



# NTNU

Det skapende universitet

## **TMA4240 Statistikk H2010**

### **Kapittel 5: Diskrete sannsynlighetsfordelinger**

#### **5.1-5.4: Uniform, binomisk, hypergeometrisk fordeling**

Mette Langaas

# Arbeidshverdag etter endt studium

- Studere et fenomen (f.eks. kvalitet av produsert maskindel, elektrisitetsprisen, forsikringspremie, holdninger til miljøvennlig energi) ved å beskrive og forstå.
- Med mål å trekke konklusjoner og gjøre beslutninger.
- Trenger data:
  - samler inn data (subjektive eller objektive) under usikkerhet,
  - studerer fenomenet fra data,
  - kan bruke en SV med tilhørende fordeling til å beskrive fenomenet.
- Hvilken fordeling?
  - se på prosessen som har “skapt dataene”.
  - se grafisk på data og studere fordelings form.

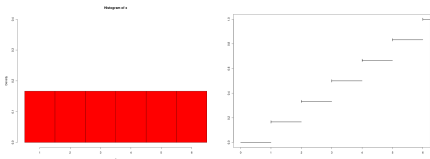
# Arbeidshverdag etter endt studium

- Derfor: kap. 5 og 6: beskrive viktige fordelinger for å
  - lære situasjoner er fordelingen passer
  - forstå hvordan  $f(x)$  fremkommer
  - se hva  $E(X)$  og  $\text{Var}(X)$  er og forstå hvorfor
  - lære å regne ut  $f(x)$ ,  $F(x) = P(X \leq x)$ ,  $P(a < X \leq b)$ .
- Deretter: anslå parametere i fordelingene og trekke konklusjoner under usikkerhet (kap. 9-11).

## 5.2 Diskret uniform fordeling

**Diskret uniform fordeling:** Hvis den stokastiske variabelen  $X$  antar verdiene  $x_1, x_2, \dots, x_k$  med lik sannsynlighet så er  $X$  diskret uniformt fordelt med fordeling

$$f(x; k) = \frac{1}{k}, \quad x = x_1, x_2, \dots, x_k$$



TEO 5.1:

$$\mu = E(X) = \frac{\sum_{i=1}^k x_i}{k} \quad \text{og} \quad \sigma^2 = \text{Var}(X) = \frac{\sum_{i=1}^k (x_i - \mu)^2}{k}$$

# Midtveiseksamen

## Eksamen 06.08.2004, oppg 1a

- Fra høsten 2004 vil det i TMA4240 bli innført tellende skriftlig midtveiseksamen.
- Denne vil bli gitt i form av en flervalgsoppgave (“multiple choice”) bestående av  $n = 20$  spørsmål som alle har  $m$  svaralternativer. Studentene må velge et svaralternativ for hvert spørsmål (det er således ikke lov å svare “blankt” på et spørsmål).
- For å få karakter bedre enn F (36%) må minst 8 spørsmål være korrekt besvart.
- Ole lurer på om han skal la være å lese til midtveiseksamen og heller velge tilfeldige svaralternativer på alle spørsmålene (han vil da ikke engang lese oppgaveteksten før han svarer). Før han bestemmer seg, ber han en studiekamerat regne ut hvor stor sannsynlighet han da har for få bedre enn F.

## Midtveiseeksamen (forts.)

- La  $X$  være antall korrekte svar Ole får på de  $n = 20$  spørsmålene.
- Forklar hvorfor vi kan anta at  $X$  er binomisk fordelt med  $n = 20$  og  $p = \frac{1}{m}$ .
- Finn sannsynligheten for at Ole får bedre enn F hvis han velger å svare tilfeldig på alle spørsmålene, dvs.  $P(X \geq 8)$ , når antall svaralternativer er  $m = 2$ . Finn også  $P(X \geq 8)$  for  $m = 4$  og  $m = 5$ .
- Hva blir forventet antall korrekte svar, dvs.  $E(X)$ , når  $m = 2, 4, 5$ ?

## 5.3 Binomisk fordeling

**Bernoulli prosess:** Et Bernoulli eksperiment (prosess) har følgende egenskaper:

1. Eksperimentet består av  $n$  gjentatte forsøk.
2. Hvert forsøk undersøker man om en hendelse  $A$  inntreffer (suksess) eller ikke ( $A'$ =fiasko).
3. Sannsynligheten for hendelsen  $A$  (suksess) kaller vi  $p$ , og denne er den samme fra forsøk til forsøk.
4. De  $n$  gjentatte forsøkene er uavhengige av hverandre.

— Dermed: et Bernoulli eksperiment kan resultere i

- hendelsen  $A$  (suksess) med sannsynlighet  $p$  og
- komplementet av hendelsen  $A$  ( $A'$ =fiasko) med sannsynlighet  $1 - p$ .

## 5.3 Binomisk fordeling (forts.)

- La den stokastiske variabelen  $X$  være antall ganger hendelsen  $A$  (suksess) inntreffer på de  $n$  uavhengige forsøkene.
- Sannsynlighetsfordelingen til  $X$  kalles *binomisk fordeling* og er gitt ved

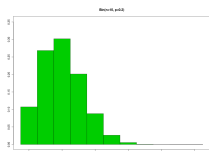
$$f(x) = b(x; n, p) = \binom{n}{x} p^x (1 - p)^{(n-x)}, \quad x = 0, 1, \dots, n$$

- Kumulativ fordeling:  $F(x) = P(X \leq x)$  finnes ved tabelloppslag.
- Eksempler:
  - antall defekte enheter i industriell prosess
  - antall pasienter med positiv effekt av medisin

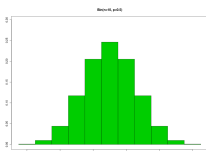


# Binomisk fordeling forts.

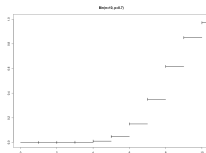
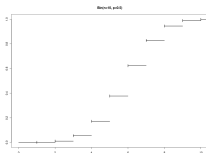
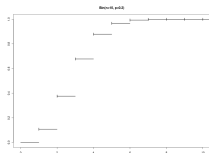
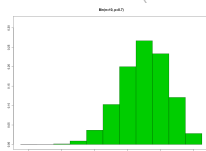
$n = 10, p = 0.2$



$n = 10, p = 0.5$



$n = 10, p = 0.7$



**TEO 5.2:** Forventning og varians i binomisk fordeling  $b(x; n, p)$  er

$$\mu = E(X) = np \quad \text{og} \quad \sigma^2 = \text{Var}(X) = np(1 - p)$$

# Urne med kuler [v1]

- Definisjon:  $p = \frac{\text{antall røde kuler}}{\text{antall kuler}}$
- Prosedyre: Utfør  $n$  ganger
  - trekk en kule tilfeldig
  - registrer fargen
  - legg kula tilbake
- Da er antallet røde kuler du trekker binomisk fordelt.

# Urne med kuler [v2]

— Definisjoner:

- $p_1 = \frac{\text{antall hvite kuler}}{\text{antall kuler}}$
- $p_2 = \frac{\text{antall sorte kuler}}{\text{antall kuler}}$
- $p_3 = \frac{\text{antall blå kuler}}{\text{antall kuler}}$
- $p_4 = \frac{\text{antall røde kuler}}{\text{antall kuler}}$

— Prosedyre: Utfør  $n$  ganger

- trekk en kule tilfeldig og registrer fargen
- legg kula tilbake

— Da er

- antallet hvite kuler og antallet sorte kuler og antallet blå kuler og antallet røde kuler som du trekker er

— multinomisk fordelt.

# Multinomisk fordeling

**Multinomisk fordeling:** Et forsøk kan resultere i

- $k$  mulige utfall  $A_1, A_2, \dots, A_k$ , med sannsynligheter
- $p_1, p_2, \dots, p_k$ .

La de stokastiske variablene  $X_1, X_2, \dots, X_k$  representere antall ganger utfallene  $A_1, A_2, \dots, A_k$  opptrer i  $n$  uavhengige forsøk.

Sannsynlighetsfordelingen til  $X_1, X_2, \dots, X_k$  kalles *multinomisk fordeling* og er gitt ved

$$f(x_1, x_2, \dots, x_k; p_1, p_2, \dots, p_k, n) = \binom{n}{x_1, x_2, \dots, x_k} p_1^{x_1} p_2^{x_2} \cdots p_k^{x_k}$$

$$\text{med } \sum_{i=1}^k x_i = n, \sum_{i=1}^k p_i = 1 \text{ og } \binom{n}{x_1, x_2, \dots, x_k} = \frac{n!}{x_1! x_2! \cdots x_k!}.$$

# Urne med kuler [v3]

- Definisjon:
  - $N$ =antall kuler
  - $k$ =antall røde kuler
- Prosedyre: Utfør  $n$  ganger
  - trekk en kule tilfeldig
  - registrer fargen
  - legg kula til side
- Da er antallet røde kuler du trekker hypergeometrisk fordelt.

# Antall fisker i dammen

- Vi vil anslå størrelsen,  $N$ , av en dyreart innenfor et område (metode fra 1896).
- Gjøre to undersøkelser:
  1. finner og merker  $k$  individ, og slipper dem ut igjen.
  2. finner så  $n$  individ, og  $x$  av disse er merket.
- Lukket populasjon: ingen død, fødsel, innflytting, utflytting.
- Andelen merkede i de to utvalgene bør da være like, bestandsanslag:

$$\frac{k}{N} = \frac{x}{n}$$
$$N = \frac{k \cdot n}{x}$$

- $X$  er hypergeometisk fordelt med parametere  $N, k, n$ .

## 5.4 Hypergeometrisk fordeling

Hypergeometrisk eksperiment: har følgende egenskaper:

1. Vi har en mengde av  $N$  enheter. Av de  $N$  enhetene så klassifiseres  $k$  som hendelsen  $A$  (suksess) og  $N - k$  som komplementet av hendelsen  $A$  ( $A'$ =fiasko).
2. Et tilfeldig utvalg av størrelse  $n$  trekkes uten tilbakelegging fra de  $N$  enhetene.

Antallet ganger,  $X$ , som hendelsen  $A$  (suksess) inntreffer er da en *hypergeometrisk stokastisk variabel*.

## 5.4 Hypergeometrisk fordeling (forts.)

**Hypergeometrisk fordeling:** En hypergeometrisk stokastisk variabel,

- $X$ , angir antallet ganger hendelsen  $A$  (suksess) inntreffer i et hypergeometrisk eksperiment
- der  $n$  enheter trekkes fra  $N$  enheter,
- der  $k$  av de  $N$  enheter er klassifisert som hendelsen  $A$  (suksess) og
- $N - k$  som komplementet av hendelsen  $A$  ( $A'$ =fiasko).

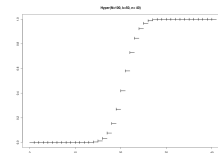
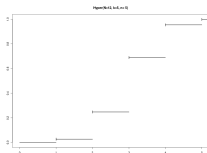
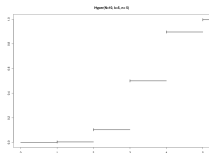
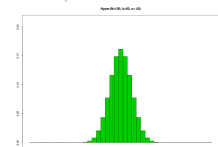
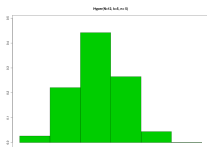
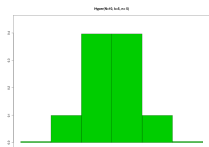
Sannsynlighetsfordelingen til  $X$  kalles en *hypergeometrisk fordeling* og er gitt ved

$$f(x) = h(x; N, n, k) = \frac{\binom{k}{x} \binom{N-k}{n-x}}{\binom{N}{n}} \quad x = 0, 1, 2, \dots, n$$



# Hypergeometrisk fordeling (forts.)

$N = 10, k = 5, n = 5$     $N = 12, k = 5, n = 5$     $N = 100, k = 50, n = 40$



**TEO 5.3:** Forventning og varians i den hypergeometriske fordelingen  $h(x; N, n, k)$  er

$$\mu = E(X) = \frac{nk}{N} \quad \text{og} \quad \sigma^2 = \text{Var}(X) = \frac{N-n}{N-1} \cdot n \cdot \frac{k}{N} \left(1 - \frac{k}{N}\right)$$

# Hypergeometisk og binomisk fordeling

- Hvis  $n$  er liten i forhold til  $N$  ( $\frac{n}{N} \leq 0.05$ ), så vil sammensetningen av de  $N$  enhetene endres lite under trekningen.
- Dermed kan  $\frac{k}{N}$  sees på som den binomiske sannsynligheten  $p$ .
- Dermed kan binomisk fordeling sees på som en “stor populasjon” versjon av hypergeometrisk fordeling.

# Urne med kuler [v4]

- Definisjoner:
  - $N$ =antall kuler
  - $a_1$ =antall hvite kuler
  - $a_2$ =antall sorte kuler
  - $a_3$ =antall blå kuler
  - $a_4$ =antall røde kuler
- Prosedyre: Utfør  $n$  ganger
  - trekk en kule tilfeldig og registrer fargen
  - legg kula til side.
- Da er
  - antallet hvite kuler
  - antallet sorte kuler
  - antallet blå kuler
  - antallet røde kuler
- som du trekker multivariat hypergeometrisk fordelt.

# Multivariabel hypergeometrisk fordeling

Multivariabelt hypergeometrisk eksperiment: har følgende egenskaper:

1. Et tilfeldig utvalg av størrelse  $n$  trekkes uten tilbakelegging fra  $N$  enheter.
2. Av de  $N$  enhetene så klassifiseres  $a_1$  i cellen  $A_1$ ,  $a_2$  i cellen  $A_2, \dots$ ,  $a_k$  i cellen  $A_k$ .
3. Av de  $n$  enhetene så klassifiseres  $x_1$  i cellen  $A_1$ ,  $x_2$  i cellen  $A_2, \dots$ ,  $a_k$  i cellen  $A_k$ .

Sannsynlighetsfordelingen til  $X_1, X_2, \dots, X_k$  kalles *multivariabel hypergeometrisk fordeling*

$$f(x_1, x_2, \dots, x_k; a_1, a_2, \dots, a_k, n) = \frac{\binom{a_1}{x_1} \binom{a_2}{x_2} \dots \binom{a_k}{x_k}}{\binom{N}{n}}$$

med  $\sum_{i=1}^k x_i = n$  og  $\sum_{i=1}^k a_i = N$ .

# Kakelotteri

- 300 lodd fordelt på 3 farger (100 av hver)
- 9 vinnerlodd, 3 av hver farge (33,66,99)
- Vi kjøper 5 lodd.
- To strategier:
  - trekk 5 lodd blant de 300 loddene
  - trekk 5 lodd av samme farge
- Hvilken strategi gir størst sjanse for å vinne?