

INF5830 – 2015 FALL

NATURAL LANGUAGE PROCESSING

Jan Tore Lønning, Lecture 3, 1.9

Today: More statistics

2

- Recap
- Probability distributions
- Categorical distributions
 - ▣ Bernoulli trial
 - ▣ Binomial distribution
- Continuous random variables/distributions
 - ▣ Normal distribution
- Sampling and sampling distribution

3

Recap

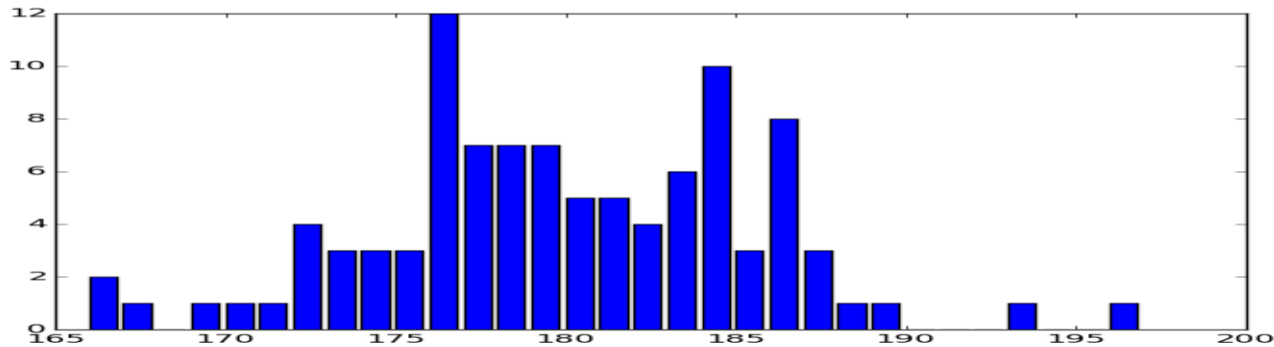
From Lecture 1: Looking at data

4

- Median, mean mode
- Median, quartiles
- Variance, standard deviation

Median, mean mode

5



□ 3 ways to define “middle”, “average”

▣ **Median:** equally many above and below, in the example: 179

▣ **Mean:** ex: 179.54

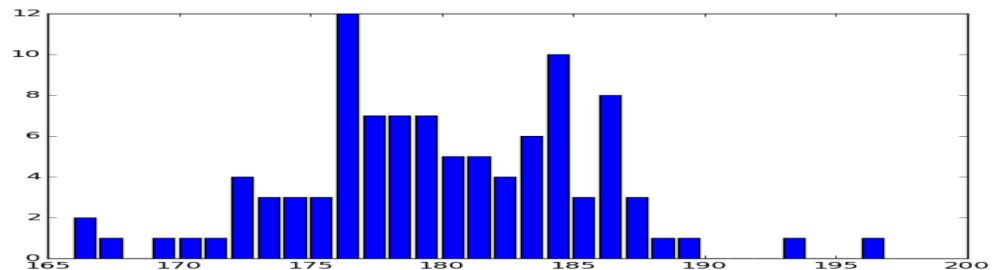
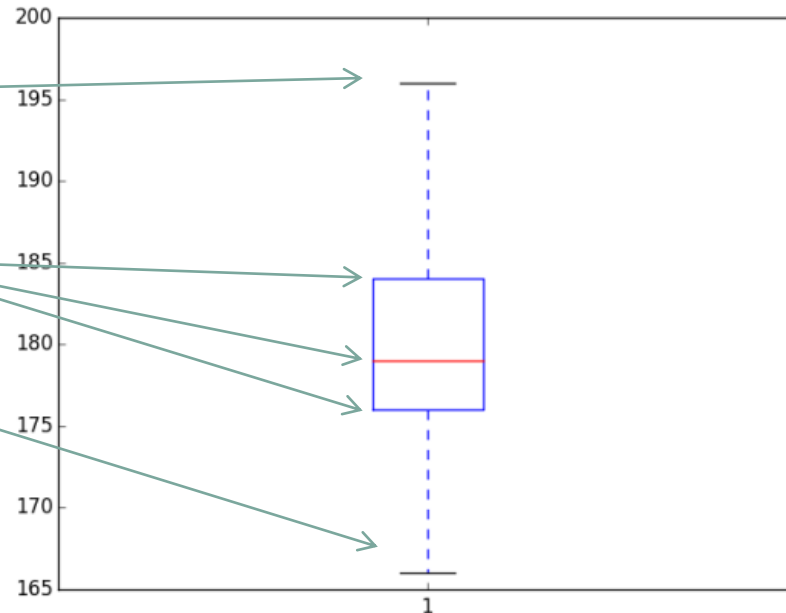
$$\blacksquare \bar{x} = (x_1 + x_2 + \dots + x_n) / n = \frac{1}{n} \sum_{i=1}^n x_i$$

▣ **Mode,** the most frequent one, ex: 176

Dispersion 1: Median, quartile

6

- Example 1:
 - Max 196
 - Quartiles:
 - 176, 179, 184
 - Min 166
- Also good for continuous data
- (The exact definition may vary when “outliers”)



Dispersion 2: Variance

7

- Mean: $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$
- Variance: $\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$
- Idea:
 - ▣ Measure how far each point is from the mean
 - ▣ Take the average
 - ▣ Square – otherwise the average would be 0
- Standard deviation: \sqrt{Var}
 - ▣ “Correct dimension and magnitude”

Beware:
For some purposes we will
later on divide by (n-1)
instead of n.
We return to that!

From Tutorial 1: Probabilities

8

- Median, mean mode
- Median, quartiles
- Variance, standard deviation

Mean of a discrete random variable

9

- The **mean** (or **expectation**) (**forventningsverdi**) of a discrete random variable X :

$$\mu_X = E(X) = \sum_x p(x)x$$

- Useful to remember

$$\mu_{(X+Y)} = \mu_X + \mu_Y$$

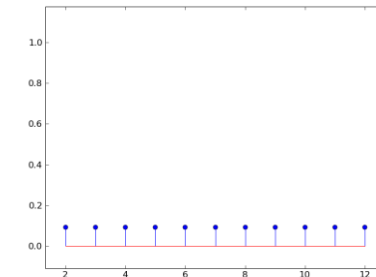
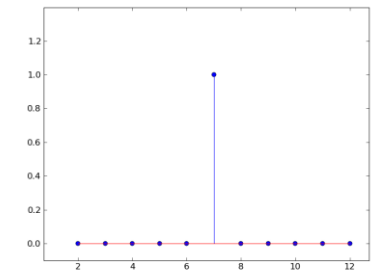
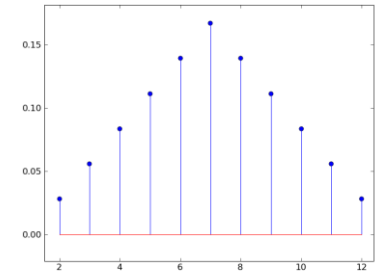
$$\mu_{(a+bX)} = a + b\mu_x$$

Examples:
One dice: 3.5
Two dices: 7
Ten dices: 35

More than mean

10

- Mean doesn't say everything
- Example
 - (1.3) The sum of the two dice, Z , i.e.
 - $p_Z(2) = 1/36, \dots, p_Z(7) = 6/36$ etc
 - (3.2) p_2 given by:
 - $p_2(7)=1$
 - $p_2(x)=0$ for $x \neq 7$
 - (3.3) p_3 given by:
 - $p_3(x)=1/11$ for $x = 2,3,\dots,12$
 - Have the same mean but are very different



Variance

11

- The **variance** of a discrete random variable X

$$\text{Var}(X) = \sigma^2 = \sum_x p(x)(x - \mu)^2$$

- Observe that

$$\text{Var}(X) = E((X - E(X))^2)$$

- It may be shown that this equals $E(X^2) - (E(X))^2$
- The **standard deviation** of the random variable

$$\sigma = \sqrt{\text{Var}(X)}$$

Summary: tutorial 1

12

- Probability space
 - Random experiment (or trial) (no: forsøk)
 - Outcomes (utfallene)
 - Sample space (utfallsrommet)
 - An event (begivenhet)
 - Bayes theorem
- Discrete random variable
 - The probability mass function, pmf
 - The cumulative distribution function, cdf
 - The mean (or expectation) (forventningsverdi)
 - The variance of a discrete random variable X
 - The standard deviation of the random variable

13

Probability distributions

Sannsynlighetsfordelinger

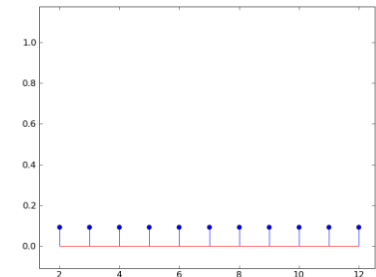
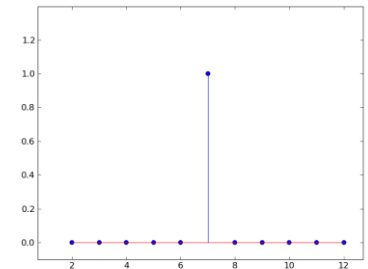
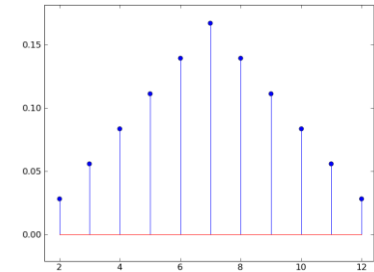
Examples of distributions

14

- (1.3) The sum of the two dice, Z , i.e.
 - $p_Z(2) = 1/36, \dots, p_Z(7) = 6/36$ etc

- (3.2) p_2 given by:
 - $p_2(7)=1$
 - $p_2(x)=0$ for $x \neq 7$

- (3.3) p_3 given by:
 - $p_3(x)=1/11$ for $x = 2,3,\dots,12$



Examples of variance

15

- Throwing one dice
 - $\mu = (1+2+\dots+6)/6=7/2$
 - $\sigma^2 = ((1-7/2)^2 + (2-7/2)^2 + \dots + (6-7/2)^2)/6 = (25+9+1)/4*3=35/12$

- (Ex 1.3) Throwing two dice: $\sigma^2 = 35/6$

- (Ex 3.2) p_2 , where $p_2(7)=1$ has variance 0

- (Ex 3.3) p_3 , the uniform distribution, has variance:
 - $((2-7)^2 + \dots + (12-7)^2)/11 = (25+16+9+4+1+0)*2/11 = 10$

Categorical distributions

- Bernoulli trial
- Binomial distribution

Bernoulli trial

17

- One experiment, two outcomes
- $\Omega_x = \{0, 1\}$
- Write p for $p(1)$
- Then $p(0) = 1 - p$

Examples:

- Flipping a fair coin, $p = 1/2$
- Rolling a dice, getting a 6, $p = 1/6$

- The mean/expectation: $0 * p(0) + 1 * p(1) = 0 + p = p$

- Variance $Var(X) = \sigma^2 = \sum_x p(x)(x - \mu)^2 =$

$$(1 - p)(0 - p)^2 + p(1 - p)^2 = p(1 - p)$$

- Standard deviation $\sigma = \sqrt{p(1 - p)}$

Bernoulli trial and binomial distribution

18

- The Bernoulli trial seems trivial, but can be used as a lego block for more interesting models
- Binomial:
 - n trials
 - let X be a random variable counting the number of successes
 - Possible values $\{0, 1, 2, \dots, n\}$
 - Consider the distribution $p(k)$
- Geometric: how many trials before the first success?

Sampling

19

Ordered sequences:

- Choose k items from a population of n items with replacement: n^k
- Without replacement (permutation):
 - : $n(n-1)(n-2)\dots(n-k+1) = \frac{n!}{(n-k)!}$

Unordered sequences

- Without replac.: $\frac{1}{k!} \left(\frac{n!}{(n-k)!} \right) = \left(\frac{n!}{k!(n-k)!} \right) = \binom{n}{k}$
 - = (The number of ordered sequences /
 - (The number of ordered sequences containing the same k elements *)
 - The number of ordered sequences containing the same $(n-k)$ elements)

Binomial distribution

20

- **Binomial distribution** (binomisk fordeling)
- Conducting n Bernoulli trials with the same probability and counting the number of successes

- Example flipping a fair coin n times, $p(k)$:
 - $n=2$: $p(0)=1/4$, $p(1)=1/2$, $p(2)=1/4$
 - $n=3$: $p(0)=1/8$, $p(1)=3/8$, $p(2)=3/8$, $p(3)=1/8$
 - $n=4$: $(1,4,6,4,1)/16$
 - $n=5$: $(1,5,10,5,1)/32$

- n :
$$p(k) = \binom{n}{k} \left(\frac{1}{2}\right)^n$$
 where
$$\binom{n}{k} = \frac{n!}{k!(n-k)!}$$

□

Binomial distribution

21

□ **Binomial distribution** (binomisk fordeling)

□ **General form:**

□ $0 < p < 1$

□ n a natural number

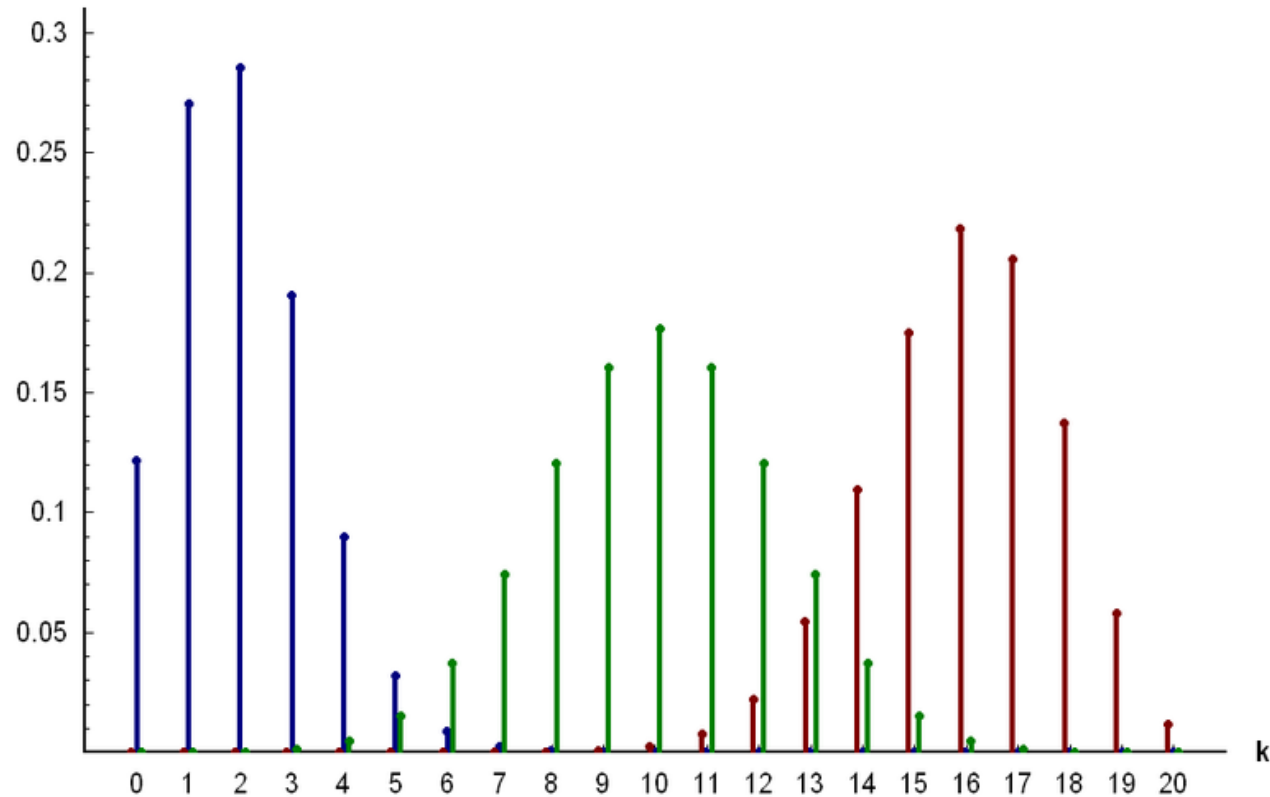
□ **B(n,p)** is given by $b(k; n, p) = \binom{n}{k} p^k (1-p)^{(n-k)}$

for $k = 0, 1, \dots, n$, where $\binom{n}{k} = \frac{n!}{k!(n-k)!}$

Binomial distribution

22

Wahrscheinlichkeit



- $n = 20$
- $p = 0.1$ (blue), $p = 0.5$ (green) and $p = 0.8$ (red)

Binomial distribution

23

- Mean/expectation, μ , of $B(n,p)$ is np
 - n Bernoulli trials
 - Each Bernoulli trial has mean p
- The variance is $np(1-p)$
 - Because the Bernoulli trials are independent
 - Each Bernoulli trial has variance $p(1-p)$

The variance of the sum of two independent random variables is the sum of their variances

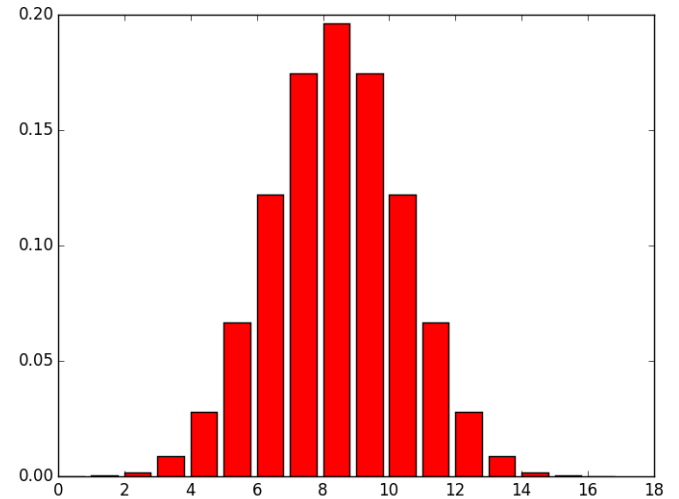
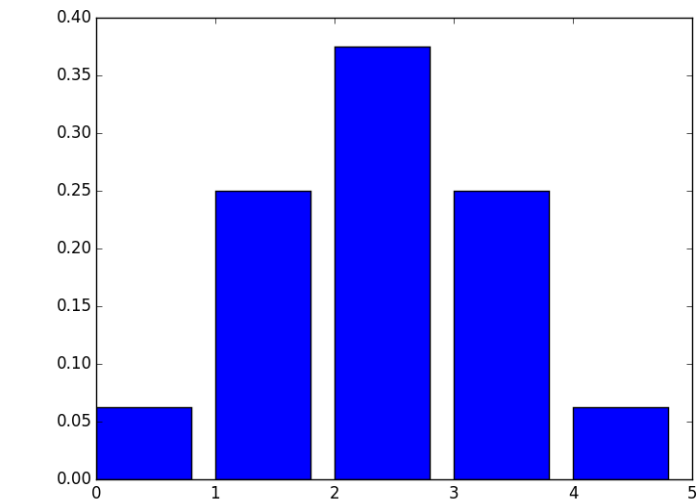
24

N=4:

p=0.5

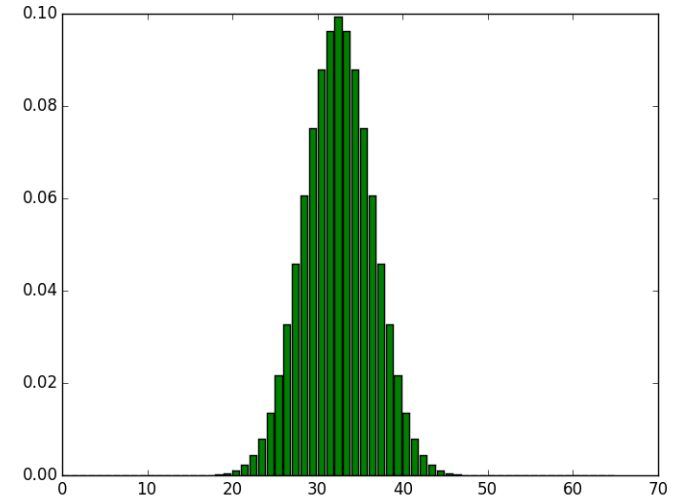
N=16:

N=64:



N	1	4	16	64	256
σ^2	0.25	1	4	16	64
σ	0.5	1	2	4	8

- The relative variation gets smaller with growing N
- The pmf graph approaches a bell shape



SciPy

25

- `import scipy`
- `from scipy import stats`
- `bin10 = stats.binom(10, 0.5) # N=10, p=0.5`
- `bin10.pmf(3) # probability mass of 3, b(3;10,0.5)`
- `bin10.cdf(3) # cumulative distribution function at 3`
- `bin10.var() # variance`
- `bin10.std() # standard deviation`
- `x = np.arange(11); plt.bar(x, bin_10.pmf(x))`

26

Continuous random variables

The normal distribution

Continuous random variables

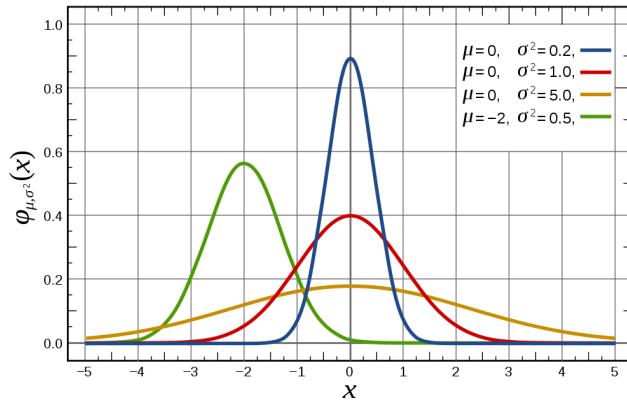
27

- $P(X=a) = 0$ for nearly all values a
- The probability mass function does not make sense
- The **cumulative distribution function**, cdf, given by $F(a) = P(X \leq a)$ makes sense
- $P(a \leq x \leq b) = F(b) - F(a)$
- To calculate expectation and variance we must use integration instead of (infinite) sums.
 - ▣ We skip the details!

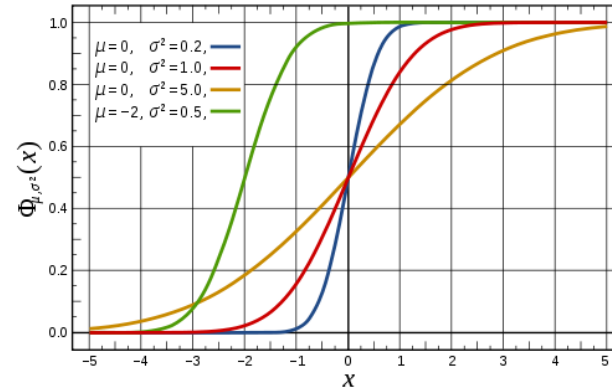
Probability density function

28

pdf



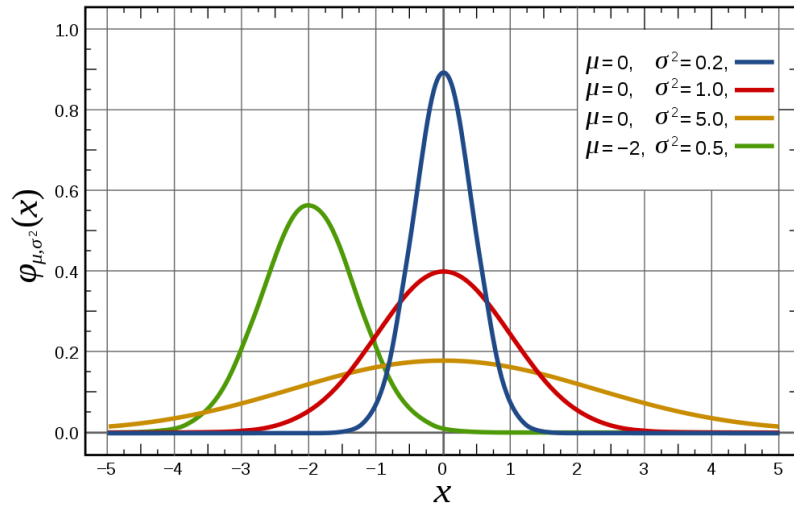
cdf



- The derivative of the cdf, F' , is called the **probability density function**, pdf (sannsynlighetstetthet)
- We draw curves for pdf-s
- The pdf has a similar relationship to the cdf in the continuous case as the pmf has in the discrete case

The normal distribution

29



z-score relates the general case to the standard case
Tells how far x is from μ in terms of standard deviations

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

Standard norm.dist.
(red curve)

General norm.dist
 $N(\mu, \sigma)$

Scary formula

(Don't have to remember)

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$$

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Important

Mean

0

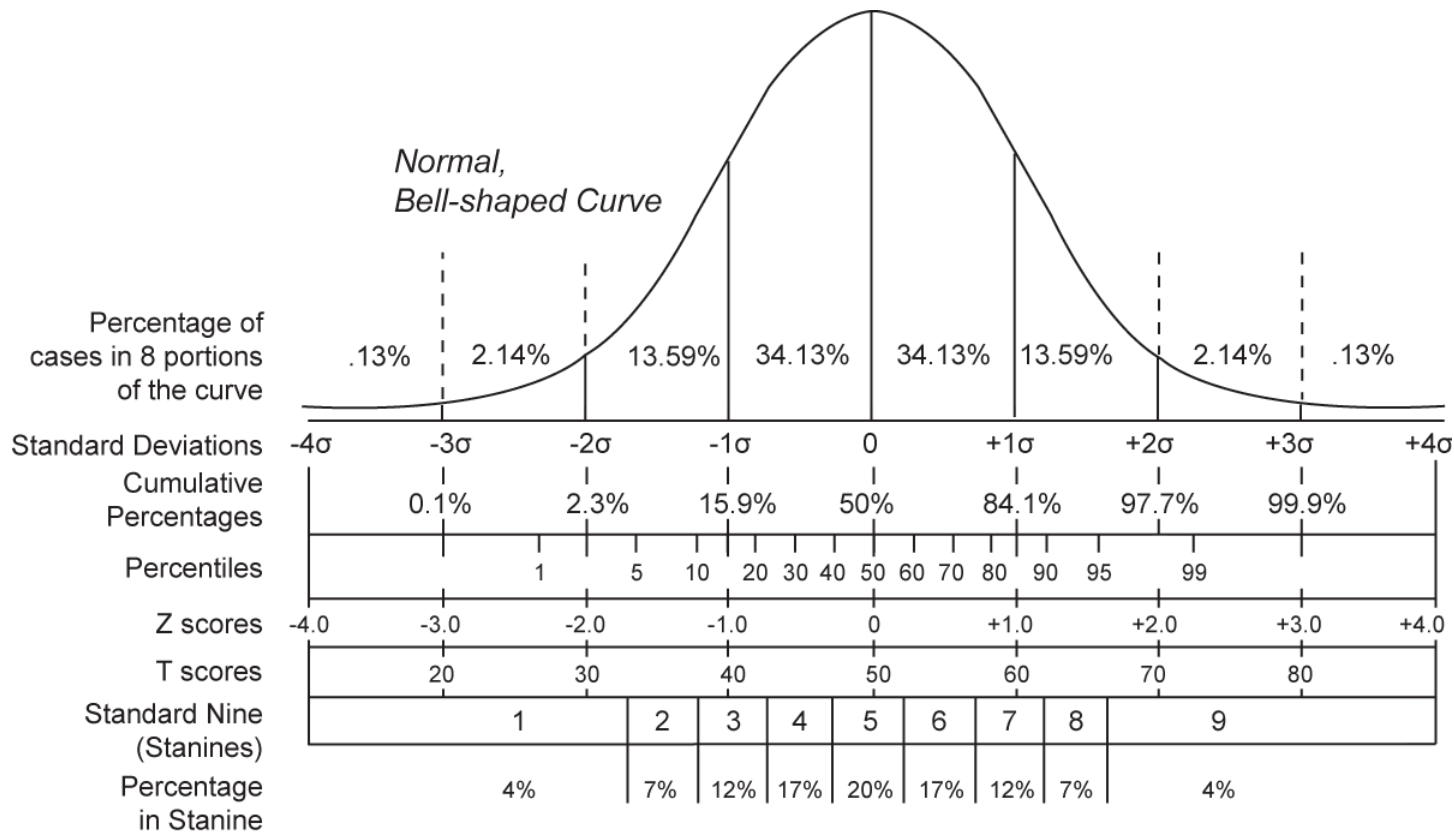
μ

Standard deviation

1

σ

68% - 95% - 99.7%



Example

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

31

□ Tallness of Norwegian young men (rough numbers):

- $\mu = 180$ cm

- $\sigma = 6$ cm

- $z = (186 - 180) / 6 = 1$
(standard deviation)

- $(100 - 68) / 2\% =$
16% are taller than 186 cm



- How many are taller than 190 cm?

- $z = (190 - 180) / 6 = 1.67$

- Prob. = 0.0475 (from table or software)

32

Sampling distribution

Utvalgsfordeling

Sampling - empirically

33

Goal:

- make assertions about a whole **population**
 - from observations of a **sample** (**utvalg**)
-
- A **simple random sample (SRS)** (**tilfeldig utvalg**):
 1. Each individual has equal chance of being chosen (**unbiased**/**forventningsrett**)
 2. Selection of the various individuals are independent
 - Not as simple as it sounds (c.f. the current election polls):
 - ▣ Various methods to rescue
 - ▣ E.g. choose from known groups, weigh by group size (gender, age, home town, etc.)

Sampling in Language Technology

34

- You want to take a simple random sample of words from a corpus?
 - ▣ Can you use the n first sentences?
 - ▣ Can you use a random sample of n sentences?
- How can you build a corpus (sample) which gives a random sample of Norwegian texts?

Sampling distributions – Example

35

- Height: X
 - assume $N(180, 6)$
 - ($\text{Var}=36$)
- Randomly choose 100.
- Add their heights:
 $S = X_1 + X_2 + \dots + X_n$
- A new random variable
(all such samples)
 - $\text{Exp}(S) = n \cdot \mu = 18000$ (cm)
 - $\text{Var}(S) = 100 \cdot \text{Var}(X) = 3600$
 - $\sigma_S = 10 \times \sigma_X = 60$ (cm)



Source: Wikipedia

Sampling distributions – Example

36

- Height: X
 - assume $N(180, 6)$
 - ($\text{Var}=36$)
- Randomly choose 100.
- Add their heights:
 $S = X_1 + X_2 + \dots + X_n$
- A new random variable
(all such samples)
 - $\text{Exp}(S) = n \cdot \mu = 18000$ (cm)
 - $\text{Var}(S) = 100 \cdot \text{Var}(X) = 3600$
 - $\sigma_S = 10 \times \sigma_X = 60$ (cm)

- The mean of the samples:
 - $\bar{X} = S/n$
- A new random variable
(all such means of samples of 100)
 - $\text{Exp}(S) = \mu = 180$ (cm)
 - $\sigma_{\bar{X}} = \frac{1}{100} \times \sigma_S = 0.6$ (cm)

Sampling distributions

37

□ Let

- X be a random variable for a population with exp: μ , std: σ
- Let $S = X_1 + X_2 + \dots + X_n$, i.e. each X_i equals X
- Let : $\bar{X} = S/n$

□ Then:

- $\text{Exp}(S) = n \cdot \mu$
- $\text{Exp}(\bar{X}) = \mu$
- $$\text{Var}(S) = \sigma_S^2 = n \times \text{Var}(X) = n \times \sigma_X^2$$
- $$\text{Var}(\bar{X}) = \sigma_{\bar{X}}^2 = \frac{1}{n^2} \times \text{Var}(S) = \frac{1}{n} \times \sigma_X^2$$
- $$\sigma_{\bar{X}} = \frac{1}{\sqrt{n}} \times \sigma_X$$

Effect of sample size

38

Sample size	1	4	16	100	400	1600
Standard dev.	6	3	1.5	0.6	0.3	0.15

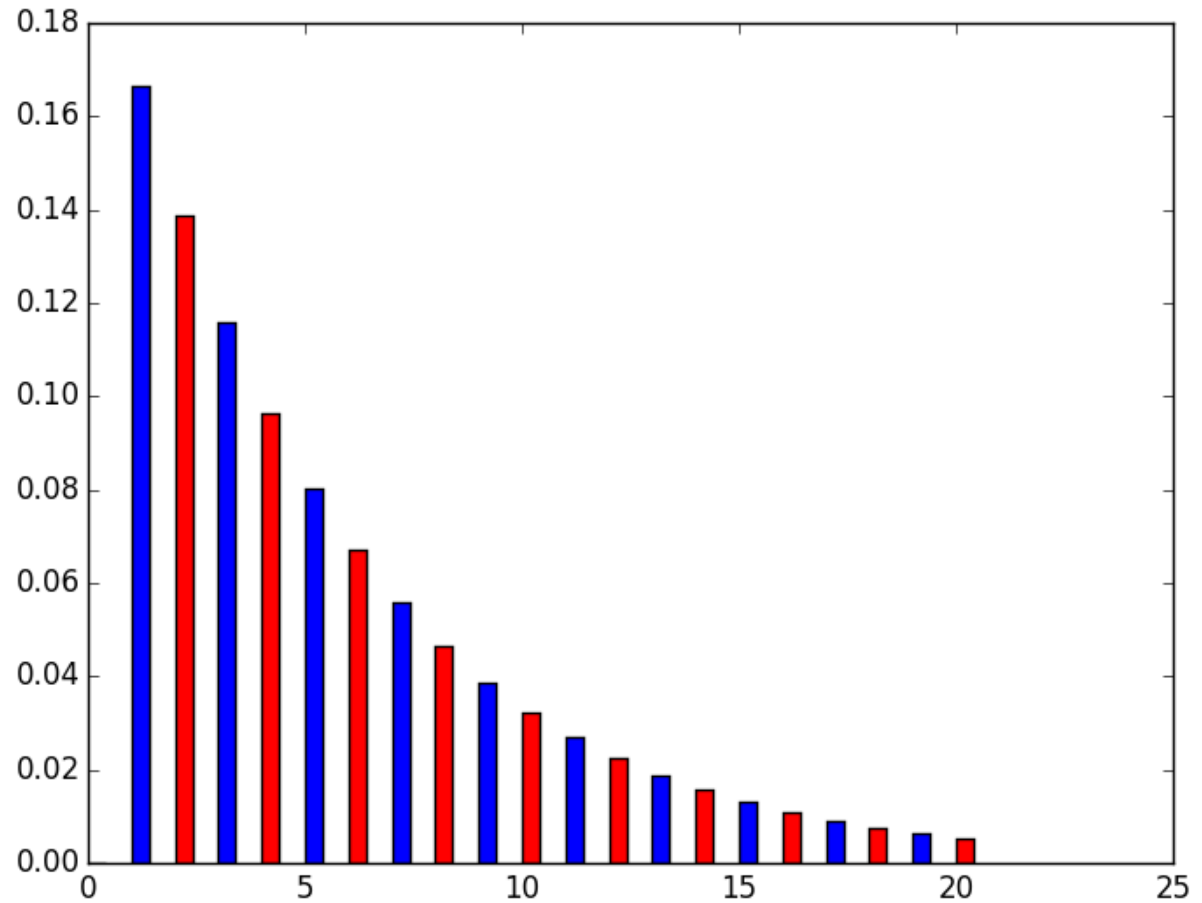
The form of the distribution

39

- If the X_i -s are independent and normally distributed, then \bar{X} is normally distributed (as expected)
- (More surprisingly) Even though the X_i -s are not normally distributed: for large n -s, the sample distribution is approximately normal
- = Central Limit Theorem

Example

40



□ Throwing a dice until you get 6

□ $pmf(n) = \frac{1}{6} \left(\frac{5}{6}\right)^{(n-1)}, n \geq 1$

□ $\mu = 6$

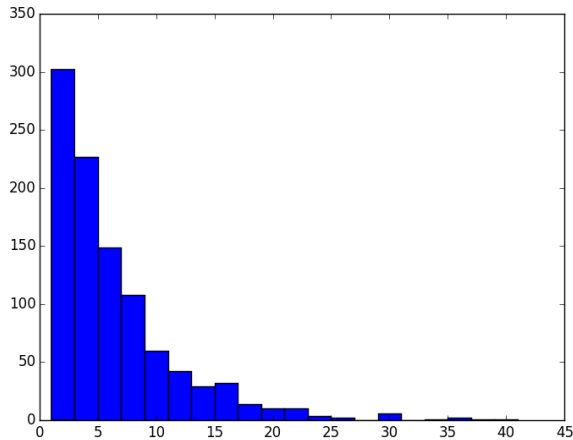
□ $pmf(n) = \frac{1}{6} \left(\frac{5}{6}\right)^{(n-1)}, n \geq 1$

□ $\mu = 6$

Example: throwing the dice until a 6

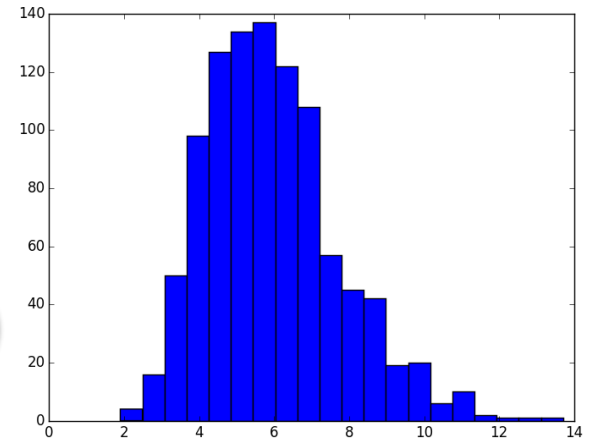
41

Number of samples: 1000

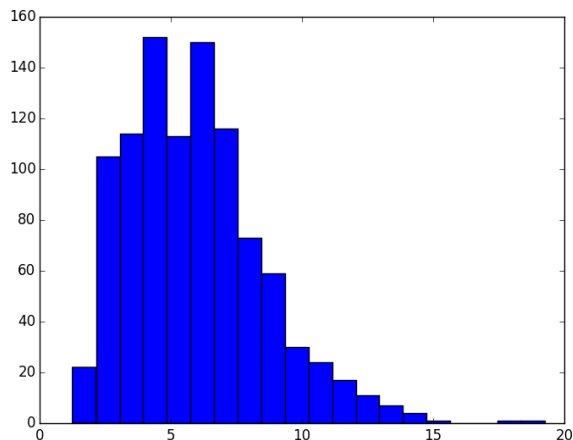


Sample size

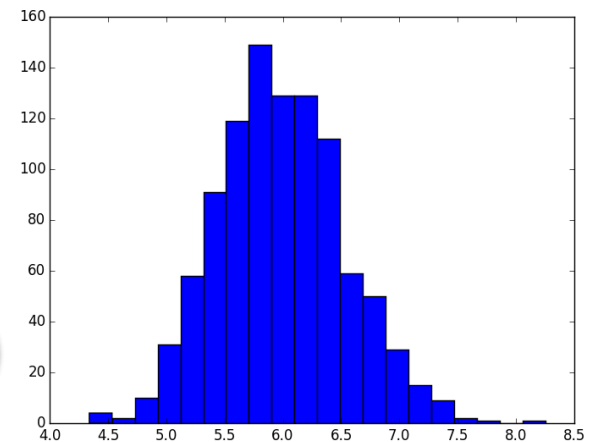
1



10



4



100

Binomial distribution

$$b(k; n, p) = \binom{n}{k} p^k (1-p)^{(n-k)}$$

Population: all Bernoulli trials with probability p .

Sample: n such trials

Example: Throwing a dice n times, counting the number of 6-s (success)

- Number of successes: X
- Random variable over all series of n trials
- **Binomial distribution** (binomisk fordeling): $B(n, p)$
- $E(X) = np$
- $\text{Var}(X) = np(1-p)$
- $\sigma_X = \sqrt{np(1-p)}$
- Approximated by $N(np, \sqrt{np(1-p)})$ for large n

- Proportion of success: $\hat{p} = X/n$
- $E(\hat{p}) = E(X/n) = np/n = p$
- $\text{Var}(\hat{p}) = \sigma_X^2 / n^2 = np(1-p) / n^2 = p(1-p) / n$
- $\sigma_{\hat{p}} = \sqrt{\frac{p(1-p)}{n}} = \frac{\sigma_Y}{\sqrt{n}}$
- Approximated by $N(p, \sqrt{p(1-p)/n})$ for large n

Rule of thumb:

$$np > 10 \text{ and}$$

$$n(1-p) > 10$$