

# Statistical methods in NLP

## Estimation



**UNIVERSITY OF  
GOTHENBURG**

Richard Johansson

January 28, 2016



















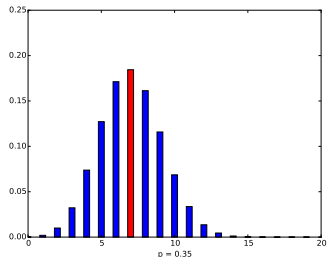
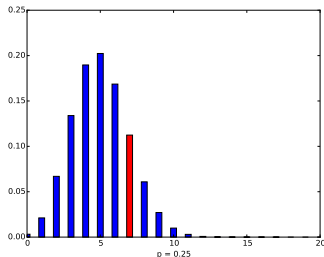






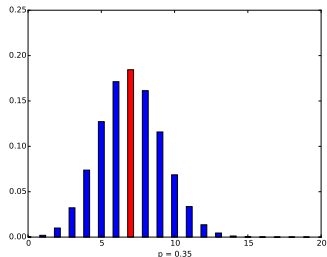
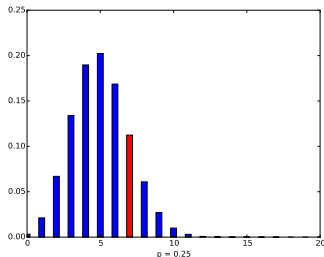
# maximizing the likelihood

- ▶ we classified 20 documents with 7 errors; what's the MLE of the error rate?



# maximizing the likelihood

- ▶ we classified 20 documents with 7 errors; what's the MLE of the error rate?



- ▶ it can be shown that the value of  $p$  that gives us the maximum of  $L(p)$  is

$$p_{MLE} = \frac{x}{n}$$











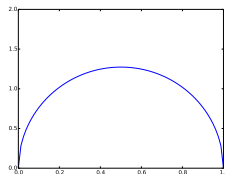






# MAP with a Dirichlet prior

- ▶ the **Dirichlet** distribution is often used as a prior in MAP estimates
- ▶ assume that we pick  $n$  words randomly from a corpus
  - ▶ the words come from a vocabulary with the size  $V$
  - ▶ we saw the word *armchair*  $x$  times out of  $n$
- ▶ with a Dirichlet prior with a concentration parameter  $\alpha$ , the MAP estimate is



$$p_{MAP} = \frac{x + (\alpha - 1)}{n + V \cdot (\alpha - 1)}$$

- ▶ for instance, with  $\alpha = 2$ , we get

$$p_{MAP} = \frac{x + 1}{n + V}$$



interval estimates

- ▶ if we get some estimate by ML, can we say something about how reliable that estimate is?
- ▶ a **confidence interval** for the parameter  $p$  with significance value  $\alpha$  is an interval  $[p_{low}, p_{high}]$  so that

$$P(p_{low} \leq p \leq p_{high}) \geq \alpha$$

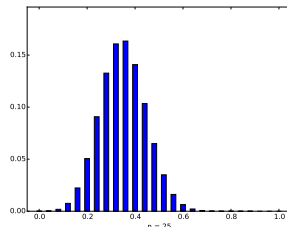
- ▶ for instance: with 95% probability, the error rate of the spam filter is in the interval  $[0.05, 0.08]$ 
  - ▶ that is: it is between 0.05 and 0.08



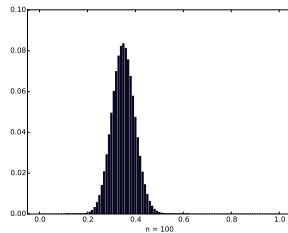
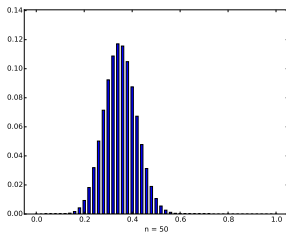
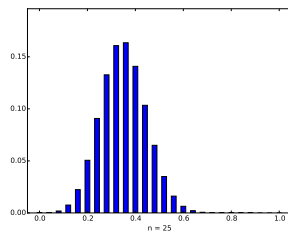
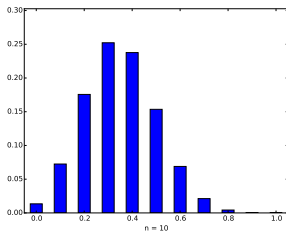


# the distribution of our estimator

- ▶ our ML or MAP estimator applied to randomly selected samples is a random variable with a distribution
- ▶ this distribution depends on the **sample size**
  - ▶ large sample  $\rightarrow$  more concentrated distribution
- ▶ we will reason about this distribution to show how a confidence interval can be found

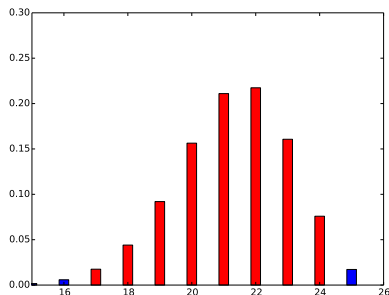


# estimator distribution and sample size ( $p = 0.35$ )



## the distribution of ML estimates of heads probabilities

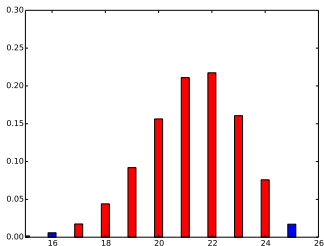
- ▶ if the true  $p$  value is 0.85 and we toss the coin 25 times, what results do we get?



- ▶ 95% of the experiments give a result between 17 and 24
- ▶ ...so 95% of the estimates will be between 0.68 and 0.96

# confidence interval for the ML estimation of a probability

- ▶ assume we toss the coin  $n$  times, with  $x$  'heads'
- ▶ cookbook recipe for computing an approximate 95% confidence interval  $[p_{low}, p_{high}]$ :
  - ▶ first estimate  $p^* = x/n$  as usual
  - ▶ to compute the lower bound  $p_{low}$ :
    1. let  $X$  be a binomially distributed r.v. with parameters  $n, p^*$
    2. find  $x_{low} =$  the 2.5% percentile of  $X$
    3.  $p_{low} = x_{low}/n$
  - ▶ for the upper bound  $p_{high}$ , use the 97.5% percentile instead



## in Scipy

- ▶ assume we got 'heads'  $x$  times out of  $n$
- ▶ recall that we use ppf to get the percentiles!

$$p\_est = x / n$$

```
rv = scipy.stats.binom(n, p_est)
```

```
p_low = rv.ppf(0.025) / n
```

```
p_high = rv.ppf(0.975) / n
```

```
print(p_low, p_high)
```

example: political polling

- ▶ I ask 38 randomly selected Gothenburgers about whether they support the congestion tax in Gothenburg
- ▶ 22 of them say yes
- ▶ an approximate 95% confidence interval for the popularity of the tax is  $0.421 - 0.737$

```
number_yes = 22
total_number = 38
p_est = number_yes / total_number

rv = scipy.stats.binom(total_number, p_est)
p_low = rv.ppf(0.025) / total_number
p_high = rv.ppf(0.975) / total_number

print(p_low, p_high)
```

don't forget your common sense

- ▶ I ask 14 MLT students about whether they support the congestion tax, 11 of them say yes
- ▶ will I get a good estimate?
- ▶ in NLP: the “WSJ fallacy”



