

# PT = POPULAR DISTRIBUTIONS

## I - DISCRETE DISTRIBUTIONS

- Bernoulli trials = two possible outcomes "success" or "failure"  
 $\Omega = \{S, F\}$  and put  $X(S) = 1$  (S)  $X(F) = 0$  (F)

Let  $p =$  probability of success

Then  $P(X=1) = p_X(1) = p = 1 - p_X(0)$

We write  $X \sim B(p)$

the Bernoulli distribution has a single parameter.

x Moments:

•  $EX = 1 p_X(1) + 0 p_X(0) = p$

•  $EX^2 = 1^2 p_X(1) + 0^2 p_X(0) = p \Rightarrow \text{Var } X = p(1-p)$   
 $(\geq 0)$

$$X \sim B(p) \quad p_X(1) = p = 1 - p_X(0)$$

$$EX = p$$

$$\text{Var } X = p(1-p)$$

- Binomial distribution = consider a random experiment consisting of  $n$  independent Bernoulli trials with proba of success  $p$ .

•  $\Omega =$  set of all sequences of the form  $\omega = \underbrace{SSFS \dots FS}_n$

•  $\mathcal{F} =$  power set of  $\Omega$

•  $P$  such that  $P(\{\omega\}) = p^{\# \text{successes}} (1-p)^{\# \text{failures}}$   
 $\uparrow$   
 since we have independent events

We are interested in  $X :=$  number of successes.

The distribution of  $X$  is given by:

(2)

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

Digression. Consider a set  $A = \{a_1, \dots, a_n\}$  of  $n$  distinct objects, and take a sample of size  $k$  from it. Call it  $(a_{j_1}, \dots, a_{j_k})$

$\uparrow$  there are  $n$  ways to choose  $j_1$   $\leftarrow$   $(n-1)$  ways to choose  $j_2$   $\dots$

the number of possible samples of size  $k$  is

$$n(n-1) \times \dots \times (n-k+1) = \frac{n!}{(n-k)!}$$

In this case, the order of selection is important = for example, samples  $(a_1, a_6, a_3)$  and  $(a_6, a_1, a_3)$  are different.

If we are interested in the number of unordered samples, we need to remove samples counted multiple times. To do so, we need to count the number of permutations of  $(a_{j_1}, \dots, a_{j_k}) \rightarrow$  there are  $k!$  of them.

Conclusion: the number of possible unordered samples of size  $k$  taken from a set of  $n$  distinct objects is

$$\frac{n!}{k!(n-k)!} = \binom{n}{k}$$

Back to the distribution of  $X$ , we see that

$$p_X(x) = P(X=x) = \binom{n}{x} p^x (1-p)^{n-x}$$

We say that  $X$  has a binomial distribution with parameters  $n$  and  $p$ , and we write  $X \sim \text{Bi}(n, p)$ .

Remark:  $\sum_{x=0}^n p_X(x) = \sum_{x=0}^n \binom{n}{x} p^x (1-p)^{n-x} = (p+1-p)^n = 1$

x Moments:  $EX = \sum_{x=0}^n \binom{n}{x} x p^x (1-p)^{n-x}$  (3)

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

$$= \sum_{x=1}^n \frac{(n-1)!}{(x-1)!(n-x)!} n p p^{x-1} (1-p)^{n-x} \quad \text{let } y=x-1$$

$$= n p \sum_{y=0}^{n-1} \frac{(n-1)!}{y!(n-1-y)!} p^y (1-p)^{n-1-y}$$

$$= n p = (p+1-p)^{n-1} = 1$$

$EX^2 = \sum_{x=0}^n \binom{n}{x} x^2 p^x (1-p)^{n-x}$

$$= \sum_{x=2}^n \binom{n}{x} x(x-1) p^x (1-p)^{n-x} + EX$$

$$= n(n-1) p^2 \sum_{x=2}^n \frac{(n-2)!}{(x-2)!(n-x)!} p^{x-2} (1-p)^{n-x} \quad \text{let } y=x-2$$

$$= n(n-1) p^2 \sum_{y=0}^{n-2} \binom{n-2}{y} p^y (1-p)^{n-2-y} + EX$$

$$= n(n-1) p^2 + n p = (p+1-p)^{n-2} = 1$$

$Var X = EX^2 - (EX)^2 = n(n-1)p^2 + np - n^2p^2 = np(1-p)$

$X \sim Bi(n, p) \quad p_X(x) = \binom{n}{x} p^x (1-p)^{n-x}$   
 $EX = np \quad x = 0, \dots, n$   
 $Var X = np(1-p)$

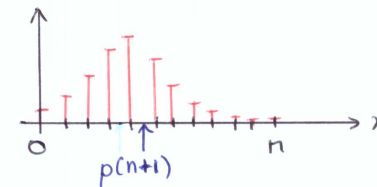
x Shape of the binomial distribution. Let's look at the ratio of successive binomial probabilities:

$$r(x) = \frac{p_X(x)}{p_X(x-1)} = \frac{\binom{n}{x} p^x (1-p)^{n-x}}{\binom{n}{x-1} p^{x-1} (1-p)^{n-x+1}} = \frac{\frac{n+1-x}{x} - 1}{\frac{1}{p} - 1}$$

compare this ratio with 1

If  $\frac{n+1}{x} > \frac{1}{p}$  i.e. if  $x < p(n+1)$ ,  $r(x) > 1$  pmf is  $\uparrow$  (4)

If  $\frac{n+1}{x} < \frac{1}{p}$  i.e. if  $x > p(n+1)$ ,  $r(x) < 1$  pmf is  $\downarrow$   
 $\Rightarrow$  the binomial distribution has a single peak, and typically looks like this.



Multinomial distribution = generalization of the Binomial distribution. We consider  $n$  independent trials, each trial consisting in the success of exactly one of  $k$  categories. The probability of success of category  $i$  is  $p_i$ , such that  $\sum_{i=1}^k p_i = 1$ .

We are interested in  $X = (X_1, \dots, X_k)$ , where  $X_i$  denotes the number of times category  $i$  is observed out of  $n$  trials.

Let  $x = (x_1, \dots, x_k)$ . Then

$$p_X(x) = P(X_1 = x_1, \dots, X_k = x_k) \propto p_1^{x_1} \dots p_k^{x_k}$$

we need to count the number of ways of selecting such outcomes.

Digression. the number of ways of partitioning a set of  $n$  distinct objects into  $k$  cells with  $r_1$  objects in the first cell,  $r_2$  objects in the second cell, ... is

$$\binom{n}{r_1, \dots, r_k} = \binom{n}{r_1} \binom{n-r_1}{r_2} \dots \binom{n-r_1-\dots-r_{k-1}}{r_k}$$

Choose  $r_1$  objects out of  $n$

then choose  $r_2$  objects out of  $n-r_1$

... / ...

$$= \frac{n!}{r_1!(n-r_1)!} \times \frac{(n-r_1)!}{r_2!(n-r_1-r_2)!} \times \dots \times \frac{(n-r_1-\dots-r_{k-1})!}{r_k!(n-r_1-\dots-r_k)!} \quad (5)$$

$\sum_{i=1}^k r_i = n$

$$\binom{n}{r_1, \dots, r_k} = \frac{n!}{r_1! \dots r_k!}$$

Thus  $p_X(x) = \binom{n}{x_1, \dots, x_k} p_1^{x_1} \dots p_k^{x_k}$

We say that  $X$  has a multinomial distribution, and we write  $X \sim MN(n, p_1, \dots, p_k)$ .

Can check that the marginal distribution of  $X_i$  is  $Bi(n, p_i)$ , so that  $EX_i = np_i$ ,  $Var X_i = np_i(1-p_i)$ .

Check also that  $Cov(X_i, X_j) = -np_i p_j$  for  $i \neq j$ .

$$X \sim MN(n, p_1, \dots, p_k) \quad p_X(x) = \binom{n}{x_1, \dots, x_k} p_1^{x_1} \dots p_k^{x_k},$$

$$EX_i = np_i \quad \begin{cases} \sum_{i=1}^k x_i = n \\ x_i \geq 0 \end{cases}$$

$$Var X_i = np_i(1-p_i)$$

$$Cov(X_i, X_j) = -np_i p_j$$

• Geometric distribution. Consider (an infinite) sequence of Bernoulli trials, each with probability of success  $p$ .

Put  $X :=$  number of failures before the first success.

EX:  $w_1 = FFSSFSS \dots \quad X(w_1) = 2$   
 $w_2 = SFFSFF \dots \quad X(w_2) = 0$

Clearly,  $p_X(x) = P(X=x) = p(1-p)^x$   
 "x failures followed by a success"

It is indeed a probability distribution since  $\sum_{x=0}^{\infty} p_X(x) = p \sum_{x \geq 0} (1-p)^x = p \frac{1}{1-(1-p)} = 1$ .

We say that  $X$  has a geometric distribution, and we write  $X \sim G(p)$ .

• Moments.  $EX = \sum_{x \geq 0} x p (1-p)^x$   
 $= p(1-p) \sum_{x \geq 1} x (1-p)^{x-1}$

digression. Put  $f(y) = \sum_{x \geq 0} y^x$ . Then for any  $y$  such that  $|y| < 1$ , we can exchange summation and derivation:

$$f'(y) = \left( \sum_{x \geq 0} y^x \right)' = \sum_{x \geq 0} (y^x)' = \sum_{x \geq 1} x y^{x-1}$$

$$\left( \frac{1}{1-y} \right)' = \frac{1}{(1-y)^2} \Rightarrow \sum_{x \geq 1} x y^{x-1} = \frac{1}{(1-y)^2}$$

$$\Rightarrow EX = p(1-p) \frac{1}{(1-(1-p))^2} = \frac{1-p}{p}$$

•  $EX^2 = \sum_{x \geq 0} x^2 p (1-p)^x$   
 $= p(1-p)^2 \sum_{x \geq 2} x(x-1)(1-p)^{x-2} + EX$  ↗  $x^2 = x(x-1) + x$

digression: apply the same trick to the second order derivative:  $\sum_{x \geq 2} x(x-1)y^{x-2} = \frac{2}{(1-y)^3}$

$$= \frac{2(1-p)^2}{p^2} + \frac{1-p}{p}$$

•  $\text{Var } X = EX^2 - (EX)^2 = \frac{1-p}{p^2}$  (7)

$X \sim g(p) \quad p_X(x) = p(1-p)^x, \quad x \geq 0$

$EX = \frac{1-p}{p}$

$\text{Var } X = \frac{1-p}{p^2}$

Alternatively, you may define  $X$  to be the total number of trials until the occurrence of the first success: we are shifting the distribution to the right. The mean is then  $1/p$ , but the variance is unchanged.

× lack of memory. First, we compute  $P(X \geq x) = \sum_{i=x}^{\infty} p(1-p)^i$   
 $= \sum_{j \geq 0} p(1-p)^{j+x} = p(1-p)^x \sum_{j \geq 0} (1-p)^j = (1-p)^x$   
 ( $j = i - x$ )

Next, we are interested in

$P(X \geq x+y \mid X \geq x) = \frac{P(X \geq x+y)}{P(X \geq x)}$

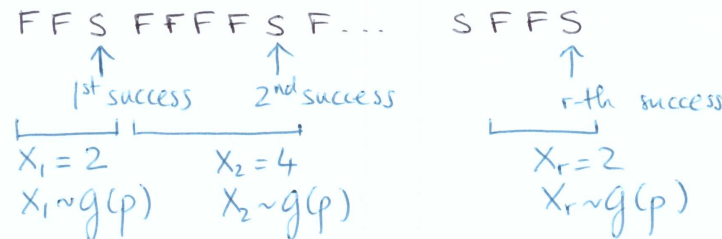
$= \frac{(1-p)^{x+y}}{(1-p)^x}$

$= (1-p)^y$   
 $= P(X \geq y)$

← independent of  $x$ .

information that there were no successes in the first  $x$  trials does not carry any information regarding the waiting time until the first success in the future.  
 "past has no effect on future" "process forgets about its past".

• Negative binomial distribution. Similar setup as for the geometric distribution, except that we are interested here in the total number of failures before the occurrence of the  $r$ -th success. Denote this number by  $X$ . (8)



$\Rightarrow X = X_1 + \dots + X_r$ , where  $X_i \sim g(p), 1 \leq i \leq r$

× Distribution of  $X$ . Consider the event  $\{X = x\}, x \geq 0$



Rearranging terms, we see that there are  $\binom{x+r-1}{r-1}$  of such arrangements, leaving the final  $r$ -th success on its own. The probability of one such sequence is  $(1-p)^x p^r$

Thus  $p_X(x) = P(X=x) = \binom{x+r-1}{r-1} p^r (1-p)^x, x \geq 0$ .

We say that  $X$  has a negative binomial distribution, and we write  $X \sim NB(r, p)$ .

The reason why it is called negative binomial is that the pmf is similar to each term in the Maclaurin's series expansion of the binomial function  $(1-x)^{-r}$  raised to the negative power  $-r$ :

$h(x) = (1-x)^{-r} = \sum_{k \geq 0} \frac{h^{(k)}(0)}{k!} x^k = \sum_{k \geq 0} \binom{r+k-1}{r-1} x^k$

Indeed,  $h^{(1)}(x) = r(1-x)^{-(r+1)}, h^{(2)}(x) = r(r+1)(1-x)^{-(r+2)}$

... so that  $h^{(k)}(0) = r(r+1)\dots(r+k-1) = \frac{(r+k-1)!}{(r-1)!}$  (9)

Consequence: probabilities sum to one:

$$\sum_{x \geq 0} p_x(x) = \sum_{x \geq 0} \binom{x+r-1}{r-1} p^r (1-p)^x = p^r h(1-p) = 1$$

x Moments. Making use of the representation  $X = X_1 + \dots + X_r$ , where  $X_1, \dots, X_r$  are iid  $g(p)$ , we immediately see that

$$\bullet EX = r EX_1 = \frac{r(1-p)}{p}$$

$$\bullet \text{Var } X = r \text{Var } X_1 = \frac{r(1-p)}{p^2}$$

$$X \sim \text{NB}(r, p) \quad p_x(x) = \binom{x+r-1}{r-1} p^r (1-p)^x, \quad x \geq 0$$

$$EX = r(1-p)/p$$

$$\text{Var } X = r(1-p)/p^2$$

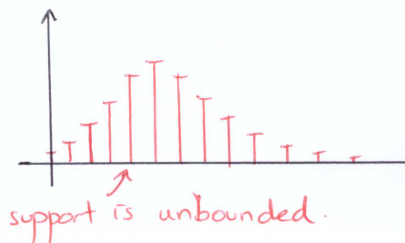
x Shape of the NB distribution.

Consider the ratio  $r(x) = \frac{p_x(x)}{p_x(x-1)} = \dots = \left(\frac{r-1}{x} + 1\right)(1-p)$   
 Compare this ratio with 1.

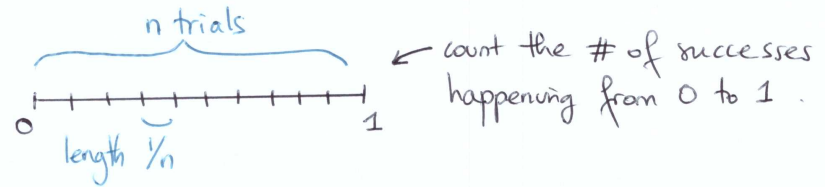
If  $x < \left(\frac{1}{p} - 1\right)(r-1)$ , then  $r(x) > 1$  and the pmf is  $\uparrow$

If  $x > \left(\frac{1}{p} - 1\right)(r-1)$ , then  $r(x) < 1$  and the pmf is  $\downarrow$ .

$\Rightarrow$  The NB distribution contains a single peak "unimodal".



• Poisson distribution.  $\equiv$  Binomial RV in continuous time = counts the number of successes occurring in continuous time ( $\Rightarrow$  still a discrete distribution). (10)



Model: the longer you wait, the more likely you obtain a success. More formally, we assume that the probability of success is proportional to the length of the interval:  $P(\text{success}) = \frac{\lambda}{n}$ ;  $\lambda > 0$ .

Then  $X :=$  number of successes  $\sim \text{Bi}(n, \frac{\lambda}{n})$ .

Its distribution is

$$P(X=x) = \binom{n}{x} \left(\frac{\lambda}{n}\right)^x \left(1 - \frac{\lambda}{n}\right)^{n-x}$$

In the limit, as  $n \rightarrow \infty$ ,

$$\begin{aligned} \lim_{n \rightarrow \infty} P(X=x) &= \lim_{n \rightarrow \infty} \frac{n!}{x!(n-x)!} \frac{\lambda^x}{n^x} \left(1 - \frac{\lambda}{n}\right)^{n-x} \\ &= \lim_{n \rightarrow \infty} \underbrace{\frac{n(n-1)\dots(n-x+1)}{n^x}}_{\downarrow 1} \frac{\lambda^x}{x!} \underbrace{\left(1 - \frac{\lambda}{n}\right)^n}_{\downarrow e^{-\lambda}} \underbrace{\left(1 - \frac{\lambda}{n}\right)^{-x}}_{\downarrow 1} \\ &= \frac{\lambda^x}{x!} e^{-\lambda} \end{aligned}$$

Let  $X$  be such that  $p_x(x) = \frac{\lambda^x}{x!} e^{-\lambda}$ ,  $x \geq 0$ .

We say that  $X$  has a Poisson distribution with parameter  $\lambda$ , and we write  $X \sim P(\lambda)$ .

It is indeed a distribution since

$$\sum_{x \geq 0} p_X(x) = e^{-\lambda} \sum_{x \geq 0} \frac{\lambda^x}{x!} = e^{-\lambda} e^{\lambda} = 1$$

Taylor series expansion of  $e^{\lambda}$ .

x Moments . . .  $EX = \sum_{x \geq 0} x \frac{\lambda^x}{x!} e^{-\lambda}$   
 $= \lambda e^{-\lambda} \sum_{x \geq 1} \frac{\lambda^{x-1}}{(x-1)!} = e^{-\lambda} \lambda e^{\lambda} = \lambda$

⇒ The parameter  $\lambda$  can be interpreted as the mean number of occurrences of success ("arrivals") in a unit interval

•  $EX^2 = \sum_{x \geq 0} x^2 \frac{\lambda^x}{x!} e^{-\lambda} = \dots$  same same  $\dots = \lambda^2 + \lambda$

•  $\text{Var } X = EX^2 - (EX)^2 = \lambda$

⇒ For the Poisson distribution, the mean and variance are governed by a single parameter  $\lambda$ ; which can be suboptimal in applications where the observed values are overdispersed for example.

$X \sim \mathcal{P}(\lambda), \quad p_X(x) = \frac{\lambda^x}{x!} e^{-\lambda}, \quad x \geq 0$   
 $EX = \text{Var } X = \lambda$

x Approximating Binomial by Poisson.

Poisson arises in a limit as  $n \rightarrow \infty$ ; i.e. when  $n$  is large.

Since proba of success  $p = \frac{\lambda}{n}$ ;  $p$  must be small, and

when  $n$  is large &  $p$  is small,  $\binom{n}{x} p^x (1-p)^{n-x} \approx \frac{\lambda^x}{x!} e^{-\lambda}$   
 $Bi(n, p) \approx \mathcal{P}(np)$

x Remark on the Poisson-Binomial approximation.

First, note that if  $X_i \sim \mathcal{P}(\lambda_i), \quad X_1, \dots, X_n$  independent, then  $S := X_1 + \dots + X_n \sim \mathcal{P}(\lambda_1 + \dots + \lambda_n)$

(direct calculation by checking  $n=2$ , then induction, or using moment generating functions).

Lemma: Let  $(X, Y)$  be a bivariate random vector. Then  $\forall B \in \sigma$ -algebra

$$|\mathbb{P}(X \in B) - \mathbb{P}(Y \in B)| \leq \mathbb{P}(X \neq Y)$$

proof =  $\mathbb{P}(X \in B) = \mathbb{P}(X \in B, X=Y) + \mathbb{P}(X \in B, X \neq Y)$   
 $= \mathbb{P}(Y \in B, X=Y) + \mathbb{P}(X \in B, X \neq Y)$   
 $= \mathbb{P}(Y \in B, X=Y) + \mathbb{P}(Y \in B, X \neq Y)$   
 $- \mathbb{P}(Y \in B, X \neq Y) + \mathbb{P}(X \in B, X \neq Y)$

↓  
 $\mathbb{P}(X \in B) - \mathbb{P}(Y \in B) = \mathbb{P}(X \in B, X \neq Y) - \mathbb{P}(Y \in B, X \neq Y)$   
 $\leq \mathbb{P}(X \in B, X \neq Y)$   
 $\leq \mathbb{P}(X \neq Y)$

Since the argument is symmetric, we can take  $|\cdot|$ .

Lemma: Let  $X = X_1 + \dots + X_n$   
 $Y = Y_1 + \dots + Y_n$   
 Then  $\mathbb{P}(X \neq Y) \leq \sum_{i=1}^n \mathbb{P}(X_i \neq Y_i)$

proof = If  $X \neq Y$  then at least one of the  $X_i, Y_i$  must differ, so that  $\{X \neq Y\} \subset \bigcup_{i=1}^n \{X_i \neq Y_i\}$   
 $\Rightarrow \mathbb{P}(X \neq Y) \leq \mathbb{P}\left(\bigcup_{i=1}^n \{X_i \neq Y_i\}\right) \leq \sum_{i=1}^n \mathbb{P}(X_i \neq Y_i)$   
 (subadditivity of  $\mathbb{P}$ )

Next, consider the bivariate  $(X, Y)$  with joint distribution

(13)

$y \backslash x$	0	1	2	3	...	$\mathbb{P}(Y=y)$
0	$e^p(1+p)^{-p}$	0	$p^2 e^p / 2!$	$p^3 e^p / 3!$	...	$1-p$
1	$p - p e^{-p}$	$p e^{-p}$	0	0	...	$p$
$\mathbb{P}(X=x)$	$e^{-p}$	$p e^{-p} / 1!$	$p^2 e^{-p} / 2!$	$p^3 e^{-p} / 3!$	...	1

Note that for  $p \leq 0.8$ , all terms are  $\geq 0$

Then  $X \sim \mathcal{P}(p)$   
 $Y \sim \mathcal{B}(p)$

Note that  $\mathbb{P}(X \neq Y) = 1 - \mathbb{P}(X=Y)$   
 $= 1 - \mathbb{P}(X=Y=0) - \mathbb{P}(X=Y=1)$   
 $= 1 + p - (1+2p)e^{-p}$

since  $(1-p) \leq e^{-p}$   $\hookrightarrow$   $\leq 1 + p - (1+2p)(1-p)$   
 $= 2p^2$

Theorem: let  $Y = Y_1 + \dots + Y_n$ ,  $Y_i \sim \mathcal{B}(p_i)$ ,  $p_i \leq 0.8$   
 $X \sim \mathcal{P}(p_1 + \dots + p_n)$

Then  $\forall B$ ,  $|\mathbb{P}(X \in B) - \mathbb{P}(Y \in B)| \leq 2 \sum_{i=1}^n p_i^2$

proof = let  $(X'_i, Y'_i) \sim$  distributed as  $\uparrow$  distrib of  $Y$  is often called the Poisson-Binomial distribution' in the table above with  $p=p_i$

Then  $X' = X'_1 + \dots + X'_n \sim \mathcal{P}(\sum p_i)$   
 $Y' = Y'_1 + \dots + Y'_n$  has the same distribution as  $Y$ .

$\downarrow$

$|\mathbb{P}(X \in B) - \mathbb{P}(Y \in B)| = |\mathbb{P}(X' \in B) - \mathbb{P}(Y' \in B)|$   
 $\leq \sum_{i=1}^n \mathbb{P}(X'_i \neq Y'_i) \leq 2 \sum_{i=1}^n p_i^2$   $\square$

Note that in the special case that  $p_i = p \forall i$ , (14)  
then  $Y \sim \mathcal{B}(n, p)$   
 $X \sim \mathcal{P}(np)$ , and the bound becomes  $2np^2$

$\Rightarrow \forall B \quad |\mathbb{P}(X \in B) - \mathbb{P}(Y \in B)| \leq 2np^2$

We recover the previous result (made more precise here):  
namely, with  $np = \lambda = \text{constant}$ ,

"distance"  $(\mathcal{B}(n, p), \mathcal{P}(\lambda)) \leq 2 \frac{\lambda^2}{n} \rightarrow 0$   
as  $n \rightarrow \infty$

Total variation distance:  $\sup_B |\mathbb{P}(X \in B) - \mathbb{P}(Y \in B)| \leq 2np^2$   
 $d_{TV}(X, Y)$

Remark: In fact, it is possible to replace the factor 2 in the bound by  $\min(1, \lambda^{-1})$ , so that

$\sup_B |\mathbb{P}(X \in B) - \mathbb{P}(Y \in B)| \leq \min(1, \lambda^{-1}) \sum_{i=1}^n p_i^2$   
 $\uparrow$   $X \sim \mathcal{P}(\sum p_i)$   $\uparrow$   $Y = Y_1 + \dots + Y_n, Y_i \sim \mathcal{B}(p_i)$

One possible way to get this bound is via STEIN'S METHOD = a very general method for quantifying the error made when approximating the distribution of a RV of interest with a reference distribution (here Poisson, but it can be something else: exponential, normal (see later)).

## II- ABSOLUTELY CONTINUOUS DISTRIBUTIONS

(15)

- Uniform distribution = has constant pdf over the set of possible values of  $X$ . For  $a < b$ ,

$$f_X(x) = \begin{cases} 1/(b-a) & \text{if } x \in (a, b) \\ 0 & \text{otherwise} \end{cases}$$

$$F_X(x) = \begin{cases} 0 & \text{if } x < a \\ (x-a)/(b-a) & \text{if } x \in (a, b) \\ 1 & \text{if } x > b \end{cases}$$

We write  $X \sim U(a, b)$ .

Straightforward calculation shows that  $EX = \frac{a+b}{2}$ , and  $\text{Var } X = (b-a)^2/12$ .

- Exponential distribution = continuous version of the geometric distribution. It models the waiting time until the occurrence of the first event.



By time  $t$ , we have completed  $nt$  trials.

$$\mathbb{P}(\text{success}) = \frac{\lambda}{n}$$

Put  $T =$  time to the first success.

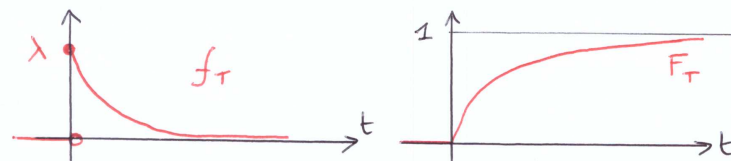
$$\text{Then } \mathbb{P}(\text{no success in } [0, t]) = \left(1 - \frac{\lambda}{n}\right)^{nt} \xrightarrow{n \rightarrow \infty} e^{-\lambda t}$$

$$\text{Thus } \mathbb{P}(T > t) = e^{-\lambda t}$$

$$F_T(t) = \mathbb{P}(T \leq t) = 1 - e^{-\lambda t}, \quad t \geq 0.$$

We say that  $T$  has an exponential distribution with parameter  $\lambda$ , and we write  $T \sim \text{Exp}(\lambda)$ .

The pdf of  $T$  is  $f_T(t) = \lambda e^{-\lambda t}, \quad t \geq 0.$



(16)

x Moments. •  $ET = \int_0^{\infty} \lambda t e^{-\lambda t} dt$  ↗ integration by parts

$$= \lambda \left\{ \left[ -\frac{t}{\lambda} e^{-\lambda t} \right]_0^{\infty} + \frac{1}{\lambda} \int_0^{\infty} e^{-\lambda t} dt \right\}$$

$$= \int_0^{\infty} e^{-\lambda t} dt = \frac{1}{\lambda}$$

The mean waiting time is inversely proportional to  $\lambda$ .

$$\begin{aligned} \bullet ET^2 &= \lambda \int_0^{\infty} t^2 e^{-\lambda t} dt \\ &= \lambda \left\{ \left[ -\frac{t^2}{\lambda} e^{-\lambda t} \right]_0^{\infty} + \frac{2}{\lambda} \int_0^{\infty} t e^{-\lambda t} dt \right\} \\ &= 2 \int_0^{\infty} t e^{-\lambda t} dt = \frac{2}{\lambda^2} \end{aligned}$$

$$\bullet \text{Var } T = ET^2 - (ET)^2 = \frac{1}{\lambda^2}.$$

$$T \sim \text{Exp}(\lambda) \quad f_T(t) = \lambda e^{-\lambda t}, \quad t \geq 0$$

$$ET = \frac{1}{\lambda}$$

$$\text{Var } T = \frac{1}{\lambda^2}$$

x Lack of memory property:

$$\mathbb{P}(T > x+y | T > x) = \frac{\mathbb{P}(T > x+y)}{\mathbb{P}(T > x)} = \frac{e^{-\lambda(x+y)}}{e^{-\lambda x}} = e^{-\lambda y} = \mathbb{P}(T > y)$$

↖ Compare with the geometric distribution



Typical application = lifetime of an electronic component: (17)  
 the probability of lasting an additional  $y$  hours, given that it lasted  $x$  hours, is the same as the probability of lasting  $y$  hours when first put into operation.

↑ Is this a realistic assumption for modeling purposes?

Here, the failure 'rate'  $\lambda$  is constant with time. What if we want to model situations where  $\lambda$  increases with time? decreases with time?

→ We capture this using the HASARD FUNCTION  $\lambda(t)$ , defined as

$$\mathbb{P}(T \in (t, t+dt) \mid T > t) = \frac{\mathbb{P}(T \in (t, t+dt))}{\mathbb{P}(T > t)}$$

↑ probability a success occurs in the infinitesimal interval  $(t, t+dt)$  ... given you made it up to time  $t$ .

$$= \frac{f_T(t) dt}{1 - F_T(t)} =: \lambda(t) dt$$

$$\lambda(t) = \frac{f_T(t)}{1 - F_T(t)} = \text{instantaneous failure rate}$$

↑ Note that for the exponential distribution,  
 $\lambda(t) = \frac{\lambda e^{-\lambda t}}{e^{-\lambda t}} = \lambda = \text{constant}$ .

⇒ Suppose now that the failure rate increases with time in a geometric way  $\lambda(t) = Ct^\alpha$ , with  $\alpha > 0$ . What is the associated distribution of  $T$ ?

Put  $G(t) = 1 - F(t)$  (18)

$$\Rightarrow f(t) = F'(t) = -G'(t)$$

$$\text{and } \lambda(t) = \frac{f(t)}{G(t)} = -\frac{G'(t)}{G(t)} = -\frac{d}{dt} [\log G(t)]$$

$$\text{Thus } G(t) = \exp \left\{ -\int_0^t \lambda(u) du \right\}$$

$$f(t) = \lambda(t) \exp \left\{ -\int_0^t \lambda(u) du \right\}$$

$$\text{Put } \lambda(t) = Ct^\alpha = \frac{k}{\lambda^k} t^{k-1}$$

$$\text{We eventually get } f(t) = \frac{k}{\lambda} \left(\frac{t}{\lambda}\right)^{k-1} \exp\left\{-\left(\frac{t}{\lambda}\right)^k\right\}$$

Weibull distribution  $(\lambda, k > 0)$

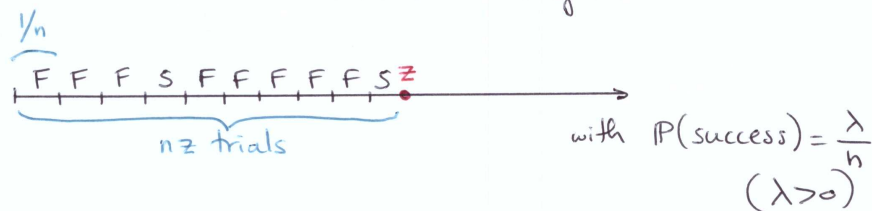
↑ so you can also model a failure rate that decreases with time

Example = Chance of a nuclear power station having a major accident in any year is proportional to its age. Suppose we keep building power stations at a rate one per year, until we have a major accident. Let  $T$  be the time of the first accident. What is (approximately) the distribution of  $T$ ?

↳ one nuclear power station aged  $t$  has an accident rate equal to  $\lambda t$ . After  $t$  years, there are  $t$  power stations  $\Rightarrow$  failure rate is proportional to  $t^2 \Rightarrow$  Weibull distribution with  $k=3$ .

• Gamma distribution. We saw that the exponential distribution is the continuous-time version of the geometric distribution. Similarly, the gamma distribution is the continuous-time version of the negative binomial distribution. (19)

Gamma distribution  $\equiv$  models the waiting time until the  $r$ -th occurrence of an event.



Let  $Z =$  waiting time until the  $r$ -th event  
The distribution of  $Z$  arises as a limit (as  $n \rightarrow \infty$ ) of

$$\begin{aligned} & \mathbb{P}(\text{wait more than } nz \text{ trials for the } r\text{-th success}) \\ &= \mathbb{P}(\text{less than } r \text{ successes in the first } nz \text{ trials}) \\ &= \sum_{k=0}^{r-1} \mathbb{P}(k \text{ successes in the first } nz \text{ trials}) \\ &= \sum_{k=0}^{r-1} \binom{nz}{k} \left(\frac{\lambda}{n}\right)^k \left(1 - \frac{\lambda}{n}\right)^{nz-k} \\ &= \sum_{k=0}^{r-1} \frac{(nz)!}{n^k (nz-k)!} \frac{\lambda^k}{k!} \underbrace{\left(1 - \frac{\lambda}{n}\right)^{nz}}_{e^{-\lambda z}} \underbrace{\left(1 - \frac{\lambda}{n}\right)^{-k}}_1 \\ & \quad \underbrace{\frac{nz(nz-1)\dots(nz-k+1)}{n^k}}_{\rightarrow z^k} \\ & \rightarrow \sum_{k=0}^{r-1} \frac{(\lambda z)^k}{k!} e^{-\lambda z} \Rightarrow \mathbb{P}(Z > z) = \sum_{k=0}^{r-1} \frac{(\lambda z)^k}{k!} e^{-\lambda z} \\ & \quad \cdot F_Z(z) = 1 - \sum_{k=0}^{r-1} \frac{(\lambda z)^k}{k!} e^{-\lambda z} \end{aligned}$$

The pdf of  $Z$  is (20)

$$\begin{aligned} f_Z(z) &= F'_Z(z) = \frac{d}{dz} \left\{ 1 - \sum_{k=0}^{r-1} \frac{(\lambda z)^k}{k!} e^{-\lambda z} \right\} \\ &= \underbrace{\lambda e^{-\lambda z}}_{k=0} + \sum_{k=1}^{r-1} \left\{ -\frac{z^{k-1} \lambda^k}{(k-1)!} e^{-\lambda z} + \frac{z^k \lambda^{k+1}}{k!} e^{-\lambda z} \right\} \\ &= \lambda e^{-\lambda z} + \left\{ \underbrace{-\lambda e^{-\lambda z}}_{\text{cancel each other}} + \underbrace{z \lambda^2 e^{-\lambda z} + \dots + \lambda^r z^{r-1} e^{-\lambda z}}_{\text{same } k=1} \right\} \\ & \quad \text{a telescoping sum! Only the last term remains} \end{aligned}$$

$$f_Z(z) = \frac{\lambda^r z^{r-1}}{(r-1)!} e^{-\lambda z}, \quad z > 0, \lambda > 0$$

Here,  $r$  is an integer. We can generalize the expression of the gamma density by introducing the GAMMA FUNCTION ( $\rightarrow$  justifies also the name of the distribution)  $\rightarrow$

$$\text{let } t > 0 \quad \Gamma(t) = \int_0^{\infty} y^{t-1} e^{-y} dy$$

$$\text{Note that } \Gamma(1) = \int_0^{\infty} e^{-y} dy = 1.$$

Moreover, if  $t > 1$ ,

$$\begin{aligned} \Gamma(t) &= \int_0^{\infty} y^{t-1} e^{-y} dy \quad \text{integration by parts} \\ &= \left[ -y^{t-1} e^{-y} \right]_0^{\infty} + (t-1) \int_0^{\infty} y^{t-2} e^{-y} dy \end{aligned}$$

$$\Gamma(t) = (t-1) \Gamma(t-1)$$

Consequence: if  $n$  is a positive integer,  $\Gamma(n) = (n-1)!$   
 $\Rightarrow$  Gamma function aka the generalized factorial

A random variable  $X$  is said to have a gamma distribution with parameters  $r, \lambda > 0$  if its pdf is

$$f(x) = \frac{\lambda^r x^{r-1}}{\Gamma(r)} e^{-\lambda x}, \quad x > 0$$

We write  $X \sim \mathcal{G}(r, \lambda)$

Note that  $f$  is indeed a pdf since  $f \geq 0$ , and

$$\begin{aligned} \int_0^{\infty} f(x) dx &= \frac{1}{\Gamma(r)} \int_0^{\infty} \lambda^r x^{r-1} e^{-\lambda x} dx \quad \rightarrow y = \lambda x \\ &= \frac{1}{\Gamma(r)} \int_0^{\infty} \lambda^r \frac{y^{r-1}}{\lambda^{r-1}} e^{-y} \frac{dy}{\lambda} \\ &= \frac{1}{\Gamma(r)} \int_0^{\infty} y^{r-1} e^{-y} dy = 1 \\ &= \Gamma(r) \text{ by definition.} \end{aligned}$$

x Moments We directly compute the  $k$ -th ( $k \geq 1$ ) moment of  $X$ :

$$\begin{aligned} EX^k &= \int_0^{\infty} x^k \frac{\lambda^r x^{r-1}}{\Gamma(r)} e^{-\lambda x} dx \quad \rightarrow u = \lambda x \\ &= \int_0^{\infty} \frac{u^k}{\lambda^k} \frac{\lambda^r}{\Gamma(r)} \frac{u^{r-1}}{\lambda^{r-1}} e^{-u} \frac{du}{\lambda} \\ &= \frac{1}{\lambda^k \Gamma(r)} \int_0^{\infty} u^{(k+r)-1} e^{-u} du = \frac{\Gamma(r+k)}{\lambda^k \Gamma(r)} \end{aligned}$$

With  $k=1$ ,  $\Gamma(r+1) = r!$   
 $\Gamma(r) = (r-1)! \Rightarrow EX = \frac{r}{\lambda}$

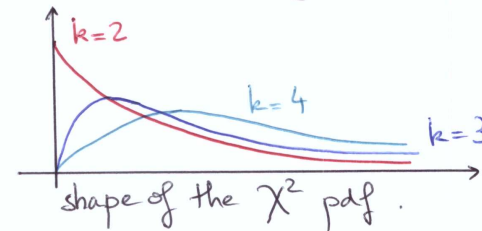
Mean waiting time to the first event is  $1/\lambda$  (= mean of  $\text{Exp}(\lambda)$ )  $\Rightarrow$  Mean waiting time to the  $r$ -th event is  $r/\lambda$ . It all makes sense.

Also,  $EX^2 = \frac{r(r+1)}{\lambda^2} \Rightarrow \text{Var } X = \frac{r}{\lambda^2}$

Remark: The chi-square distribution is a gamma distribution with  $\lambda = 1/2$  and  $r = k/2$ ,  $k$  being a positive integer. Its pdf is:

$$f_X(x) = \frac{1}{2^{k/2} \Gamma(k/2)} x^{k/2-1} e^{-x/2}, \quad x \geq 0$$

We say that  $X$  has a chi-square distribution with  $k$  degrees of freedom, and we write  $X \sim \chi^2(k)$ .  
 We have  $\begin{cases} EX = k \\ \text{Var } X = 2k \end{cases}$



It has numerous applications in statistics for reason that will become clear later.

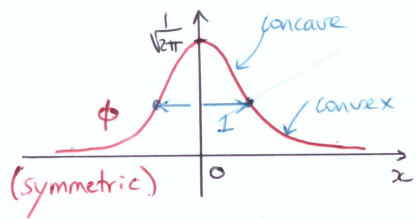
Normal distribution We say that a RV  $X$  has a normal distribution with parameters  $\mu$  and  $\sigma^2$  if its pdf is given by

$$f_X(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}\right\}, \quad \begin{array}{l} x \in \mathbb{R} \\ \mu \in \mathbb{R} \\ \sigma > 0 \end{array}$$

and we write  $X \sim \mathcal{N}(\mu, \sigma^2)$ .

An important special case is the standard normal distribution  $\mathcal{N}(0, 1)$ , whose pdf is

$$\phi(x) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{x^2}{2}\right\}$$



(23)

$$\begin{aligned}\phi(x) &= \frac{1}{\sqrt{2\pi}} e^{-x^2/2} \\ \phi'(x) &= -x \phi(x) \\ \phi''(x) &= -\phi(x) - x \phi'(x) \\ &= (x^2 - 1) \phi(x) \\ &= 0 \text{ at } x = \pm 1\end{aligned}$$

$\Rightarrow \phi$  changes curvature at  $x = \pm 1$

The standard normal cdf is denoted  $\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-u^2/2} du$ .

Remark = let  $Z \sim \mathcal{N}(0, 1)$   
 $X := \sigma Z + \mu$ .

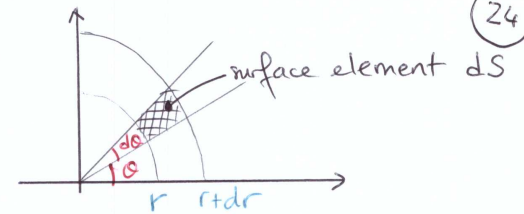
Then  $F_X(x) = \mathbb{P}(X \leq x) = \mathbb{P}\left(Z \leq \frac{x-\mu}{\sigma}\right)$   
 $= \int_{-\infty}^{(x-\mu)/\sigma} \phi(u) du$

$$\begin{aligned}\Rightarrow f_X(x) &= \left(\frac{x-\mu}{\sigma}\right)' \phi\left(\frac{x-\mu}{\sigma}\right) \\ &= \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right\}, x \in \mathbb{R} \\ \Rightarrow X &\sim \mathcal{N}(\mu, \sigma^2).\end{aligned}$$

Summary:  $Z \sim \mathcal{N}(0, 1)$   
 $X := \mu + \sigma Z \sim \mathcal{N}(\mu, \sigma^2)$   
 $\Rightarrow \mu = \text{location parameter}$   
 $\sigma = \text{scale parameter}.$

$\phi$  is indeed a pdf  $\therefore \phi \geq 0$  (obvious)  
 $\cdot \int \phi(x) dx = 1$  (not so obvious)

Use polar coordinates:



$$dS \approx \text{length} \times \text{width} = (dr)(r d\theta)$$

double integral:

$$\Rightarrow \iint f(x, y) dx dy = \iint_{r>0, 0<\theta<2\pi} f(r \cos \theta, r \sin \theta) r dr d\theta$$

cartesian coordinates      polar coordinates

The easiest way to check that  $\int \phi(x) dx = 1$  is to calculate

$$\begin{aligned}\left(\int \phi(x) dx\right)^2 &= \left(\int \phi(x) dx\right) \left(\int \phi(y) dy\right) \\ &= \iint_{\mathbb{R}^2} \phi(x) \phi(y) dx dy \\ &= \frac{1}{2\pi} \iint_{\mathbb{R}^2} \exp\left\{-\frac{x^2+y^2}{2}\right\} dx dy \\ &= \frac{1}{2\pi} \iint e^{-\frac{r^2}{2}} r dr d\theta \quad \begin{array}{l} x = r \cos \theta \\ y = r \sin \theta \\ x^2 + y^2 = r^2 \end{array} \\ &= \frac{1}{2\pi} \int_0^{2\pi} \left\{ \int_0^{\infty} r e^{-\frac{r^2}{2}} dr \right\} d\theta \\ &= \frac{1}{2\pi} \int_0^{2\pi} d\theta \\ &= 1 \quad \Rightarrow \int \phi(x) dx = 1 \quad \blacksquare\end{aligned}$$

x Moments: For  $Z \sim \mathcal{N}(0,1)$ ,

(25)

$$E Z = \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} x e^{-x^2/2} dx = 0$$

↑  
even function

$$E Z^2 = \int_{\mathbb{R}} \frac{1}{\sqrt{2\pi}} x^2 e^{-x^2/2} dx$$

$$= \frac{1}{\sqrt{2\pi}} \left\{ \left[ -x e^{-x^2/2} \right]_{\mathbb{R}} + \int_{\mathbb{R}} e^{-x^2/2} dx \right\} = 1$$

$$\Rightarrow \text{Var } Z = E Z^2 - (E Z)^2 = 1$$

Thus, for  $X: \mu + \sigma Z \sim \mathcal{N}(\mu, \sigma^2)$ ,  $\begin{cases} E X = \mu \\ \text{Var } X = \sigma^2 \end{cases}$   
μ = location & σ = scale.

Lemma: If  $Z \sim \mathcal{N}(0,1)$ , then  $Y := Z^2 \sim \chi^2(1)$

proof = The support of  $Y$  is  $[0, \infty)$ . Let  $y \geq 0$ .

$$F_Y(y) = P(Y \leq y) = P(Z^2 \leq y)$$

$$= P(-\sqrt{y} \leq Z \leq \sqrt{y})$$

$$= 2 \int_0^{\sqrt{y}} \frac{1}{\sqrt{2\pi}} e^{-u^2/2} du$$

↗  $v = u^2$   
↘  $dv = 2u du$

$$= 2 \int_0^y \frac{1}{\sqrt{2\pi}} e^{-\frac{v}{2}} \frac{dv}{2\sqrt{v}}$$

$$= \int_0^y \frac{1}{\sqrt{2\pi v}} e^{-v/2} dv$$

$$\Rightarrow f_Y(y) = F'_Y(y) = \frac{1}{\sqrt{2\pi}} y^{-1/2} e^{-y/2} \Rightarrow Y \sim \chi\left(\frac{1}{2}, 1\right) = \chi^2(1)$$

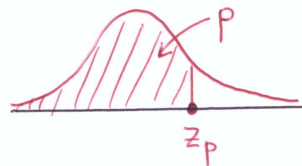
More generally, one can show (e.g. using CHF) that (26)

if  $Z_1, \dots, Z_k$  are  $k$  independent  $\mathcal{N}(0,1)$  RVs, then  $Y := Z_1^2 + \dots + Z_k^2 \sim \chi^2(k)$

↖ hence the name of the parameter of the  $\chi^2$  distribution "degrees of freedom".

x Remark: Percentiles of the  $\mathcal{N}(0,1)$  distribution.

$\forall p \in (0,1)$ , the (100p)-th percentile  $z_p$  (aka the p-QUANTILE) is defined as  $\int_{-\infty}^{z_p} \phi(u) du = p$



↳ Most used values in statistics:

- $p = 0.95$       $z_p \approx 1.645$
- $p = 0.975$     $z_p \approx 1.960$

x Higher-order moments. We derive the  $k$ -th order central moment  $\mu_k := E(X-\mu)^k$  of  $X \sim \mathcal{N}(\mu, \sigma^2)$ .

↖ Note that central moments can be obtained from the  $k$ -th moment  $\mu_k := E X^k$  using the binomial formula:

$$E(X-\mu)^k = E \sum_{i=0}^k \binom{k}{i} X^i (-\mu)^{k-i}$$

$$= \sum_{i=0}^k \binom{k}{i} (-\mu)^{k-i} \mu_i$$

And vice-versa, using  $\mu_k = E[(X-\mu) + \mu]^k$ .

$$\mu_k = \frac{1}{\sigma \sqrt{2\pi}} \int_{\mathbb{R}} (x-\mu)^k \exp\left\{-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}\right\} dx$$

$$v_k = \frac{1}{\sigma\sqrt{2\pi}} \int_{\mathbb{R}} \sigma^k z^k e^{-z^2/2} \sigma dz = \sigma^k \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} z^k e^{-z^2/2} dz \quad (27)$$

$z = \frac{x-\mu}{\sigma}$

↑  
integral = 0 for k odd since the integrand is an odd function.

Take  $k = 2p$ ,  $p \geq 1$

$$v_{2p} = \sigma^{2p} \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} z^{2p} e^{-z^2/2} dz = \sigma^{2p} \sqrt{\frac{2}{\pi}} \int_0^{\infty} z^{2p} e^{-z^2/2} dz$$

$$u = \frac{z^2}{2}$$

$$z = \sqrt{2u}$$

$$dz = \frac{du}{\sqrt{2u}}$$

$$= \sigma^{2p} \sqrt{\frac{2}{\pi}} \int_0^{\infty} (2u)^p e^{-u} \frac{du}{\sqrt{2u}}$$

$$= \sigma^{2p} \frac{2^p}{\sqrt{\pi}} \int_0^{\infty} u^{p-\frac{1}{2}} e^{-u} du$$

$$= \sigma^{2p} \frac{2^p}{\sqrt{2\pi}} \Gamma\left(p + \frac{1}{2}\right)$$

$$= \sigma^{2p} \frac{2^p}{\sqrt{2\pi}} \frac{(2p)!}{4^p p!} \sqrt{\pi}$$

$$= \frac{\sigma^{2p}}{2^p} \frac{(2p)!}{p!}$$

Making use of the expression

$$\Gamma(p) \Gamma\left(p + \frac{1}{2}\right) = \frac{(2p)!}{4^p p!} \sqrt{\pi}$$

Summary:

$$v_k = E(X-\mu)^k = \begin{cases} 0 & \text{if } k = 2p+1 \\ \frac{\sigma^{2p}}{2^p} \frac{(2p)!}{p!} & \text{if } k = 2p \end{cases}$$

Particular cases =

$$v_2 = \sigma^2$$

$$v_4 = 3\sigma^4$$

Higher order moments are given in term of the second moment  $\sigma^2$

Higher order central moments contain additional information about the shape of the distribution, other than its location and spread. (28)

▶ The third central moment is an indicator of SKEWNESS.

$$v_3 = E(X-\mu)^3$$

↪ If the distribution of  $X$  has a long positive tail, then  $(X-\mu)^3$  has large positive values (with high probability) but less often large negative values  $\Rightarrow$  on average  $(X-\mu)^3$  is positive and  $v_3 \geq 0$ . We say that  $X$  has a "heavy right tail".



A heavy right tail distribution (with  $v_3 > 0$ ) looks typically like this.

And similarly, a "heavy left tail" typically implies  $v_3 < 0$ .

If the distribution is symmetric,  $v_3 = 0$ . So  $v_3$  characterize somehow the asymmetry of the distribution.

▶ Peakedness & Tail Thickness = captured by the coefficient of KURTOSIS, defined by  $kurt(X) = \frac{v_4}{\sigma^4} - 3$

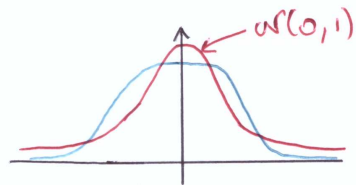
$$= E\left(\frac{X-\mu}{\sigma}\right)^4 - 3$$

Remove location and scale effects;  $\mu = \frac{E(X)}{1}$   
 $\sigma = \sqrt{\text{Var } X}$

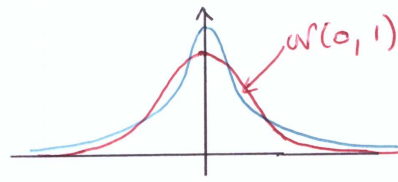
We remove 3 so that  $kurt X = 0$  for  $X \sim \mathcal{N}(\mu, \sigma^2)$  (29)  
 The normal distribution is taken as a reference distribution.

→ If  $kurt X < 0$ , the distribution of  $X$  is likely to be "flatter" and have "lighter tails" than the normal distribution.

→ If  $kurt X > 0$ , the distribution of  $X$  is likely to be "more peaked" and have "heavier tails" than the normal distribution.



$kurt X < 0$



$kurt X > 0$

Example: The pareto distribution is a heavy-tailed distribution: large values are more likely to occur than under a normal distribution. The pdf and cdf of a Pareto RV are

$$f(x) = \frac{\alpha k^\alpha}{x^{\alpha+1}} ; F(x) = 1 - \left(\frac{k}{x}\right)^\alpha, x \geq k > 0$$

The mean value of a Pareto RV might be infinite: provided  $\alpha > 1$ ,

$$E(X^p) = \int_k^{\infty} x^p \frac{\alpha k^\alpha}{x^{\alpha+1}} dx = \alpha k^\alpha \int_k^{\infty} x^{p-\alpha-1} dx = \frac{\alpha k^\alpha}{\alpha-p},$$

which yields

$$\begin{cases} E X = \frac{\alpha k}{\alpha-1}, & \alpha > 1 \\ \text{Var } X = \frac{\alpha k^2}{(\alpha-1)^2(\alpha-2)}, & \alpha > 2 \end{cases}$$

• Multivariate normal distribution. A random vector  $X \in \mathbb{R}^d$  has a multivariate normal distribution if its joint pdf is given by

$$f_X(x) = \frac{1}{(2\pi)^{d/2} \sqrt{|\Sigma|}} \exp\left\{-\frac{1}{2}(x-\mu)^t \Sigma^{-1}(x-\mu)\right\},$$

where  $\mu \in \mathbb{R}^d, x \in \mathbb{R}^d$   
 $\Sigma \in \mathbb{R}^{d \times d}, \Sigma^{-1} = \text{symmetric \& positive definite.}$   
 ( $\forall z \neq 0, z^t \Sigma^{-1} z > 0$ )

x Shape of the multivariate  $\mathcal{N}$  distribution.

The dependence on the variable  $x$  is only through the quadratic form  $\Delta^2 = (x-\mu)^t \Sigma^{-1}(x-\mu)$ .

Remark:  $\Sigma$  can be taken symmetric, without loss of generality. Indeed, suppose that  $\Lambda$  is a  $(d \times d)$  matrix, not necessarily symmetric. Then we can always decompose  $\Lambda$  as  $\Lambda = \Lambda^a + \Lambda^s$ , where  $\Sigma^a$  and  $\Sigma^s$  have entries

$$\Lambda_{ij}^a = \frac{\Lambda_{ij} - \Lambda_{ji}}{2} ; \Lambda_{ij}^s = \frac{\Lambda_{ij} + \Lambda_{ji}}{2}$$

↑  
 antisymmetric  $\Lambda_{ij}^a = -\Lambda_{ji}^a$       symmetric since  $\Lambda_{ij}^s = \Lambda_{ji}^s$

Then

$$\Delta^2 = \sum_{i,j=1}^d (x_i - \mu_i) \Lambda_{ij} (x_j - \mu_j) \quad \mu = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_d \end{pmatrix}$$

↑  
 $\Lambda = \Sigma^{-1}$

$x = \begin{pmatrix} x_1 \\ \vdots \\ x_d \end{pmatrix}$

Substitute  $\Lambda = \Lambda^a + \Lambda^s =$  the term involving  $\Lambda^a$  vanishes  $\Rightarrow$  we can always take  $\Lambda = \Lambda^s$  symmetric.

Consider now the eigenvalues / eigenvectors of the symmetric matrix  $\Sigma$ : (31)

$$\Sigma u_i = \lambda_i u_i$$

Fact: the eigenvalues are real, and the eigenvectors can be chosen orthonormal:  $u_i^t u_j = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{otherwise} \end{cases}$

• Indeed, we first show that the eigenvalues of a symmetric matrix with real coefficients are real:

let  $u_i^+$  = conjugate transpose of  $u_i$ . Then

$$\Sigma u_i = \lambda_i u_i$$

$$u_i^+ \Sigma u_i = \lambda_i u_i^+ u_i \quad \text{--- (*)}$$

Next, consider the conjugate transpose of  $\Sigma u_i = \lambda_i u_i$ :

$$u_i^+ \Sigma^t = \lambda_i^* u_i^+$$

$$\Rightarrow u_i^+ \Sigma^t u_i = \lambda_i^* u_i^+ u_i \quad \text{--- (**)}$$

Subtracting (\*) from (\*\*) yields

$$u_i^+ \Sigma u_i - u_i^+ \Sigma^t u_i = (\lambda_i - \lambda_i^*) u_i^+ u_i$$

$\uparrow \qquad \qquad \uparrow$   
 $\Sigma = \Sigma^t$  since  $\Sigma$  is symmetrical with real coefficients.

$$\Rightarrow \lambda_i = \lambda_i^* \Rightarrow \lambda_i \in \mathbb{R}$$

• We turn our attention to the eigenvectors of  $\Sigma$ .

Since  $\Sigma u_j = \lambda_j u_j$ , we get that

$$\lambda_j u_i^t u_j = u_i^t \Sigma u_j \quad \text{--- } \Sigma \text{ is symm.}$$

$$= u_i^t \Sigma^t u_j$$

$$= (\Sigma u_i)^t u_j$$

$$\Rightarrow \lambda_j u_i^t u_j = \lambda_i u_i^t u_j \quad \text{--- (***)}$$

Consider three cases: (32)

(i)  $\lambda_i, \lambda_j \neq 0$ , and  $\lambda_i \neq \lambda_j$ .

Then necessarily  $u_i^t u_j = 0$  follows from (\*\*\*) i.e.  $u_i \perp u_j$ , and can be taken with unit length.

(ii)  $\lambda_i, \lambda_j \neq 0$ , and  $\lambda_i = \lambda_j = \lambda$ .

Then any linear combination of  $u_i$  and  $u_j$  will be an eigenvector of  $\Sigma$ , with eigenvalue  $\lambda$ , since  $\Sigma (\alpha u_i + \beta u_j) = \alpha \lambda u_i + \beta \lambda u_j = \lambda (\alpha u_i + \beta u_j)$ .

Assuming  $u_i \neq u_j$ , (otherwise

we can construct  $u_\alpha = \alpha u_i + \beta u_j$   
 $u_\beta = \gamma u_i + \delta u_j$ ,

such that  $u_\alpha \perp u_\beta$ , of unit length. It follows from the previous remark that  $u_\alpha, u_\beta$  are eigenvectors of  $\Sigma$ .

Since  $u_i, u_j$  are  $\perp$  to  $u_k$ , for  $k \neq i, j$  (from case (i)), so are  $u_\alpha$  and  $u_\beta$ .

(iii)  $\lambda_i = 0, \lambda_j \neq 0$

Then  $u_i$  belongs to the null space of  $\Sigma$ .

$$\Sigma u_i = 0 \Rightarrow u_i^t \Sigma^t = 0 \quad \text{--- } \Sigma \text{ symmetric}$$

$$\Rightarrow u_i^t \Sigma = 0$$

$$\Rightarrow u_i^t \Sigma u_j = 0$$

$$\Rightarrow \lambda_j u_i^t u_j = 0$$

Since  $\lambda_j \neq 0$ , we see that  $u_i \perp u_j$ ; and we can choose  $u_i$  of unit length.



The symmetric matrix  $\Sigma$  can be expressed as an expansion in terms of its eigenvectors & eigenvalues: (33)

$$\Sigma = \sum_{i=1}^d \lambda_i u_i u_i^t = U \Lambda U^t, \text{ where } \begin{cases} \Lambda = \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_d \end{pmatrix} \\ U = \begin{pmatrix} | & & | \\ u_1 & \dots & u_d \\ | & & | \end{pmatrix} \\ U^t U = I_d \end{cases}$$

The SPECTRAL DECOMPOSITION of  $\Sigma$ .

Note that if all the  $\lambda_i \neq 0$ , then  $\Sigma$  is invertible, and  $\Sigma^{-1} = \sum_{i=1}^d \frac{1}{\lambda_i} u_i u_i^t$ .

Indeed, put  $M := U \Lambda U^t$

Then  $U^t M U = U^t (U \Lambda U^t) U = \Lambda$

&  $U^t (\Sigma U) = U^t (U \Lambda) = U^t (U \Lambda) = \Lambda$

Since  $\Lambda$  is diagonal and  $U$  square

Thus  $U^t M U = U^t \Sigma U \Rightarrow M = \Sigma$   
since  $U U^t = I$ ;  $U$  orthogonal

Moreover,  $(U \Lambda U^t)^{-1} = (U^t)^{-1} \Lambda^{-1} U^{-1}$   
 $= U \Lambda^{-1} U^t = \sum_{i=1}^d \frac{1}{\lambda_i} u_i u_i^t$

Note that indeed both  $\Sigma$  and  $\Sigma^{-1}$  are symmetric. Moreover, both are positive definite ( $\forall z \neq 0, z^t \Sigma z \neq 0$  and  $z^t \Sigma^{-1} z \neq 0$ ), with strictly positive eigenvalues.

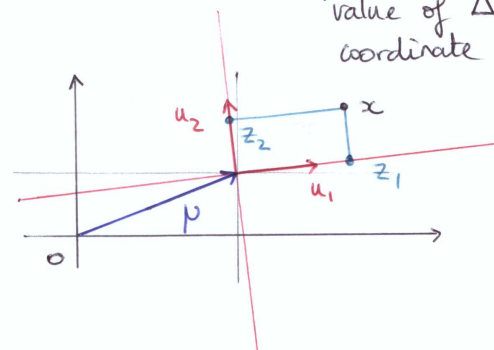
Indeed, for some  $z \neq 0$ ,  $z^t \Sigma z = (z^t U) \Lambda (U z) = v^t \Lambda v$ . Now, assuming the diagonal terms are non-zero, we see that  $v^t \Lambda v > 0$  provided all  $\lambda_i$  are  $> 0$ .

Moreover, the determinant of a positive definite matrix  $\Sigma$  is non zero, and given by the product of its eigenvalues:  $|\Sigma| = \prod_{i=1}^d \lambda_i$ . (34)

$$\begin{aligned} \Rightarrow \Delta^2 &= (x - \rho)^t \Sigma^{-1} (x - \rho) \\ &= \sum_{i=1}^d \frac{(x - \rho)^t u_i u_i^t (x - \rho)}{\lambda_i} \\ &= \sum_{i=1}^d \frac{z_i^2}{\lambda_i} \end{aligned}$$

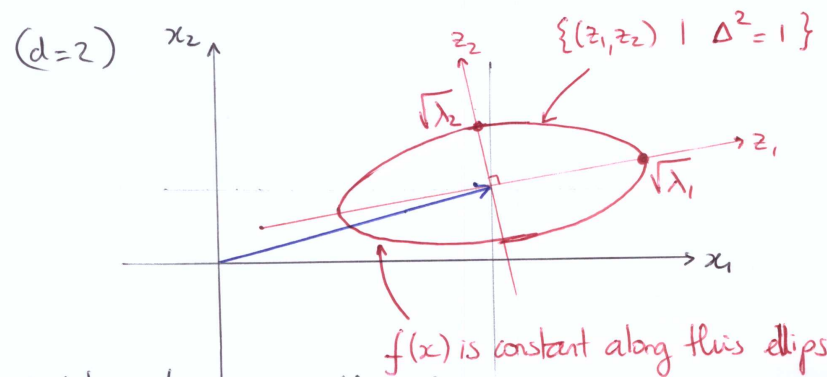
Put  $z_i := u_i^t (x - \rho)$   
 $[ \in \mathbb{R} ]$

value of  $\Delta$  in the new coordinate system.



$\Rightarrow$  the multivariate normal density is constant on surfaces where  $\Delta^2$  is constant. If all  $\lambda_i > 0$ , this surface is an ellipse:

$$(d=2) \quad \frac{z_1^2}{\lambda_1} + \frac{z_2^2}{\lambda_2} = C.$$



$f(x)$  is constant along this ellipse.

How stretched/compressed the ellipse is depends on the values of  $\lambda_i > 0$ .

In the new coordinate system,

$$f(z) = \prod_{i=1}^d \frac{1}{(2\pi\lambda_i)^{1/2}} \exp\left\{-\frac{1}{2} \frac{z_i^2}{\lambda_i}\right\}$$

= product of d independent univariate normal distrib  $\mathcal{N}(0, \lambda_i)$ . [and so  $\int f(z) dz = 1$ , as required]

x Moments.

$$E X = \frac{1}{(2\pi)^{d/2} \sqrt{|\Sigma|}} \int_{\mathbb{R}^d} x \exp\left\{-\frac{1}{2}(x-\mu)^t \Sigma^{-1}(x-\mu)\right\} dx$$

$$= \frac{1}{(2\pi)^{d/2} \sqrt{|\Sigma|}} \int_{\mathbb{R}^d} (z+\mu) \exp\left\{-\frac{1}{2} z^t \Sigma^{-1} z\right\} dz$$

=  $\mu$ , since  $\exp\{\dots\}$  is an even function of  $z$ , and the integral is taken over  $\mathbb{R}^d$ .

$$E X X^t = \frac{1}{(2\pi)^{d/2} \sqrt{|\Sigma|}} \int_{\mathbb{R}^d} x x^t \exp\left\{-\frac{1}{2}(x-\mu)^t \Sigma^{-1}(x-\mu)\right\} dx$$

$$\int_{\mathbb{R}^d} (z+\mu)(z+\mu)^t \exp\left\{-\frac{1}{2} z^t \Sigma^{-1} z\right\} dz$$

↳ The terms  $z\mu^t$  and  $\mu^t z$  vanish by symmetry

↳ The term  $\mu\mu^t$  = constant and can be taken out of the integral.

↳ It remains to evaluate the term with  $z^t z$

Recall:  $u_1, \dots, u_d$  = orthonormal basis.

$$\text{Put } z = \sum_{i=1}^d \langle z, u_i \rangle u_i = \sum_{i=1}^d \bar{z}_i u_i$$

$\uparrow$   
 $\bar{z}_i := u_i^t z$

Thus  $z z^t = \sum_{i,j} \bar{z}_i \bar{z}_j u_i^t u_j$ , and

$$\int z z^t \exp\left\{-\frac{1}{2} z^t \Sigma^{-1} z\right\} dz = \sum_{i,j} u_i u_j^t \int \bar{z}_i \bar{z}_j \exp\left\{-\sum_{k=1}^d \frac{\bar{z}_k^2}{2\lambda_k}\right\} d\bar{z}$$

Jacobian of the transformation is equal to 1

odd function, unless  $i=j$ .  
⇒ integral vanishes unless  $i=j$

$$= \sum_{i=1}^d u_i u_i^t \int \bar{z}_i^2 \exp\left\{-\sum_{k=1}^d \frac{\bar{z}_k^2}{2\lambda_k}\right\} d\bar{z}$$

even function.

$$E X X^t = \sum_{i=1}^d u_i u_i^t \int \prod_{k=1}^d \left[ \frac{\bar{z}_i^2}{(2\pi\lambda_k)^{1/2}} \exp\left\{-\frac{\bar{z}_k^2}{2\lambda_k}\right\} \right] d\bar{z} + \mu\mu^t$$

$$\left( \int \frac{1}{\sqrt{2\pi\lambda_1}} \exp\left\{-\frac{\bar{z}_1^2}{2\lambda_1}\right\} d\bar{z}_1 \right) \times \dots \times \left( \int \frac{\bar{z}_i^2}{\sqrt{2\pi\lambda_i}} \exp\left\{-\frac{\bar{z}_i^2}{2\lambda_i}\right\} d\bar{z}_i \right)$$

$$1 \times \dots \times \left( \int \frac{1}{\sqrt{2\pi\lambda_d}} \exp\left\{-\frac{\bar{z}_d^2}{2\lambda_d}\right\} d\bar{z}_d \right)$$

=  $\lambda_i$  (= variance of a  $\mathcal{N}(0, \lambda_i)$ )

$$E X X^t = \sum_{i=1}^d \lambda_i u_i u_i^t + \mu\mu^t = \Sigma + \mu\mu^t$$

$$\Rightarrow \Sigma = E X X^t - \mu\mu^t = E (X-\mu)(X-\mu)^t = \text{COVARIANCE MATRIX of } X$$

Summary = A multivariate normal RV  $X$  has density  $f(x) = \frac{1}{(2\pi)^{d/2} \sqrt{|\Sigma|}} \exp\left\{-\frac{1}{2}(x-\mu)\Sigma^{-1}(x-\mu)^t\right\}$  (37)

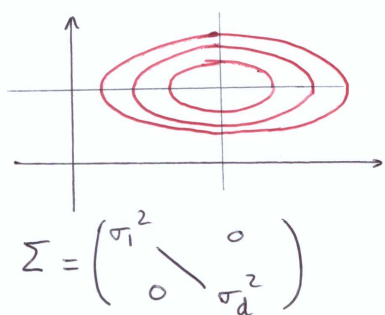
where  $EX = \mu \in \mathbb{R}^d$

$\text{Cov } X = \Sigma \in \mathbb{R}^{d \times d}$  is symm, pos. definite.

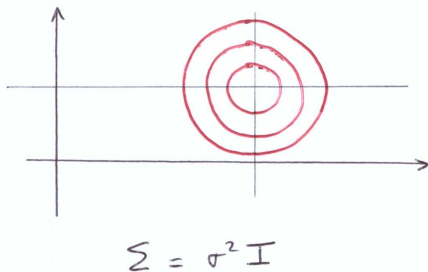
We write  $X \sim \mathcal{N}(\mu, \Sigma)$ .

↖ A general covariance matrix has  $\frac{d(d+1)}{2}$  distinct terms  
 $\Rightarrow$  Total of  $\frac{d(d+1)}{2} + d = \frac{d(d+3)}{2} = O(d^2)$  terms to describe a  $\mathcal{N}(\mu, \Sigma)$  distribution.

↳ Special cases =



The axes of the contour lines are aligned with the coordinate axes



The contours are concentric circles.

Remark: If  $\forall i \neq j \text{ Cov}(X_i, X_j) = 0$ , then  $\Sigma^{-1}$  is diagonal and the multivariate normal density factorizes as a product of  $d$  marginal distributions  $\Rightarrow$  Uncorrelation implies independence for normal RVs.

• Student's t distribution. (38)

Let  $Z \sim \mathcal{N}(0, 1)$  and  $U \sim \chi^2(k)$  independent RVs.

Put  $T = \frac{Z}{\sqrt{U/k}}$ .

$T$  take values in  $\mathbb{R}$ .

↳ We derive the distribution of  $T$ . The first step is to derive the expression of the joint distribution of  $(T, U)$ :

$$f_{(T,U)}(t,u) = \underbrace{f_{T|U}(t|u)}_{\downarrow} f_U(u)$$

On the event  $\{U = u\}$ , one has  $T = \frac{Z}{\sqrt{u/k}}$  and

$$\mathbb{P}(T \leq t \mid U = u) = \mathbb{P}\left(Z \leq t \sqrt{\frac{u}{k}} \mid U = u\right)$$

$$= \int_{-\infty}^{t \sqrt{u/k}} \phi(v) dv$$

$$\Rightarrow f_{T|U}(t|u) = \sqrt{\frac{u}{k}} \phi\left(t \sqrt{\frac{u}{k}}\right)$$

$$= \sqrt{\frac{u}{k}} \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2} \frac{ut^2}{k}\right\},$$

$\begin{matrix} u > 0 \\ t \in \mathbb{R} \end{matrix}$

It follows that

$$f_{(T,U)}(t,u) = \sqrt{\frac{u}{k}} \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2} \frac{ut^2}{k}\right\} \frac{u^{\frac{k}{2}-1}}{2^{\frac{k}{2}} \Gamma(\frac{k}{2})} e^{-u/2}$$

$$= \frac{2^{-\frac{(k+1)}{2}}}{\sqrt{\pi k} \Gamma(\frac{k}{2})} u^{\frac{k+1}{2}-1} \exp\left\{-\frac{u}{2} \left(1 + \frac{t^2}{k}\right)\right\} =: C$$

Now, integrate out the variable  $u$  to get the distribution of  $T$ :

$$\begin{aligned}
 f_T(t) &= \int_0^{+\infty} f_{(T,u)}(t, u) du \\
 &= C \int_0^{+\infty} u^{\frac{k+1}{2}-1} \exp\left\{-\frac{u}{2}\left(1+\frac{t^2}{k}\right)\right\} du \\
 &\quad \text{let } w = \frac{u}{2}\left(1+\frac{t^2}{k}\right) \\
 &\quad du = 2\left(1+\frac{t^2}{k}\right)^{-1} dw \\
 &= C \int_0^{+\infty} \left[2w\left(1+\frac{t^2}{k}\right)^{-1}\right]^{\frac{k+1}{2}-1} e^{-w} \frac{2}{1+\frac{t^2}{k}} dw \\
 &= C 2^{\frac{k+1}{2}} \left(1+\frac{t^2}{k}\right)^{-\frac{k+1}{2}} \underbrace{\int_0^{+\infty} w^{\frac{k+1}{2}-1} e^{-w} dw}_{= \Gamma\left(\frac{k+1}{2}\right)}
 \end{aligned}$$

After simplifications,

$$f_T(t) = \frac{\Gamma\left(\frac{k+1}{2}\right)}{\sqrt{\pi k} \Gamma\left(\frac{k}{2}\right)} \left(1 + \frac{t^2}{k}\right)^{-\frac{k+1}{2}}, \quad t \in \mathbb{R}$$

We say that  $T$  has a Student's  $t$  distribution with  $k$  degrees of freedom, and we write  $T \sim t(k)$

Numerous applications in Statistics

(what happens as  $k \rightarrow \infty$ ? why?)

