)^d \det(\Sigma)}} e^{-\frac{1}{2}(x-\mu)' \Sigma^{-1}(x-\mu)}$	$\mu$	$\Sigma$	
<b>Erlang</b>	$n \in \mathbb{N}^+$ , $\lambda \in \mathbb{R}^+$	$x \in \mathbb{R}_+$	$\frac{\lambda^n x^{n-1} e^{-\lambda x}}{(n-1)!}$	$\frac{n}{\lambda}$	$\frac{n}{\lambda^2}$	
<b>Cauchy</b>	$x_0 \in \mathbb{R}$ (location), $\gamma \in \mathbb{R}^+$ (scale)	$x \in \mathbb{R}$	$\frac{1}{\pi\gamma(1+(\frac{x-x_0}{\gamma})^2)}$	NDEF	NDEF	
<b>Chi-sq</b>	$d \in \mathbb{N}^+$	$x \in \mathbb{R}^+$	$\frac{1}{2^{\frac{d}{2}} \Gamma(\frac{d}{2})} x^{\frac{d}{2}-1} e^{-\frac{x}{2}}$	$d$	$2d$	
<b>T</b>	$\nu \in \mathbb{R}^+$	$x \in \mathbb{R}$	$\frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\nu\pi}\Gamma(\frac{\nu}{2})} \left(1 + \frac{x^2}{\nu}\right)^{-\frac{\nu+1}{2}}$			
<b>F</b>	$d_1 \in \mathbb{N}^+$ , $d_2 \in \mathbb{N}^+$	$x \in \mathbb{R}^+$	$\frac{\sqrt{\frac{(d_1 x)^{d_1} d_1^{d_1}}{(d_1 x + d_2)^{d_1+d_2}}}}{xB(\frac{d_1}{2}, \frac{d_2}{2})}$	$\frac{d_2}{d_2-2}$ for $d_2 > 2$	$\frac{2d_2^2(d_1+d_2-2)}{d_1(d_2-2)^2(d_2-4)}$ for $d_2 > 4$	
<b>Beta</b>	$\alpha, \beta \in \mathbb{R}^+$	$x \in [0, 1]$	$\frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}$	$\frac{\alpha}{\alpha+\beta}$	$\frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$	
<b>Gamma</b>	$\alpha \in \mathbb{R}^+$ (shape), $\beta \in \mathbb{R}^+$ (rate)	$x \in \mathbb{R}^+$	$\frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x}$	$\frac{\alpha}{\beta}$	$\frac{\alpha}{\beta^2}$	

**Beta function** :  $B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)} = \frac{\alpha+\beta}{\alpha\beta}$

**Gamma function**:  $\Gamma(x) = (x-1)! \forall x \in \mathbb{N}$

#### Expected value

→ The mean we would get running an experiment many times

- $E[X] := \sum_x x p_X(x) := \int_{-\infty}^{+\infty} x f_X(x) dx$
- **Expected value rule**:  $E[Y = g(X)] = \sum_y Y p_Y(y) = \sum_x g(x) p_X(x) \neq g(\sum_x x p_X(x)) = g(E[X])$  (in general)
- **Linearity of expectations**:  $E[aX + b] = aE[X] + b$ ;  $E[X + Y + Z] = E[X] + E[Y] + E[Z]$
- **X, Y independent**:  $E[XY] = E[X]E[Y]$ ,  $E[g(X)h(Y)] = E[g(X)]E[h(Y)]$
- **Law of Iterated Expectations**:  $E[E[X|Y]] := \sum_Y E[X|Y] p_Y(y) = E[X]$  ( $E[X|Y]$  is seen as a function  $g(Y)$ )
- Expectations of convex functions are convex
- The expectations of an indicator function is the prob that the event indicated is true

#### Variance

- $Var(X) := E[(X - \mu)^2] = \sum_x (x - \mu)^2 p_X(x) = \int_{-\infty}^{+\infty} (x - \mu)^2 f_X(x)$
- $Var(X) = E[X^2] - (E[X])^2$
- $Var(g(X)) = E[g(X)^2] - (E[g(X)])^2$
- $Var[aX + b] = a^2 Var[X]$ ;
- var of sum of r.v.:
  - X, Y independent:  $Var(X + Y) = Var(X) + Var(Y)$
  - in general:  $Var(\sum_{i=1}^n X_i) = \sum_{i=1}^n Var(X_i) + 2 \sum_{i < j} Cov(X_i, X_j) = \sum_{i,j} Cov(X_i, X_j)$
- **Law of total Variance**:  $var(X) = E[var(X|Y)] + var(E[X|Y])$  (expected value of the variances of X *within* each group Y + variance of the means *between* groups)

#### Moments

- $M_n = \int_{-\infty}^{\infty} (x - c)^n f_X(x) dx$
- $c = 0 \rightarrow$  "raw moment" (e.g.  $E[X]$ ,  $E[X^2]$ , ...)
- $c = \mu \rightarrow$  "central moment" (e.g.  $var(X)$  (second), skewness (third), kurtosis (fourth))
- "Mode"  $\rightarrow \arg\max(f_X)$
- "Median"  $\rightarrow k : \int_{-\infty}^k f_X(x) dx = \int_k^{\infty} f_X(x) dx$
- "Mean"  $\rightarrow$  the expected value

#### Covariance and correlation

- 2 r.v.:  $Cov(X, Y) := E[(X - E[X])(Y - E[Y])] = E[(X)(Y - E[Y])] = E[XY] - E[X]E[Y]$
- $Cov(aX + bY + c, eZ) = ae Cov(X, Z) + be Cov(Y, Z)$
- X random vector:
  - $Cov(X) := E[(X - E[X])(X - E[X])'] = E[XX'] - E[X]E[X]'$
  - $Cov(AX + b) = ACov(X)A'$  (all cov matrix are positive definite and so it's ok to take square roots of them)
- Correlation coeff.:  $\rho := \frac{cov(X, Y)}{\sqrt{var(X)var(Y)}} = E[\frac{X - E[X]}{\sigma_X} * \frac{Y - E[Y]}{\sigma_Y}]$  with  $-1 \leq \rho \leq +1$  and  $\sigma_X, \sigma_Y \neq 0$ 
  - $(X, Y)$  indep.  $\rightarrow cov(X, Y) = 0 \leftrightarrow \rho = 0$  (but not  $\leftarrow$ )
  - $|\rho| = 1 \leftrightarrow (X - E[X]) = c(Y - E[Y])$  (i.e. X, Y linearly correlated)
  - $\rho(aX + b, Y) = sign(a) * \rho(X, Y)$  (because of dimensionless)
  - $X = Z + V, Y = Z + W, Z, V, W$  indep.  $\rightarrow \rho(X, Y) = \frac{\sigma_Z^2}{\sigma_Z^2 + \sigma_Z\sigma_W + \sigma_V\sigma_Z + \sigma_V\sigma_W}$

#### Normal RV

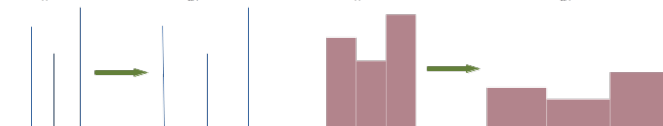
$X \sim N(\mu, \sigma^2) \rightarrow Y = aX + b \sim N(a\mu + b, a^2\sigma^2)$   
 $X \sim N(\mu, \sigma^2), Z \sim N(0, 1) \rightarrow P(a \leq X \leq b) = P(\frac{a-\mu}{\sigma} \leq \frac{X-\mu}{\sigma} \leq \frac{b-\mu}{\sigma}) = P(\frac{a-\mu}{\sigma} \leq Z \leq \frac{b-\mu}{\sigma}) = \Phi(\frac{b-\mu}{\sigma}) - \Phi(\frac{a-\mu}{\sigma})$   
 $X \sim N(0, \sigma^2) \rightarrow P(X < -a) = 1 - P(X < a)$   
 $X_i \sim N(\mu_i, \sigma_i^2), X_i$  i.i.d.  $\rightarrow Y = \sum_i X_i \sim N(\sum_i \mu_i, \sum_i \sigma_i^2)$   
 $E[Z] = 0, E[Z^2] = 1, E[Z^3] = 0, E[Z^4] = 3, E[Z^5] = 0, E[Z^6] = 15$

#### Derived Distributions

##### Function of a single R.V.

**Linear function of a r.v.**:  $Y = aX + b$

- $p_Y(y) = p_X(\frac{y-b}{a})$  (where  $\frac{y-b}{a}$  is the value of X that raises y)
- $f_Y(y) = f_{aX+b}(y) = \frac{1}{|a|} f_X(\frac{y-b}{a})$  (area must be constant)



**Monotonic**:  $Y = g(x)$  with  $g(x)$  monotonic and continuous

$$f_Y(y) = f_X(g^{-1}(y)) * |\frac{dg^{-1}}{dy}(y)|$$

##### Probability Integral Transformation

Considering as "function" the CDF, this is uniformly distributed for any r.v.:  $Y = g(X) = CDF_X(X) \sim U(0, 1)$

$$Y = F_X(x) \rightarrow F_Y(y) = P(Y \leq y) = P(F_X(x) \leq y) = P(X \leq F_X^{-1}(y)) = F_X(F_X^{-1}(y)) = y$$

For a sample  $k \sim U(0, 1)$  the corresponding value over X is  $x = F_X^{-1}(k)$ , i.e. the inverse CDF evaluated on k.

**General**:  $Y = g(X)$

$$p_Y(y) = \sum_{\{x: g(x)=y\}} p_X(x)$$

$$f_Y(y) = \frac{dF_Y}{dy}(y) \text{ with } F_Y(y) = P(g(X) \leq y) = \int_{\{x: g(x) \leq y\}} f_X(x) dx \text{ (express the CDF of Y in terms of the CDF of X and then derive to find the PDF)}$$

##### Function of multiple R.V.

**Sum of 2 independent R.V., discrete::**

$$\bullet Z = X + Y \quad p_Z(z) = \sum_x p_X(X = x) p_{Z|X}(Z = z | X = x) = \sum_x p_X(x) p_Y(z - x)$$

**Sum of 2 independent R.V., continue (convolution):**

$$\bullet Z = z \text{ in all occasions where } X = x \text{ and } Y = z - x$$

$$\bullet f_Z(z) = \int_{\max(x_{\min}, z - y_{\max})}^{\min(x_{\max}, z - y_{\min})} f_X(x) * f_Y(z - x) dx$$

**General:**  $Z = g(X, Y, \dots)$  Find (e.g. geometrically) the CDF of  $Z$  and differentiate for  $Z$  to find the PDF.

**Sum of random number of i.i.d. R.V.**  $Y = \sum_{i=1}^N X_i$  with  $X_i \forall i$  i.i.d and indep to  $N$

$$E[Y] = E[E[Y|N]] = E[N * E[X]] = E[N] * E[X]$$

$$\text{var}(Y) = E[N] * \text{var}(X) + (E[X])^2 * \text{var}(N)$$

$$X \sim \text{bern}(p); N \sim \text{bin}(m, q) \rightarrow Y \sim \text{bin}(m, pq)$$

$$X \sim \text{bern}(p); N \sim \text{pois}(\lambda) \rightarrow Y \sim \text{pois}(p\lambda)$$

$$X \sim \text{geom}(p); N \sim \text{geom}(q) \rightarrow Y \sim \text{geom}(pq)$$

$$X \sim \text{exp}(\lambda); N \sim \text{geom}(q) \rightarrow Y \sim \text{exp}(q\lambda)$$

### Order statistics

$$Y = \max(X), X \text{ i.i.d} \rightarrow F_Y(y) = P(X_i \leq y) \forall i \in [1, N] = F_X(y)^N \rightarrow f_Y(y) = N F_X(y)^{N-1} f_X(y)$$

$$Y = \min(X), X \text{ i.i.d} \rightarrow F_Y(y) = 1 - P(X_i \geq y) \forall i \in [1, N] = (1 - (1 - F_X(y))^N) \rightarrow f_Y(y) = n(1 - N F_X(y))^{N-1} f_X(y)$$

### Limits

**Properties of the sample mean:**  $\bar{X}_n = \frac{\sum_{i=1}^n X_i}{n} \rightarrow E[\bar{X}_n] = E[X], \text{var}(\bar{X}_n) = \text{var}(X)/n$  (from properties of expectation and variance)

**Markov Inequality:**  $X$  non neg r.v.,  $t > 0 \rightarrow P(X \geq t) \leq E[X]/t$

**Chebyshev Inequality:**  $X$  a r.v.,  $t > 0 \rightarrow P(|X - E[X]| \geq t) \leq \text{Var}(X)/t^2$

- proof: from Markov in. by considering a new r.v.  $Y = (X - E[X])^2$
- corollary: the prob that a r.v. is  $k$  st.dev. away from the mean is less than  $1/k^2$ , whatever its distribution
- the  $t$  is the "accuracy" and the probability itself is the "confidence" in reaching the given accuracy

**Hoffding's Inequality:**  $X_1, X_2, \dots, X_n$  i.i.d. with  $E[X_i] = \mu$  and  $X \in [a, b]$  almost surely, and  $a < b \rightarrow P(|\bar{X}_n - \mu| \geq \epsilon) \leq$

$$2e^{-\frac{2n\epsilon^2}{(b-a)^2}} \quad \forall \epsilon > 0 \text{ (with } n \text{ not necessarily large)}$$

### Def of Convergences

We are interested in r.v. whose distribution is parametrised by  $n$ , i.e. in sequences of a r.v.

**Of a deterministic sequence::** The sequence  $a_n$  converge to the value  $a$  if for any  $\epsilon > 0$  it exists a value  $n_0$  such that  $|a_n - a| \leq \epsilon \quad \forall n \geq n_0$ , i.e. whatever small we choose  $\epsilon$ , we can always find a  $n$  limit where subsequent sequences values are less than  $\epsilon$  far away to  $a$

**In distribution:** A sequence  $Y_n$  of r.v. (not necessarily indep.) converges in distribution to the r.v.  $Y$  iff  $\lim_{n \rightarrow \infty} E[f(Y_n)] = E[f(Y)]$  for any function  $f$  continuous and bounded (so, it's not true in general, and in particular you may have convergence in distribution without having  $E[Y_n]$  converging to  $E[Y]$ ) or, equivalently, that  $Y$  iff  $\lim_{n \rightarrow \infty} P(Y_n \leq y) = P(Y \leq y)$ . That is, the  $Y_n$  CDF converges to a limiting CDF. Converging of a function. The output of the function depends on  $n$ , but with less and less variations. Base for the CLT. Aka the "convergence in law" or "weak convergence".

**In probability:** A sequence  $Y_n$  of r.v. (not necessarily indep.) converges in probability to the value  $a$  if for every  $\epsilon > 0$  we have  $\lim_{n \rightarrow \infty} P(|Y_n - a| \geq \epsilon) = 0$ . Base for the Weak L.L.N.

- The "bulk" of the PMF is within a range of  $a$ , but nothing is said for the values that are not in the range  $\Rightarrow$  conv. in p. doesn't implies conv. in expectations (as expectations are sensitive to the tail of the distribution).
- To verify, first guess the  $\epsilon$  (or the "a" ???) and then find the limit of the CDF for the relevant bound goes to 0

**With probability 1 (almost sure):** A sequence  $Y_n$  of r.v. (not necessarily indep.) converges w.p. 1 to  $c$  if  $P(\lim_{n \rightarrow \infty} Y_n = c) = 1$  or, more in general,  $P(\{w : \lim_{n \rightarrow \infty} Y_n(w) = Y(w)\}) = 1$ . Base for the strong L.L.N. The closest version to a deterministic convergence, whatever is the outcome of a random experiment all they have to converge.

### Convergence theorems:

- $X_n \xrightarrow{p/a.s.} a, Y_n \xrightarrow{p/a.s.} b \Rightarrow X_n + Y_n \xrightarrow{p/a.s.} a + b; X_n * Y_n \xrightarrow{p/a.s.} a * b$
- Continuous Mapping Theorem:  $X_n \xrightarrow{d/p/a.s.} a, g$  is a continuous function,  $\Rightarrow g(X_n) \xrightarrow{d/p/a.s.} g(a)$
- Slutsky's Theorem:  $X_n \xrightarrow{d} X, Y_n \xrightarrow{p/a.s.} y$

$$\circ X_n + Y_n \xrightarrow{d} X + y$$

$$\circ X_n * Y_n \xrightarrow{d} X * y$$

**Law of large numbers:** The sample mean converge to the pop mean

- weak L.L.N.:  $\lim_{n \rightarrow \infty} P(|\bar{X}_n - E[X]| > \epsilon) = 0$ 
  - from the Chebyshev Inequality by using  $\bar{X}_n$  and taking the limit for  $n \rightarrow \infty$
  - it is a convergence in p. of  $\bar{X}_n$  to  $E[X]$ .
- strong L.L.N.:  $\lim_{n \rightarrow \infty} P(\bar{X}_n = E[X]) = 1$

Note that given  $\bar{X}_n = \frac{1}{n} \sum X_i$ :

- $\bar{X}_n \xrightarrow{n \rightarrow \infty} E[X]$  (LLN)
- $g(\bar{X}_n) \xrightarrow{n \rightarrow \infty} g(E[X])$  (cont. map. theorem)
- $Y = g(X) \Rightarrow \bar{Y}_n = \frac{1}{n} \sum Y_i = \frac{1}{n} \sum g(X_i) \xrightarrow{n \rightarrow \infty} E[Y] = E[g(X)]$  (LLN)

**Central Limit Theorem:** The distribution of the mean from i.i.d. samples converges in distribution to a Normal distribution with mean equal to the population mean and variance of the population variance divided by the sample size:

$$\bar{X}_n \sim N(E[X], \sigma_X^2/n)$$

- formally the CLT is stated in terms of  $\frac{\bar{X}_n - E[X]}{\sigma_X/\sqrt{n}} \xrightarrow{\text{dist}} N(0, 1)$
- multivariate CLT:  $\sqrt{n} * \Sigma_X^{-\frac{1}{2}} (\bar{X}_n - \mu) \xrightarrow{\text{dist}} N_d(0, I_d)$
- versions exists for identically distributed  $X_i$  or "weakly dependent" ones (dependence only local between neighbour  $X_i$ )
- $X$  integer: consider  $S_{n+1/2}$
- Approximation to the binomial:  $P(k \leq S_n \leq l) \approx \Phi\left(\frac{l + \frac{1}{2} - np}{\sqrt{n(1-p)}}$

### Bernoulli and Poisson random processes

Stochastic processes: a probabilistic phenomenon that evolves in time, i.e. an infinite sequence of r.v.

We need to characterise it with informations on the individual r.v. but also on how they relate (joint)

Bernoulli, Poisson  $\rightarrow$  Assumptions: independence ( $\rightarrow$  memoryless), time-homogeneity

	Bernoulli	Poisson
Time of arrival $\tau$	Discrete	Continuous
Arrival rate	$p$ per trial	$\lambda$ per unit time
N# of arrivals	Binomial $p_n(n; \tau, p)$	Poisson $p_n(n; \tau; \lambda)$
Interarrival time	Geometric $p_\tau(\tau; p)$	Exponential $f_\tau(\tau; \lambda)$
Time to $n^{\text{th}}$ arrival	Pascal $p_\tau(k, p)$	Erlang $f_\tau(\tau; n, \lambda)$

Fresh start: The Bernoulli or Poisson process after time  $N$ , where  $N$  is a r.v. causally determined from the history of the process, is a new Bernoulli/Poisson process with the same probabilistic characteristics as the original one.

### The poisson as approximation of the binomial

Given  $p$  the probability of a successes in a single slot and  $n$  the number of slots, the expected number of successes  $\lambda$  is given by  $\lambda = pn$ .

The poisson PDF can be seen as the limit of the Bernoulli pdf when we consider smaller and smaller time slots, keeping constant the total expected number of successes for the period (that is  $p$  - on the single period - becomes smaller and smaller and the number of periods tends to  $\infty$ ).

The poisson process can hence be seen as a limiting case of a Bernoulli process or, alternatively, as the process deriving from a sequence of exponential r.v..

In a small interval  $\delta$ , the probability of 1 success is  $\lambda\delta$  and of 0 successes is  $1 - \lambda\delta$  (and negligible probabilities for more than one).

### Merging

The process made by a sequence of r.v. functions of other sequences of r.v.

### Merging of Bernoulli processes

$$X_1^i \sim \text{Bern}(p), X_2 \sim \text{Bern}(q), X_1 \text{ indep } X_2$$

$$Y^i = X_1^i \text{ or } X_2^i \Rightarrow Y^i \sim \text{Bern}(p + q - pq)$$

$$Y^i = X_1^i \text{ and } X_2^i \Rightarrow Y^i \sim \text{Bern}(pq) \text{ (both new Bernoulli processes)}$$

The probability that observing a success in the merged process we have a success also in the original process 1 is:

$$\begin{aligned} \bullet Y^i = X_1^i \text{ or } X_2^i &\Rightarrow P(X_1^i | Y^i) = \frac{P(X_1^i, Y^i)}{P(Y^i)} = \frac{P(X_1^i)}{P(Y^i)} = \frac{p}{p+q-pq} \\ \bullet Y^i = X_1^i \text{ and } X_2^i &\Rightarrow P(X_1^i | Y^i) = 1 \end{aligned}$$

#### Merging of Poisson processes

$$X_1^i \sim \text{Poisson}(\lambda_1), X_2 \sim \text{Poisson}(\lambda_2), X_1 \text{ indep } X_2$$

Note that differently from Bernoulli case, here the change of a match is zero.

$$Y^i = X_1^i \text{ or } X_2^i \Rightarrow Y^i \sim \text{Poisson}(\lambda_1 + \lambda_2)$$

The probability that observing a success in the merged process we have a success also in the original process 1 is:

$$\bullet Y^i = X_1^i \text{ or } X_2^i \Rightarrow P(X_1^i | Y^i) = \frac{P(X_1^i, Y^i)}{P(Y^i)} = \frac{\lambda_1}{\lambda_1 + \lambda_2}$$

#### Splitting

##### Splitting of Bernoulli processes

Given a Bernoulli process  $X$  with probability  $p$  and an other independent Bernoulli process  $Y$  that assigns each success of  $X$  to  $Z_1$  with probability  $q$  and to  $Z_2$  with probability  $(1 - q)$ , we have:

$$\begin{aligned} \bullet Z_1^i &= X^i \text{ and } Y^i, Z_1^i \sim \text{Bern}(pq) \\ \bullet Z_2^i &= X^i \text{ and not } Y^i, Z_2^i \sim \text{Bern}(p(1 - q)) \end{aligned}$$

$Z_1$  and  $Z_2$  are *not* independent !

##### Splitting of Poisson processes

Given a Poisson process  $X$  with rate  $\lambda$  and an other independent Bernoulli process  $Y$  that assigns each success of  $X$  to  $Z_1$  with probability  $q$  and to  $Z_2$  with probability  $(1 - q)$ , we have:

$$\begin{aligned} \bullet Z_1^i &= X^i \text{ and } Y^i, Z_1^i \sim \text{Poisson}(\lambda q) \\ \bullet Z_2^i &= X^i \text{ and not } Y^i, Z_2^i \sim \text{Poisson}(\lambda(1 - q)) \end{aligned}$$

Note that differently from Bernoulli case, here the two processes *are* independent, as the probability of an arrival at any given *point* in time is zero.

##### Summing Poisson rv

Given  $X_1 \sim \text{Poisson}(p)$  (the distribution, not the process) and  $X_2 \sim \text{Poisson}(q)$ , and  $X_1, X_2$  i.i.d.  $\Rightarrow Y = (X_1 + X_2) \sim \text{Poisson}(p + q)$  (think as the two input r.v. as representing numbers of arrivals in disjoint time intervals).

##### Random incidence ("Inspection paradox")