

# Statistik och Dataanalys I

## Föreläsning 15 - Sannolikhetsmodeller för diskreta variabler

**Mattias Villani**



Statistiska institutionen  
Stockholms universitet



mattiasvillani.com



@matvil



@matvil



mattiasvillani

- Bernoulliförsök
- Geometrisk fördelning
- Binomialfördelning

# Bernoulliförsök

## ■ Bernoulliförsök

- 1 Bara **två möjliga utfall**: lyckas/misslyckas.
- 2 **Samma sannolikhet** för lyckas,  $p$ , i alla försök.
- 3 **Oberoende försök**.

## ■ Typexempel: **slantsingling**.

- ▶ Lyckas = Kona, Misslyckas = Klave.
- ▶ Sannolikhet  $p = 0.5$  för schysst mynt.
- ▶ Utfall på en singling beror inte på andra singlar.

## ■ Lyckas/Misslyckas är bara en benämning.

## ■ Död/Levande. Hel/Trasig. Spam/Ham.

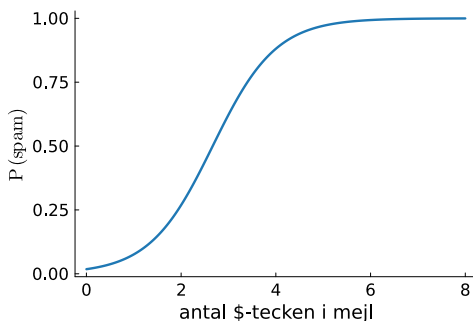


Utan återläggning  $\Rightarrow$  inte samma  $p$  i olika försök:

- ▶  $P(1:a \text{ kortet } \spadesuit) = \frac{13}{52} = \frac{1}{4}$
- ▶  $P(2:a \text{ kortet } \spadesuit) = \frac{12}{51}$  om 1:a  $\spadesuit$  eller  $\frac{13}{51}$  om 1:a  $\heartsuit, \diamondsuit, \clubsuit$ .

# Motivation - regression med binära y-variabler

- Bernoulli-fördelning med **samma sannolikhet**  $p$ .
- Spamdata: lära oss om  $p = P(\text{spam})$  från data.  $\hat{p} = 0.9$ . 🙄
- **Spam-filter**: ska datorn skicka **just detta mejl** till Spam?
- SDAll: **Logistisk regression** där spam sannolikheten  $p$  **beror på förklarande variabler**, som i regression. 🤖

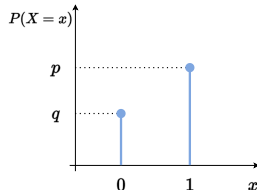


# Bernoullifördelning

- Två möjliga utfall: lyckad/misslyckad. **Binär variabel**.
- Vi kan koda **lyckat = 1**, **misslyckat = 0**.

$$X = \begin{cases} 1 & \text{om Bernoulli-försök lyckat} \\ 0 & \text{om Bernoulli-försök misslyckat} \end{cases}$$

$$P(X = x) = \begin{cases} p & \text{för } x = 1 \\ q = 1 - p & \text{för } x = 0 \end{cases}$$



## ■ Väntevärde och Varians

$$\begin{aligned} E(X) &= \mu = \sum_{\text{alla } x} x \cdot P(x) = 0 \cdot P(X=0) + 1 \cdot P(X=1) \\ &= 0 \cdot q + 1 \cdot p = p \end{aligned}$$

$$\begin{aligned} \text{Var}(X) &= \sum (x - \mu)^2 \cdot P(x) = (0 - p)^2 \cdot q + (1 - p)^2 \cdot p \\ &= p^2 q + q^2 \cdot p = pq \underbrace{(p + q)}_1 = pq \end{aligned}$$

# Geometrisk fördelning

- Email: **spam** eller **ham** (icke-spam). **Lyckas = ham**.

- ▶  $P(\text{ham}) = p = 0.1$

- ▶  $P(\text{spam}) = q = 1 - 0.1 = 0.9$ .

- Hur många mejl måste du öppna tills du får ditt första ham?

$$P(\text{första ham på fjärde mejlet}) = \overbrace{0.9 \cdot 0.9 \cdot 0.9}^{\text{gänger pga oberoende}} \cdot \underbrace{0.1}_{\text{ham}} = 0.9^3 \cdot 0.1 = 0.0729$$

3 spam

- Vad är sannolikheten för  $x$  st mejl tills första ham?

$$P(\text{första ham på } x\text{:te mejlet}) = 0.9^{x-1} \cdot 0.1$$

- **Geometrisk slumpvariabel** från Bernoulliförsök

$X$  = antal försök **tills första lyckade** inträffar

- **Geometrisk fördelning**

$$P(X = x) = q^{x-1} p, \quad \text{för } x = 0, 1, 2, \dots$$



**$X$  inkluderar** försöket där du först lyckas.

Wikipedia kallar detta för **för-första-gången-fördelning**.

# Geometrisk fördelning

## Geometrisk fördelning

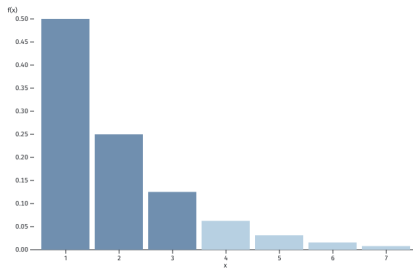
$p$  :    
Kvantil:  

Om  $X \sim \text{Geo}(0.5)$  så gäller att

$$E(X) = \frac{1}{p} = 2.00$$

$$\text{Var}(X) = \frac{1-p}{p^2} = 0.250$$

$$P(X \leq 3) = 0.8750$$



# Geometrisk fördelning i R

- $X \sim \text{Geom}(p = 0.4)$ . Sannolikheten  $p$  kallas `prob` i R.

Beräkning	R kommando
$P(X = 2)$	<code>dgeom(x = 2, prob = 0.4)</code>
$P(X \leq 2)$	<code>pgeom(q = 2, prob = 0.4)</code>
Kvantil	<code>qgeom(p = 0.5, prob = 0.4)</code>
10 slumpstal	<code>rgeom(n = 10, prob = 0.4)</code>

⚠ R använder Wikipedias definition av geometrisk fördelning.  $X$  räknar **antalet misslyckade försök innan** första lyckade. Fix:

```
y = rgeom(n = 100, prob = 0.5) # y is number of trials BEFORE first success
x = y + 1                      # x is number of trials INCLUDING first success
```

- Se programkoden [geometric.R](#) på kurssidan.



# Binomialfördelning

## ■ Geometrisk fördelning:

- ▶ Hur många Bernoulli-försök tills första lyckade?
- ▶ Antal försök är slumpmässigt.

## ■ Binomialfördelning:

- ▶ Hur många lyckade i  $n$  Bernoulli-försök med sannolikhet  $p$ .
- ▶ Antal försök  $n$  är förbestämt och fixerat.
- ▶ Antal lyckade är slumpmässigt.

## ■ Vi skriver $X \sim \text{Bin}(n, p)$ och säger:

## ■ “ $X$ är binomialfördelad med parametrar $n$ och $p$ .”

## ■ Binomial: summan av $n$ oberoende Bernoullivariabler

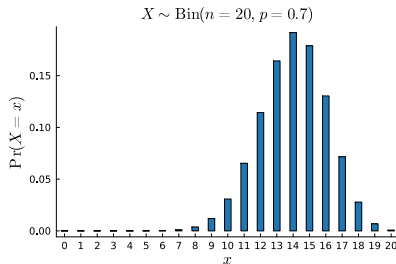
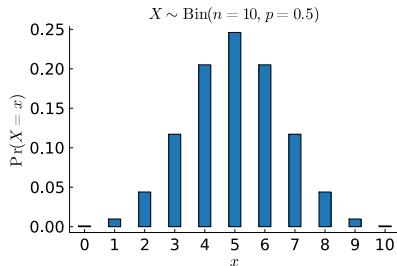
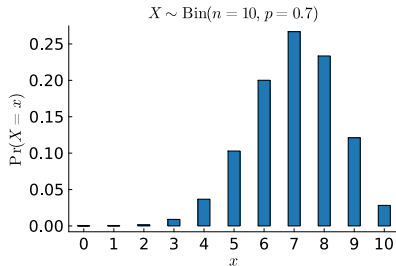
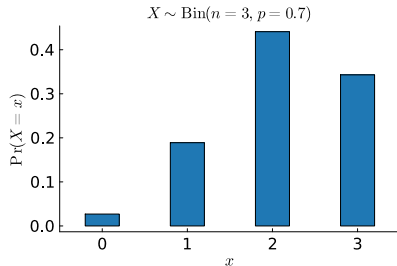
$$X = X_1 + X_2 + \dots + X_n$$

## ■ Exempel: $n = 3$ försök med resultat:

$X_1 = 1$  (Krona första),  $X_2 = 1$  (Krona andra) och  $X_3 = 0$  (Klave tredje).

$$X = 1 + 1 + 0 = 2 \text{ st lyckade (Krona).}$$

# Binomialfördelning



# Binomialfördelning - väntevärde

- Väntevärde i en binomialfördelning? 🤪

$$E(X) = \sum_{x=0}^n x \cdot P(x)$$

**Väntevärde - summa av slumpvariabler.**

$$E(X_1 + X_2 + \dots, X_n) = E(X_1) + E(X_2) + \dots + E(X_n)$$

- Väntevärde för varje Bernoulli-variabel:  $E(X_i) = p$ .

- **Väntevärde för  $X \sim \text{Bin}(n, p)$**

$$E(X) = E(X_1) + E(X_2) + \dots + E(X_n) = \underbrace{p + p + \dots + p}_{n \text{ st}} = np$$

# Binomialfördelning - varians

- Varians i en binomialfördelning? 🤔🤔🤔

$$\text{Var}(X) = \sum_{x=0}^n (x - \mu)^2 \cdot P(x)$$

**Varians - summa av oberoende slumpvariabler.**


$$\text{Var}(X_1 + X_2 + \dots, X_n) = \text{Var}(X_1) + \text{Var}(X_2) + \dots + \text{Var}(X_n)$$


- Bernoulliförsök är oberoende. ✓
- Varians för varje Bernoulli-variabel:  $\text{Var}(X_i) = pq$ .
- **Varians för  $X \sim \text{Bin}(n, p)$**


$$\text{Var}(X) = \text{Var}(X_1) + \dots + \text{Var}(X_n) = \underbrace{pq + pq + \dots + pq}_{n \text{ st}} = npq$$

# Binomialfördelning - interaktivt

## Binomialfördelningen

$n$  :  

$p$  :  

Kvantil:  

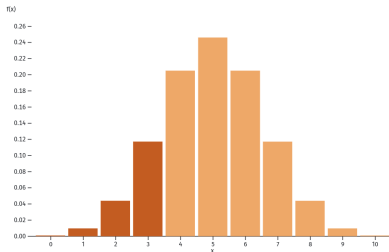
Visa  
normalapproximation ☐

Om  $X \sim \text{Binom}(10, 0.5)$  så gäller att

$$E(X) = np = 5.00$$

$$\text{Var}(X) = np(1-p) = 2.50$$

$$\text{Exakt: } P(X \leq 3) = 0.1719$$



# Binomialfördelningens sannolikheter

- Om  $X \sim \text{Bin}(n, p)$  - vad är egentligen  $P(X = x)$ ?
- Sannolikheten att få  $\{1, 1, 0\}$  i  $n = 3$  försök?

$$p \cdot p \cdot q = p^2 q^1$$

- Det finns dock **flera sätt att få  $X = 2$**  i  $n = 3$  försök:

1:a försök	2:a försök	3:e försök	$X$	$P(X = x)$
1	1	0	2	$p^2 q$
1	0	1	2	$p^2 q$
0	1	1	2	$p^2 q$

- Eftersom dessa tre olika sätt att få  $X = 2$  är **disjunkta**:

$$P(X = 2) = 3 \cdot p^2 q$$

- På samma sätt

$$P(X = 0) = P(\{0, 0, 0\}) = 1 \cdot q^3$$

$$P(X = 1) = P(\{1, 0, 0\}, \{0, 1, 0\}, \{0, 0, 1\}) = 3 \cdot p q^2$$

$$P(X = 2) = P(\{1, 1, 0\}, \{1, 0, 1\}, \{0, 1, 1\}) = 3 \cdot p^2 q$$

$$P(X = 3) = P(\{1, 1, 1\}) = 1 \cdot p^3$$

# Binomialfördelningens sannolikheter

- **Sannolikhetsfördelning**  $X \sim \text{Bin}(3, p)$

$x$	0	1	2	3
$P(x)$	$q^3$	$3 \cdot pq^2$	$3 \cdot p^2q$	$p^3$

- Kolla att summan av alla sannolikheter är ett:

$$q^3 + 3 \cdot pq^2 + 3 \cdot p^2q + p^3 = (p + q)^3 = 1^3 = 1$$

- Allmänna fallet  $X \sim \text{Bin}(n, p)$

$$P(X = x) = {}_nC_x \cdot p^x q^{n-x}$$

- ${}_nC_x$  är antalet sätt ordna  $x$  st 1:or bland  $n$  observationer.

## Kombinationer och permutationer

Hur många sätt att välja $k$ element bland $n$ element?		
	med återläggning	utan återläggning
med ordning	$n^k$	${}_nP_k = \frac{n!}{(n-k)!}$
utan ordning	ej på kurs	${}_nC_k = \frac{n!}{(n-k)!k!}$

# Approximera binomialfördelning med normal

- Om  $X \sim \text{Bin}(n, p)$  så

$$E(X) = \mu = np$$

och

$$SD(X) = \sigma = \sqrt{npq}$$

- **Normalapproximation** av binomialfördelning

$$X \stackrel{\text{approx}}{\sim} N(np, \sqrt{npq})$$

- Approximationen är tillräckligt bra om


$$np \geq 10 \text{ och } nq \geq 10$$

- Man kan också göra en **kontinuitetskorrektur** som korrigerar för att vi approximerar en diskret fördelning (binomial) med en kontinuerlig (normal), se SDM-boken kapitel 15.5.



# Normalapproximation av binomial - interaktivt

## Binomialfördelningen

$n$  :    
 $p$  :    
Kvantil:  

Visa  
normalapproximation ☒

Om  $X \sim \text{Binom}(10, 0.5)$  så gäller att

$$E(X) = np = 5.00$$

$$\text{Var}(X) = np(1-p) = 2.50$$

Exakt:

$$P(X \leq 3) = 0.1719$$

Normal approx:

$$P(X \leq 3) = 0.1030$$

Normal approx med kontinuitetskorrektion:

$$P(X \leq 3) = 0.1714$$

