

SF1901: Sannolikhetslära och statistik

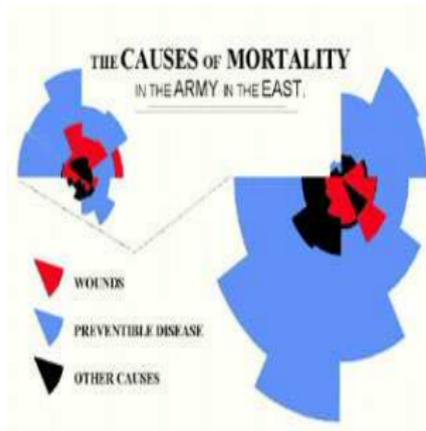
Föreläsning 1.

Jan Grandell & Timo Koski

19.01.2016



Många tänker på tabeller ¹ när de hör ordet "statistik". Här avses dock med *statistik* läran om hur man från observationer eller analyser under osäkerhet drar slutsatser och beskriver dessa slutsatser på ett korrekt sätt.



¹Florence Nightingale (i bilden till höger) var även en framstående statistiker, tabellen till vänster har uppfunnits av henne

- aritmetiskt medelvärde, standardavvikelse, relativ frekvens
- utfall, utfallsrum, händelse, omöjlig (tom) händelse
- union av två händelser, snitt av två händelser
- komplementhändelse, De Morgans regler
- sannolikhet (snl), räknereglerna för sannolikhet
- Ex. på snl: $\frac{\text{antalet gynnsamma fall}}{\text{totalantalet fall}}$, kombinatorik

Tabell: ogrupperade data

Man undersökte 35 tändsticksaskar och noterade för varje ask hur många tändstickor den innehöll. Följande värden ($x_i, i = 1, \dots, 35$) erhöles:

51	52	49	51	52	51	53
52	48	52	50	53	49	50
51	53	51	52	50	51	53
53	55	50	49	53	50	51
51	52	48	53	50	49	51

Detta är ett exempel på *ogrupperade* data och **VARIATION**.

Hur många tändstickor tillverkas i Sverige ?

Varje dag tillverkar Swedish Match cirka 5 miljoner tändsticksaskar, vilket motsvarar omkring 250 miljoner tändstickor.

Ställer man askarna på varandra, räcker en minuts produktion till en pelare högre än Eiffeltornet i Paris. Läger man askarna på rad på E4:an, räcker en dags produktion från Jönköping till Stockholm. Stickornas sammanlagda längd skulle räcka till Australien.



En stor ogrupperad datamängd är svår att överskåda. Resultatet i exemplet ovan kan sammanfattas i en *frekvenstabell* av *grupperade* data.

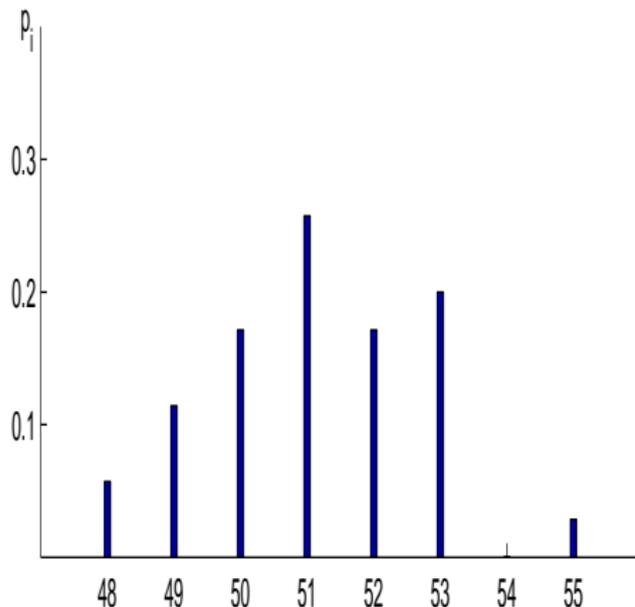
- *absoluta frekvenserna* f_i för de olika förekommande värdena
- de *relativa frekvenserna* $p_i = f_i/n$. (Här avser $i = 1$ lägsta klassen, $i = 2$ nästa klass o.s.v..) $n =$ antalet data (=35).

Table : Frekvenstabell för antal tändstickor i tändsticksaskar.

Klass	Absolut frekvens	Relativ frekvens (%)
i	f_i	$100 p_i$
48	2	5.7
49	4	11.4
50	6	17.1
51	9	25.7
52	6	17.1
53	7	20.0
54	0	0.0
55	1	2.9
S:a	35	100.0

Stolpdiagram (Variation)

Större åskådlighet får man genom ett *stolpdiagram* med de relativa frekvenserna p_i inritade.



Lägesmått (Variation)

Låt allmänt x_1, \dots, x_n vara de data som skall bearbetas. Som lägesmått används ofta (aritmeriska) *medelvärde*

$$\bar{x} = \frac{1}{n} \sum_{j=1}^n x_j = \frac{x_1 + x_2 + \dots + x_n}{n}.$$

I exemplet blir $\sum_{j=1}^{35} x_j = 1789$ och $\bar{x} = 1789/35 = 51.1143$.



Spridningsmått (Variation)

Som spridningsmått används ofta (*stickprovs*)*variansen*

$$s^2 = \frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{x})^2$$

eller (*stickprovets*) *standardavvikelse*

$$s = \sqrt{\frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{x})^2},$$

Kvadratrotstrutdragningen motiveras av att samma enhet då erhålles som för de givna värdena.

I kursens formelsamling hittar man

$$\begin{aligned}s^2 &= \frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{x})^2 \\ &= \frac{1}{n-1} \left(\sum_{j=1}^n x_j^2 - n\bar{x}^2 \right)\end{aligned}$$

I exemplet blir $\sum_{j=1}^{35} x_j^2 = 91533$ och $\bar{x}^2 = 2612.7$ och standardavvikelsen blir $s = \sqrt{\frac{1}{34} (91533 - 35 \cdot 2612.67)} \approx 1.62$.

Med *kovariansen* mellan x - och y -värdena i en datamängd $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ menas

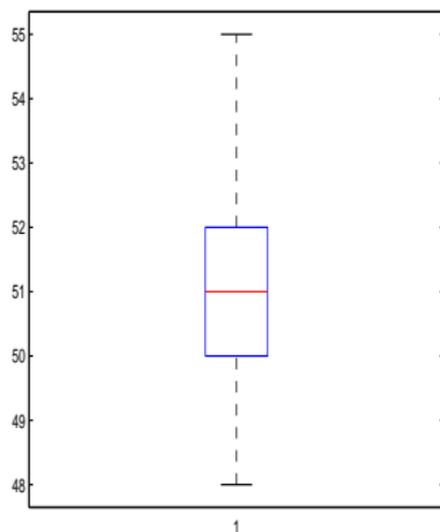
$$c_{xy} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

och med *korrelationskoefficienten* menas

$$r = \frac{c_{xy}}{s_x s_y},$$

där s_x och s_y är stickprovsstandardavvikelserna för x - respektive y -data.

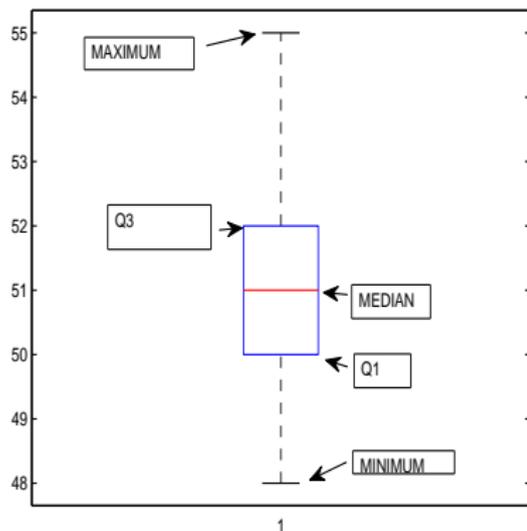
BOXPLOT för tändsticksdata



51	52	49	51	52	51	53
52	48	52	50	53	49	50
51	53	51	52	50	51	53
53	55	50	49	53	50	51
51	52	48	53	50	49	51

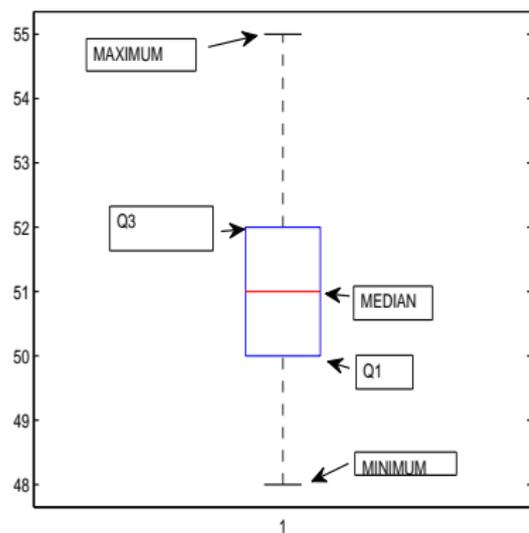
BOXPLOT a.k.a box-and-whisker diagram

A **boxplot** is a graph of a data set that consists of a line from the minimum value to the maximum value and a box with lines drawn at the first quartile Q_1 , the median and the third quartile Q_3



BOXPLOT

första kvartilen $Q_1 = 25\%$ av värdena ligger under, medianen = mittpunkten, 50% ligger under, the tredje kvartilen $Q_3 = 75\%$ ligger under



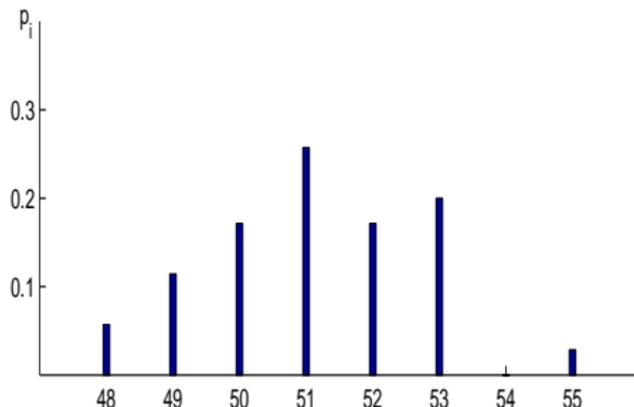
statistics ... uses powerful computers and sophisticated mathematical models to hunt for meaningful patterns and insights in vast troves of data. The applications are as diverse as improving Internet search and online advertising, culling gene sequencing information for cancer research and analyzing sensor and location data to optimize the handling of food shipments.

mathematical model = sannolikhetsmodell (slumpmodell)



Sannolikheten för händelsen " $50 \leq$ antalet tändstickor i en ask ≤ 53 "

$$p_{50} + p_{51} + p_{52} + p_{53} = 0.171 + 0.257 + 0.171 + 0.20 \approx 0.80$$



Vi skriver **sannolikheten för händelsen** " $50 \leq$ antalet tändstickor i en ask ≤ 53 "

$$P(50 \leq \text{antalet tändstickor i en ask} \leq 53) = 0.80$$

Detta är en sannolikhet som baserar sig på observationer. Om inget ändras i produktionsprocessen förväntar vi oss att om vi får en ny tändsticksask och checkar antalet tändstickor i den, så har vi 80% chans att händelsen " $50 \leq$ antalet tändstickor i en ask ≤ 53 " inträffar.

Sannolikheten

$$Pr(50 \leq \text{antalet tändstickor i en ask} \leq 53) = 0.80$$

baserade sig på 35 tändsticksaskar. De relativa frekvenserna kommer att ändras vi kollar nya tändsticksaskar, men de kommer att stabilisera sig.

Vi postulerar att det finns ett tal som kallas sannolikheten för en händelse. Därmed menas helt enkelt att *man tilldelar varje händelse ett visst tal*. Om händelsen är A , betecknas talet med $P(A)$, *sannolikheten för A* . (P som i Probability.)

Beträffande talet $P(A)$ gäller allmänt, att man söker välja det så att den relativa frekvensen vid ett någorlunda stort antal försök kommer i närheten av $P(A)$. Om vi säger att $P(A) = 0.80$ kan vi ge detta uttalande den påtagliga men samtidigt vaga *frekvenstolkningen*: Vid ett stort antal försök blir den relativa frekvensen av händelsen A nog ungefär lika med 0.80.

Vi betraktar nu slumpförsök medelst allmänna beteckningar

- Varje möjligt resultat ω av ett slumpförsök kallas ett *utfall*, eller ett elementärt utfall.

Vi betraktar nu slumpförsök medelst allmänna beteckningar

- Varje möjligt resultat ω av ett slumpförsök kallas ett *utfall*, eller ett elementärt utfall.
- Mängden av alla *utfall*, eller resultat, kallar vi utfallsrummet och betecknar det med Ω .

Vi betraktar nu slumpförsök medelst allmänna beteckningar

- Varje möjligt resultat ω av ett slumpförsök kallas ett *utfall*, eller ett elementärt utfall.
- Mängden av alla *utfall*, eller resultat, kallar vi utfallsrummet och betecknar det med Ω .
- En händelse A är en mängd av utfall, dvs en delmängd av Ω , $A \subset \Omega$.

Ofta kan valet av utfallsrum bero på situationen eller den fråga vi vill studera.

Exempel: $\Omega =$ de fem miljoner tändsticksaskarna producerade under en given dag.

$\Omega =$ de 35 miljoner tändsticksaskarna producerade under en sjudagars period.

Exempel: tärningskast

Slutförsök: kast av en tärning. Utfall: $\omega =$ antalet ögon

$$\Omega = \{\text{etta, tvåa, trea, fyra, femma, sexa}\}$$

Exempel på en händelse $A =$ udda antal ögon $= \{\text{etta, trea, femma}\}$. Vi säger att *händelsen A inträffar*, om vi kastar en tärning (=genomför ett slutförsök) och får etta eller trea eller femma. (Slut på exemplet.)



Definition

Mängden av alla utfall, eller resultat, kallar vi utfallsrummet och betecknar det med Ω .

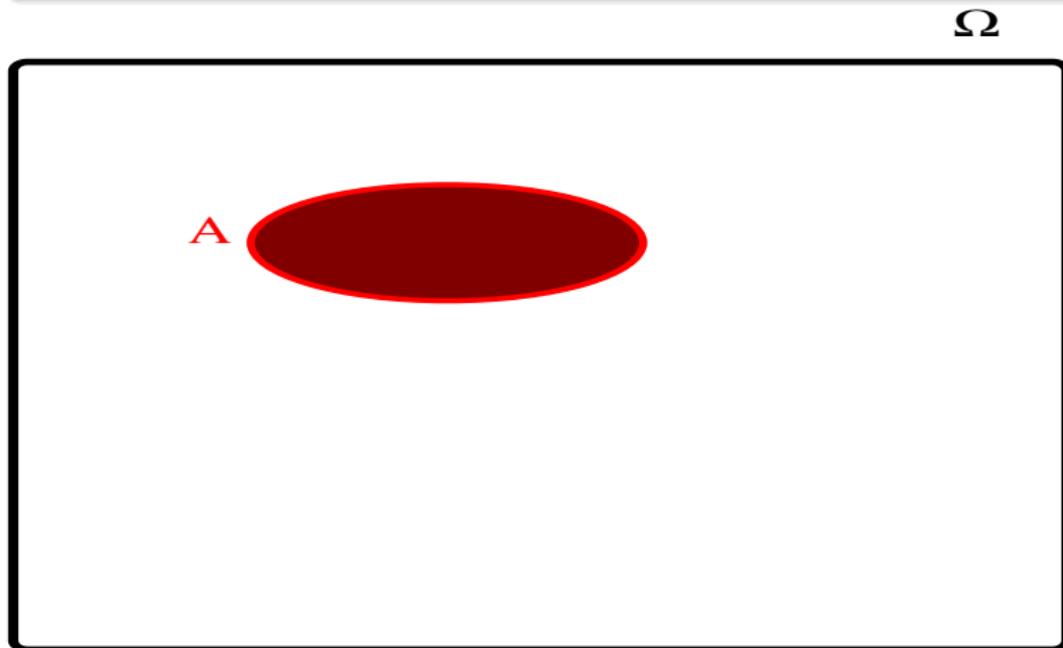
Ω

ω

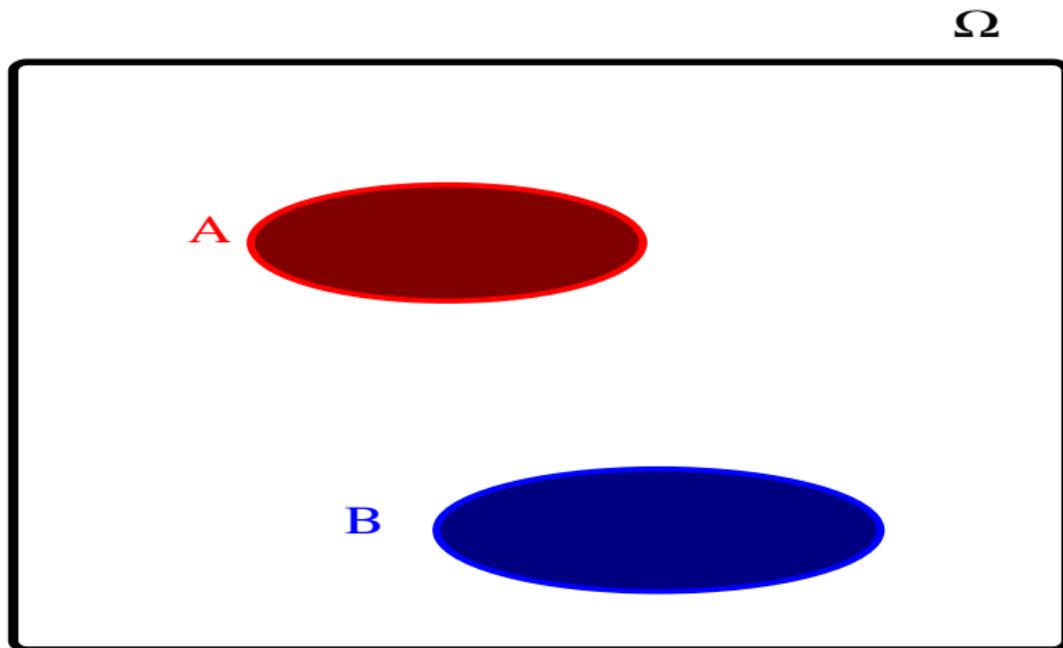


Definition

En händelse A är en mängd av utfall, dvs en delmängd av Ω , $A \subset \Omega$.



Venn diagram; två händelser



Händelser $A \cap B$

Låt oss nu anta att vi är intresserade av två händelser A och B definierade på samma försök. Här är några exempel på vad som kan inträffa, och hur vi matematiskt kan uttrycka detta:

" A inträffar", A

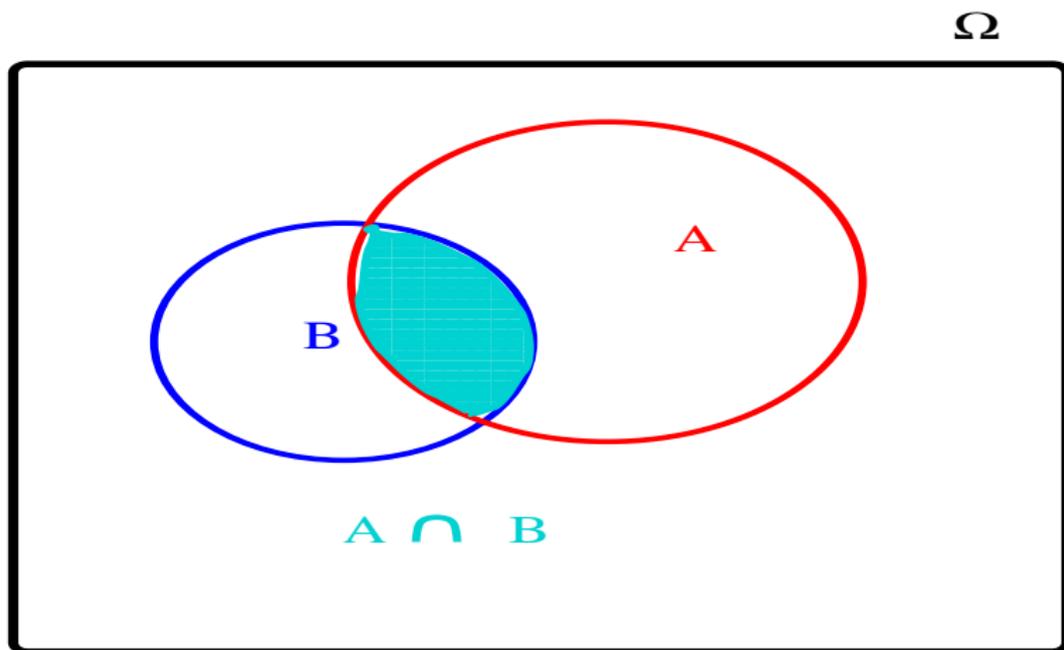
" A och B inträffar" eller " A snitt B inträffar", $A \cap B$

Exempel

$$\Omega = \{ \text{etta, tvåa, trea, fyra, femma, sexa} \}$$

$A = \text{udda antal ögon} = \{ \text{etta, trea, femma} \}$. $B = \{ \text{femma, sexa} \}$,
 $A \cap B = \{ \text{femma} \}$.

Venndiagram $A \cap B$



" A eller B inträffar" eller " A union B inträffar", $A \cup B$

Obs! $A \cup B$ betyder att minst en av A eller B inträffar, så $A \cap B$ kan mycket väl inträffa. I matematik betyder "eller" och/eller!

$A \cup B$ betyder att minst en av A eller B inträffar

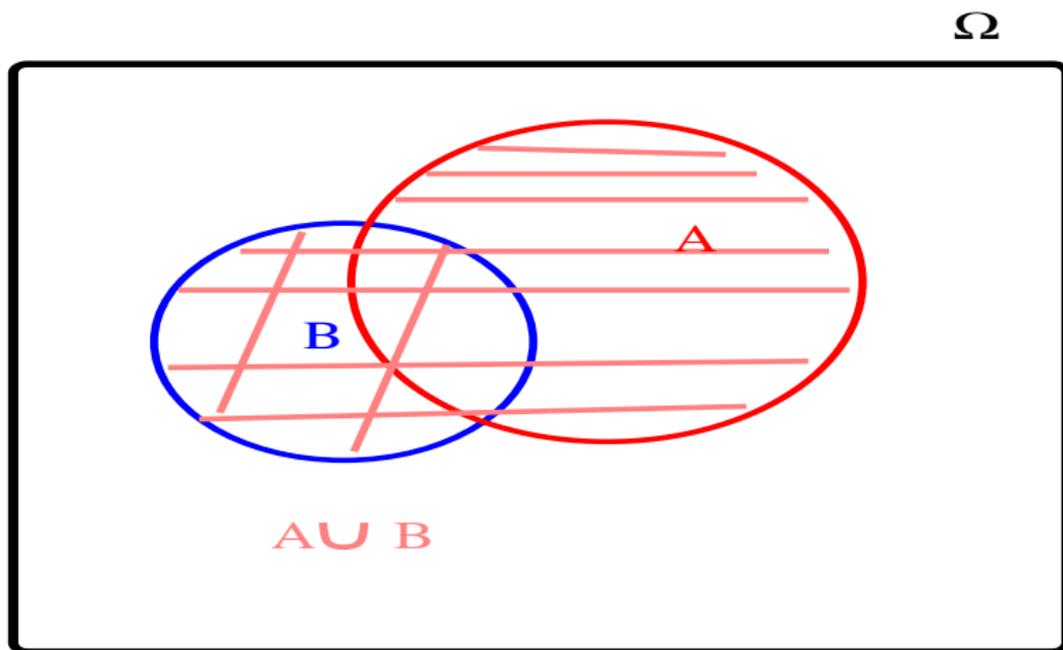
Exempel

$$\Omega = \{ \text{etta, tvåa, trea, fyra, femma, sexa} \}$$

$A = \text{udda antal ögon} = \{ \text{etta, trea, femma} \}$. $B = \{ \text{femma, sexa} \}$,

$A \cup B = \{ \text{etta, trea, femma, sexa} \}$.

Venn diagram $A \cup B$



" A inträffar inte", A^* .

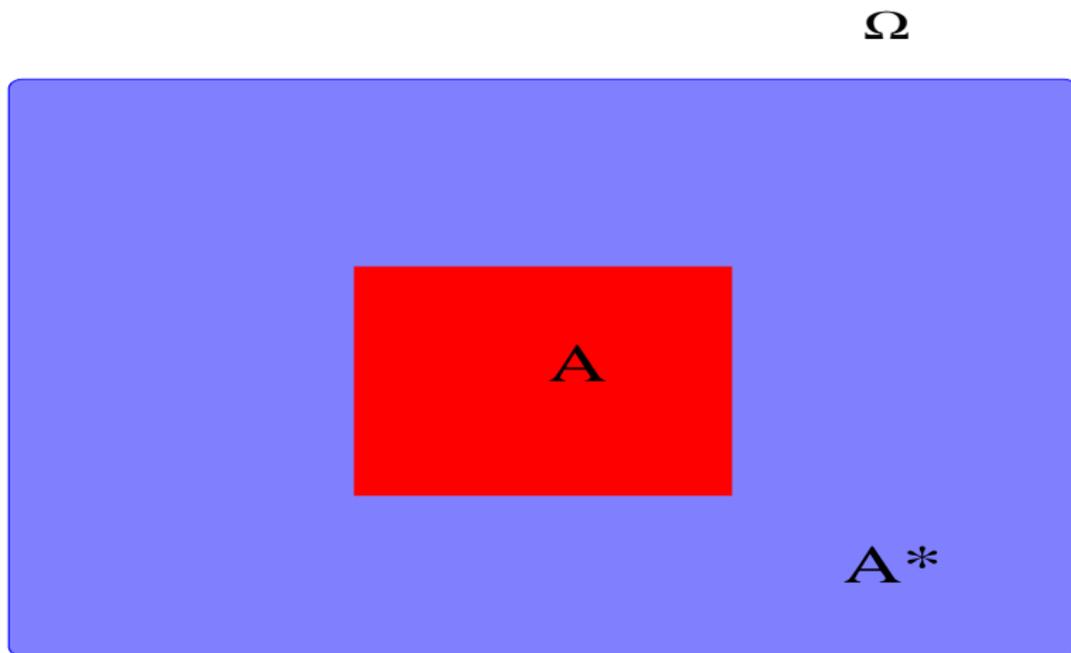
Exempel

$$\Omega = \{ \text{etta, tvåa, trea, fyra, femma, sexa} \}$$

$A = \text{udda antal ögon} = \{ \text{etta, trea, femma} \}.$

$A^* = \{ \text{tvåa, fyra, sexa} \} = \text{jämnt antal ögon}.$

Venn diagram A^*



"tomma mängden" \emptyset

Om A och B utesluter varandra, dvs. omöjligt kan inträffa samtidigt, så säger vi att A och B är *disjunkta* eller oförenliga, dvs.

$$A \cap B = \emptyset$$

där \emptyset är "tomma mängden" eller "den omöjliga händelsen".

$$\Omega^* = \emptyset$$



"tomma mängden" \emptyset

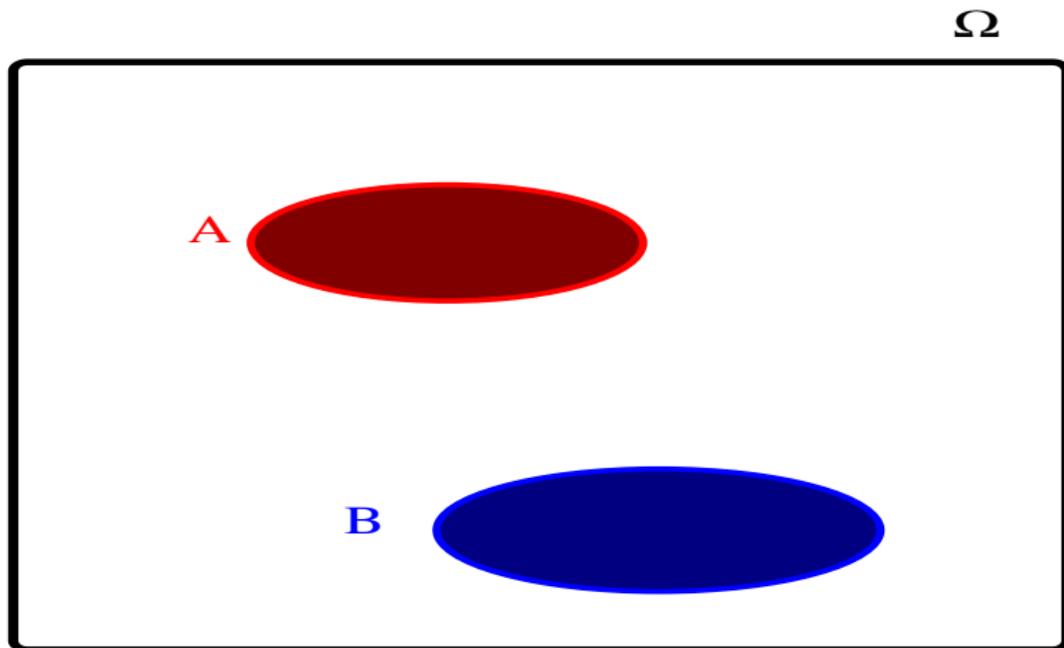
A och B utesluter varandra, dvs. omöjligt kan inträffa samtidigt.

Exempel

$$\Omega = \{ \text{etta, tvåa, trea, fyra, femma, sexa} \}$$

$$A = \{ \text{etta, trea, femma} \}. \quad B = \{ \text{fyra, sexa} \}, \quad A \cap B = \emptyset.$$

Venn diagram; $A \cap B = \emptyset$



Sats

$$A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$$

$$A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$$

Sats

$$(A \cap B)^* = A^* \cup B^*$$

$$(A \cap B)^* = A^* \cup B^*$$

Exempel

$$\Omega = \{ \text{etta, tvåa, trea, fyra, femma, sexa} \}$$

$$A = \{ \text{etta, trea, femma} \}. \quad B = \{ \text{femma, sexa} \}, \quad A \cap B = \{ \text{femma} \}$$

$$(A \cap B)^* = \{ \text{etta, tvåa, trea, fyra, sexa} \}$$

$$\begin{aligned} A^* \cup B^* &= \{ \text{tvåa, fyra, sexa} \} \cup \{ \text{etta, tvåa, trea, fyra} \} \\ &= \{ \text{etta, tvåa, trea, fyra, sexa} \} \end{aligned}$$

Drill på Venndiagram

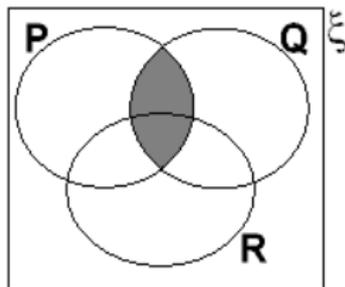


Diagram A

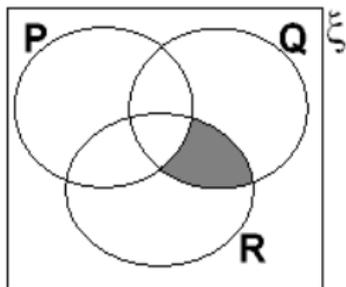


Diagram B

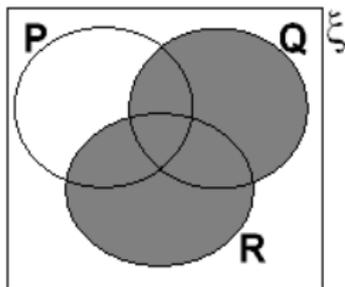


Diagram C

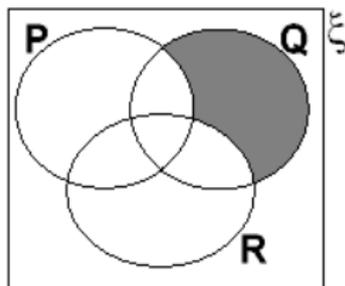


Diagram D

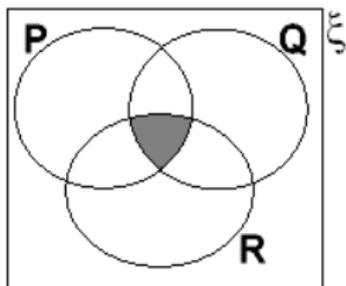


Diagram E

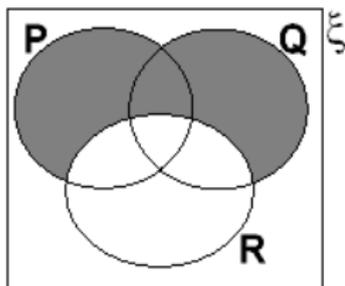


Diagram F

Drill på Venndiagram

- a) Vilket diagram svarar mot $Q \cup R$?
- b) Vilket diagram svarar mot $P \cap Q \cap R$?
- c) Vilket diagram svarar mot $Q \cap R \cap P^*$?
- d) Vilket diagram svarar mot $P \cap Q$?
- e) Vilket diagram svarar mot $P^* \cap R^* \cap Q$?
- f) Vilket diagram svarar mot $(P \cup Q) \cap R^*$?
- g) Vilket diagram svarar mot $Q \cap (P \cup R)^*$?



Har vi många händelser kan vi, precis som med summa- och produkt-tecken, använda ett förkortat skrivsätt:

$$\bigcup_{1}^n A_i = A_1 \cup A_2 \cup \dots \cup A_n \quad \text{och} \quad \bigcap_{1}^n A_i = A_1 \cap A_2 \cap \dots \cap A_n$$

Slumpförsök: tärningskast och $P(A)$

Låt oss säga att vi kastar en tärning, och är intresserade av händelsen
 $\{\text{vi får en sexa}\}$.

Alla håller nog med om att, om det är en just tärning, att den sannolikheten är $\frac{1}{6}$. Symboliskt kan vi skriva

$$A = \{\text{vi får en sexa}\} \quad \text{och} \quad P(A) = \frac{1}{6}.$$

Är det överhuvudtaget meningsfullt att tala om sannolikheter, och om så är fallet, hur skall man tolka dessa?



Vi skall tolka detta som att om man kastar tärningen många gånger, så blir den relativa frekvensen för ungefär $\frac{1}{6}$. Allmänt sett, om vi har ett försök och en händelse A och gör försöket n gånger, så gäller

$$f_n(A) = \frac{\text{antalet gånger } A \text{ inträffar}}{n} \rightarrow P(A) \text{ då } n \text{ växer.}$$

Kolmogorovs axiomsystem (1933):

Ett *sannolikhetsmått* P är en funktion av händelser, sådan att:

(a) $0 \leq P(A) \leq 1$;



Kolmogorovs axiomsystem (1933):

Ett *sannolikhetsmått* P är en funktion av händelser, sådan att:

(a) $0 \leq P(A) \leq 1$;

(b) $P(\Omega) = 1$



Kolmogorovs axiomsystem (1933):

Ett *sannolikhetsmått* P är en funktion av händelser, sådan att:

- (a) $0 \leq P(A) \leq 1$;
- (b) $P(\Omega) = 1$
- (c) om A_1, A_2, \dots är disjunkta händelser, så gäller

$$P\left(\bigcup_1^{\infty} A_i\right) = \sum_1^{\infty} P(A_i).$$

Kolmogorovs axiomsystem (1933):

Ett *sannolikhetsmått* P är en funktion av händelser, sådan att:

- (a) $0 \leq P(A) \leq 1$;
- (b) $P(\Omega) = 1$
- (c) om A_1, A_2, \dots är disjunkta händelser, så gäller

$$P\left(\bigcup_1^{\infty} A_i\right) = \sum_1^{\infty} P(A_i).$$

- (a) och (b) kan ses som en kalibrering så att P stämmer med intuitionen (det blir lättare då) och (c) (som är det "viktiga" axiomet) betyder att P är ett mått.



Kolmogorovs axiomsystem: ett specialfall

Om A_1, A_2, \dots är disjunkta händelser, så gäller

$$P\left(\bigcup_1^{\infty} A_i\right) = \sum_1^{\infty} P(A_i).$$

Ett specialfall: Betrakta A och B sådana att $A \cap B = \emptyset$, då fås

$$P(A \cup B) = P(A) + P(B)$$



Regler för sannolikhetskalkyl (1)

Sats

$$P(A^*) = 1 - P(A).$$

Bevis. Vi ska ge ett mycket formellt bevis, för att illustrera axiomsystemet: Eftersom A och A^* disjunkta och $A \cup A^* = \Omega$, så fås enligt (c) och (b) ovan

$$P(A) + P(A^*) = P(\Omega) = 1 \quad \Rightarrow \quad P(A^*) = 1 - P(A).$$



Följdsats

$$P(\emptyset) = 0.$$

Sats

$$P(A \cup B) = P(A) + P(B) - P(A \cap B).$$

Bevis. Satsen följer med hjälp av Venn-diagram, och observationen att $P(A) + P(B)$ "mäter" $A \cap B$ två gånger. □

Regler för sannolikhetskalkyl (3)

Om $A \cap B = \emptyset$, så fås $P(A \cap B) = P(\emptyset) = 0$, dvs.

$$P(A \cup B) = P(A) + P(B).$$

Detta följer av det ovan visade

Sats

$$P(A \cup B) = P(A) + P(B) - P(A \cap B).$$



Antag att Ω består av m (möjliga) elementarutfall $\omega_1, \dots, \omega_m$, var och en med samma sannolikhet att inträffa, dvs

$$P(\omega_k) = \frac{1}{m} \quad k = 1, \dots, m.$$

Betrakta en händelse A , $A \subset \Omega$. Antag att A innehåller g (gynnsamma) elementarutfall. Då gäller

$$P(A) = \frac{g}{m}.$$

Problemet med den klassiska sannolikhetsdefinitionen, i mera komplicerade situationer, är att hitta en uppdelning av Ω i lika sannolika elementarutfall och att beräkna m och g . I många – de flesta – situationer är det inte alls möjligt att göra detta.

För att beräkna m och g behöver vi några kombinatoriska grundbegrepp. Vi inleder med en s.k. *multiplikationsprincip*.

Sats

Om en åtgärd kan utföras på n_1 och en annan (nästa) åtgärd kan utföras på n_2 sätt, så finns det $n_1 \cdot n_2$ sätt att utföra de båda åtgärderna.



Enligt multiplikationsprincipen finns det n^k olika sätt att plocka ut k st. av n st. föremål om varje föremål som har plockats ut stoppas tillbaka och om vi tar hänsyn till i vilken ordning de plockas ut.

Exempel

Antalet PIN-koder = $10^4 = 10000$.

Enligt multiplikationsprincipen finns det n^k olika sätt att plocka ut k st. av n st. föremål om varje föremål som har plockats ut stoppas tillbaka och om vi tar hänsyn till i vilken ordning de plockas ut.

Exempel

Föremålen är $\heartsuit, \diamondsuit, \spadesuit$, dvs. $n = 3$. Vi tar $k = 2$. $\heartsuit\heartsuit, \heartsuit\diamondsuit, \heartsuit\spadesuit, \diamondsuit\diamondsuit, \diamondsuit\heartsuit, \diamondsuit\spadesuit, \spadesuit\spadesuit, \spadesuit\heartsuit, \spadesuit\diamondsuit$.

Kombinatoriska grundbegrepp: permutation

n st. föremål kan permuteras eller ordnas på

$$n! = n \cdot (n - 1) \dots 2 \cdot 1$$

olika sätt.

$$0! = 1$$



Kombinatoriska grundbegrepp: permutation

n st. föremål kan permuteras eller ordnas på

$$n! = n \cdot (n - 1) \dots 2 \cdot 1 \quad \text{utläses: } n \text{ fakultet}$$

olika sätt.

$$0! = 1$$

Exempel

$$n = 3, 3! = 6. \quad \heartsuit \diamondsuit \spadesuit, \heartsuit \spadesuit \diamondsuit, \diamondsuit \heartsuit \spadesuit, \diamondsuit \spadesuit \heartsuit, \spadesuit \heartsuit \diamondsuit, \spadesuit \diamondsuit \heartsuit.$$

$$n \cdot (n - 1) \dots (n - k + 1)$$

kallas kallas antalet permutationer av k element bland n (=antalet sätt at välja k element bland n utan återläggning och med hänsyn till ordningen). Vi har

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

$$n \cdot (n - 1) \dots (n - k + 1)$$

kallas kallas antalet permutationer av k element bland n (=antalet sätt at
välja k element bland n utan återläggning och med hänsyn till ordningen).

Vi har

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

Exempel

$n = 3, k = 2$. $\heartsuit\spadesuit, \spadesuit\heartsuit, \spadesuit\clubsuit, \clubsuit\spadesuit, \heartsuit\clubsuit, \clubsuit\heartsuit$.

Kombinatoriska grundbegrepp: binomialkoefficienterna

Låt x vara antalet sätt att plocka ut k st. av n st. föremål om vi ej tar hänsyn till i vilken ordning de plockas ut. Då gäller enligt multiplikationsprincipen att

$$\underbrace{n \cdot (n-1) \dots (n-k+1)}_{k \text{ utan återläggning och med ordning}} = \underbrace{k!}_{\text{antalet sätt att ordna } k \text{ element}} \cdot x$$

\Rightarrow

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



KTH Matematik

Vi ger beteckningen $x = \binom{n}{k}$. Alltså: Det finns

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

olika sätt att plocka ut k st. av n om vi ej tar hänsyn till i vilken ordning de plockas ut.

Vi ger beteckningen $x = \binom{n}{k}$. Alltså: Det finns

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

olika sätt att plocka ut k st. av n om vi ej tar hänsyn till i vilken ordning de plockas ut.

Exempel

$n = 3, k = 2$. $\heartsuit\spadesuit, \heartsuit\clubsuit, \spadesuit\clubsuit$.

Två urnmodeller: Dragning utan återläggning

I en urna finns kulor av två slag: v vita och s svarta. Drag n kulor ur urnan slumpmässigt och så att en kula som dragits inte stoppas tillbaka. dvs dragning utan återläggning.

Sätt $A =$ "Man får k vita kulor i urvalet".



Två urnmodeller: Draging utan återläggning

Välj Ω : Alla uppsättningar om n kulor utan hänsyn till ordning.
Då fås:

$$m = \binom{v+s}{n} \quad \text{och} \quad g = \binom{v}{k} \binom{s}{n-k}$$

och således

$$P(A) = \frac{\binom{v}{k} \binom{s}{n-k}}{\binom{v+s}{n}}.$$

Två urnmodeller: Dragning med återläggning

Samma modell som i fallet med dragning utan återläggning, men kulorna stoppas tillbaka igen efter det att man observerat dess färg, och urnan skakas om för nästa dragning.

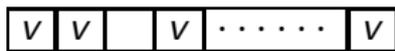
Välj Ω : Alla uppsättningar om n kulor med hänsyn till ordning:

$$m = (v + s)^n.$$



Two urn models: Drawing with replacement

Assume that we have chosen k white and $n - k$ black balls. These can be placed on $\binom{n}{k}$ places:



Number of ways to choose k white = v^k . Number of ways to choose $n - k$ black = s^{n-k} .

Två urnmodeller: Draging med återläggning

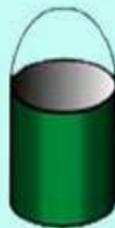
Detta ger $g = \binom{n}{k} v^k s^{n-k}$ och således får vi

$$P(A) = \frac{\binom{n}{k} v^k s^{n-k}}{(v+s)^n} = \binom{n}{k} \left(\frac{v}{v+s} \right)^k \left(\frac{s}{v+s} \right)^{n-k}.$$





Probability: Given the information in the pail, what is in your hand?



Statistics: Given the information in your hand, what is in the pail?



²A.N.Kolmogorov, 1903 - 1987