# Topics in Natural Language Processing

Shay Cohen

Institute for Language, Cognition and Computation

University of Edinburgh

Lecture 3

# Administrativia

There is an online forum available now on Piazza. You are able to:

- Ask your peers questions about the material

- Ask your peers general questions about NLP

# Last class

Maximum likelihood estimation:

$\theta \in \Theta \quad p(w \mid \theta) \qquad \underbrace{w_1 \ldots w_n}_{data}$

$$L(\theta, w_1 \ldots w_n) = \frac{1}{n} \sum_{i=1}^{n} \log p(w_i \mid \theta)$$

$$\theta^* = \underset{\theta}{argmax} \; L(\theta, w_1 \ldots w_n)$$

count of "ugh"

$$\theta^* = \frac{a}{n}$$

$$\widehat{\Theta} = [0, 1]$$

$$\Omega = \{1, \ldots, d\} \qquad \theta_i = \frac{count(i \; in \; w_1 \ldots w_n)}{n}$$

$$\Theta \in \mathbb{R}^d \qquad \text{Multinomial distribution} \quad MLE$$

# Likelihood maximisation

By doing maximum likelihood maximisation we:

- Choose the parameters that make the data most probable,

  or, from an information-theoretic perspective:

- Choose the parameters that make the encoding of the data most succinct (bit-wise),

  in other words, we

- Minimize the cross-entropy between the empirical distribution and the model we choose.

# Today's class

- The Bayesian paradigm

- If there is time: structure in NLP - or "what is our $\Omega$?"

# Some history

- History: 1700s. Seminal ideas due to Thomas Bayes and Pierre-Simon Laplace

# Bayes' rule

What is Bayes' rule? $p(X=x, Y=y)$

$p(X=x \mid Y=y)$

$p(Y=y \mid X=x)$

$$\frac{p(X=x \mid Y=y)\, p(Y=y)}{p(X=x)} = p(Y=y \mid X=x)$$

Reminder: What does Statistics do? Invert the relationship between model and data.

Bayes' rule does the same with random variables.

# Bayes' rule

What is Bayes' rule?

Reminder: What does Statistics do? Invert the relationship between model and data.

Bayes' rule does the same with random variables.

What if our model parameters were one random variable and our data were another random variable?

# Prior beliefs about models

We have a parameter space $\Theta$ and prior beliefs $p(\theta)$.

Our $\theta$ is now a random variable.

From the chain rule: $p(w, \theta) = p(\theta)p(w|\theta)$

"$X$" "$Y$"

"prior"

Prior   likelihood

$$p(\theta|w) = \frac{p(\theta) p(w|\theta)}{p(w)}$$

"posterior"

Bayes rule

# Posterior inference

$$p(\theta \mid w) = \frac{p(w \mid \theta)\,p(\theta)}{p(w)}$$

basic posterior inference

$$p(w) =$$

$$\int_\theta p(\theta\mid w)\,d\theta = 1 \qquad \int \frac{p(w\mid\theta)\,p(\theta)}{p(w)}\,d\theta = 1$$

$$\frac{1}{p(w)} \int_\theta p(w\mid\theta)\,p(\theta)\,d\theta = 1$$

$$p(w) = \int_\theta p(w\mid\theta)\,p(\theta)\,d\theta = 1$$

evidence, normalization constant

# Priors

Our prior beliefs are considered in inference. There is no "correct" prior.

Is that a good or bad thing?

# Priors

Our prior beliefs are considered in inference. There is no "correct" prior.

Is that a good or bad thing?

- Frequentists: probability is the frequency of an event

- Bayesians: probability denotes the state of our knowledge about an event
    - Subjectivists: probability is a personal belief
    - Objectivists: minimise human's influence on decision making
- In practice: NLP use of Bayesian theory is largely driven by computation

# Back to pre-historic languages



Language with two words: "argh" and "blah"

Our $\Omega$ is $\{\mathrm{argh}, \mathrm{blah}\}$.

Our $\Theta$ is $[0, 1]$.

Define $I(w) = 1$ if $w = \mathrm{argh}$ and $0$ if $w = \mathrm{blah}$.

Then, $p(w|\theta) = \theta^{I(w)}(1-\theta)^{(1-I(w))}$.
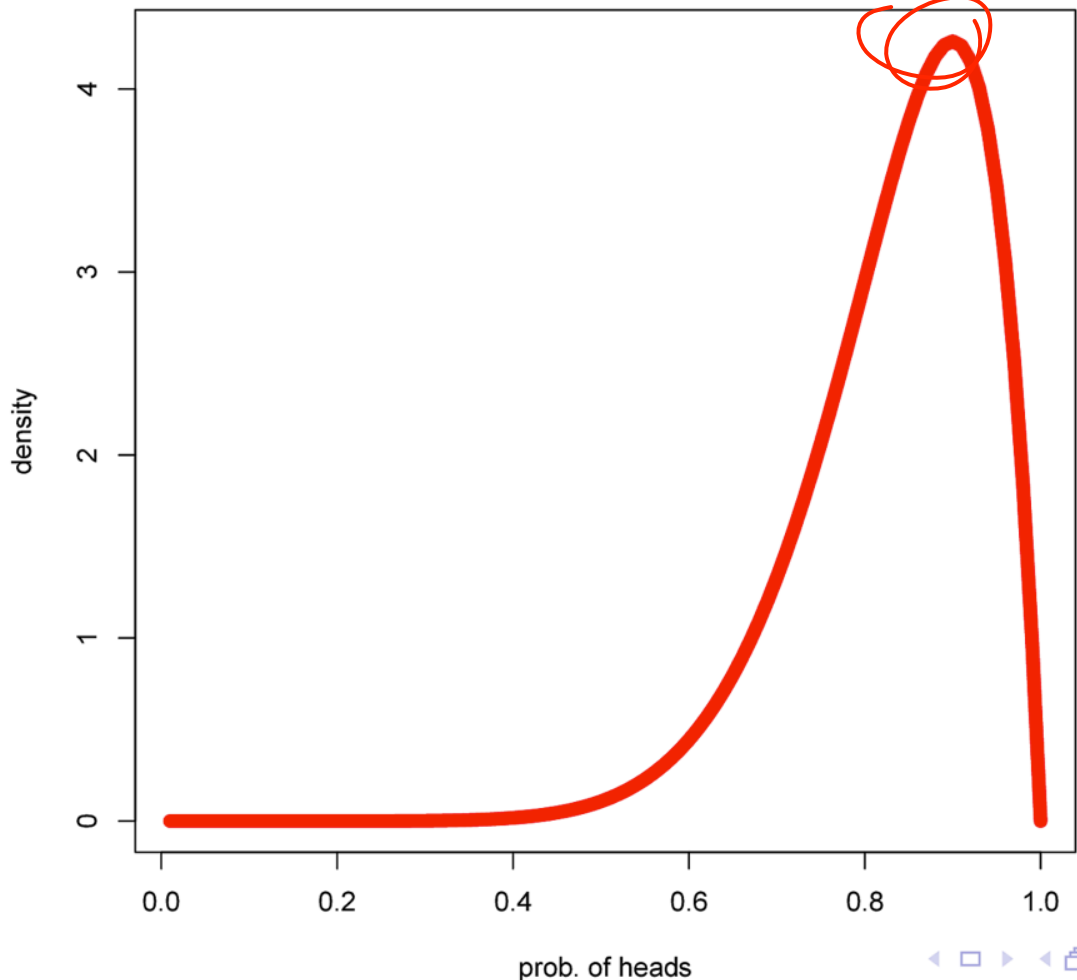
# Uniform prior, 0.7 prob. for argh

# Posterior with 10 datapoints, truth is 0.7 prob. for argh



$p(\theta | w_1 \cdots w_{10})$

density

prob. of heads
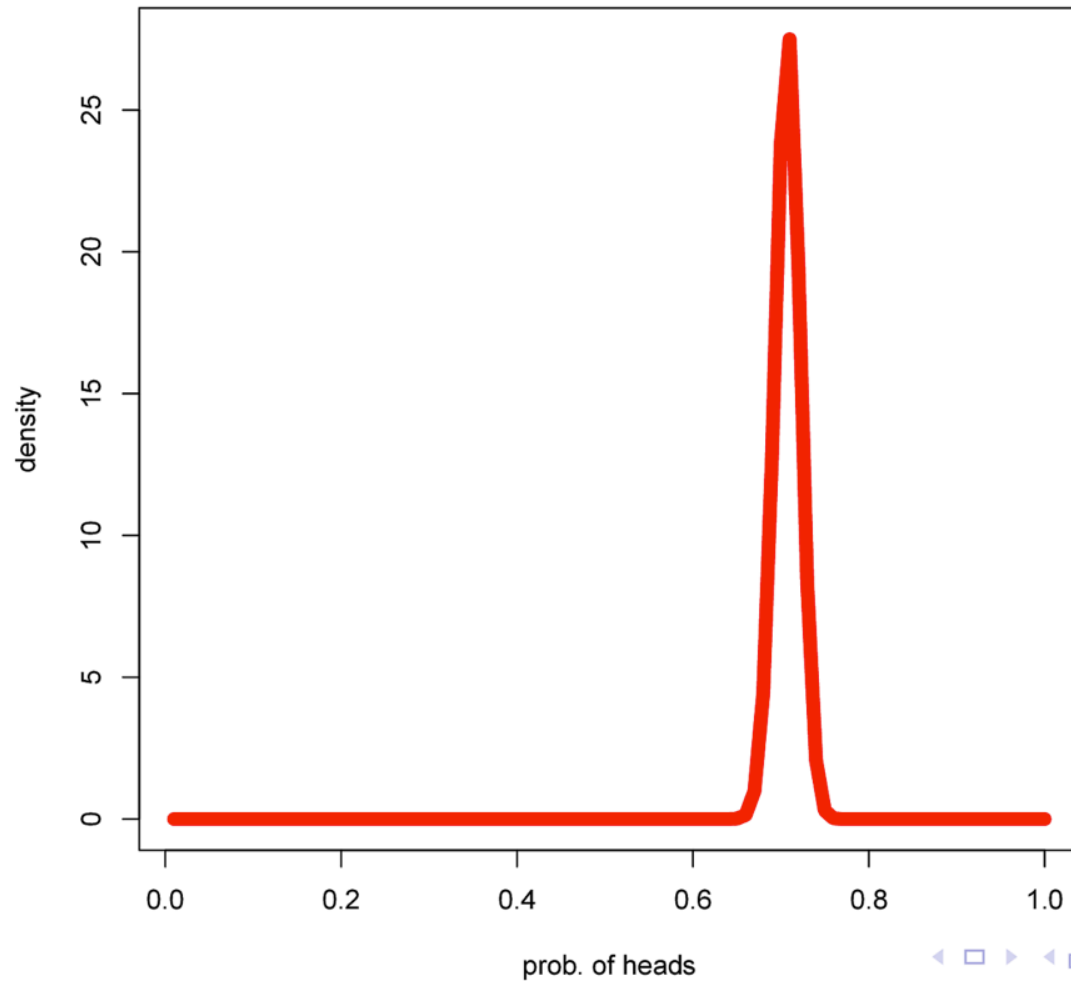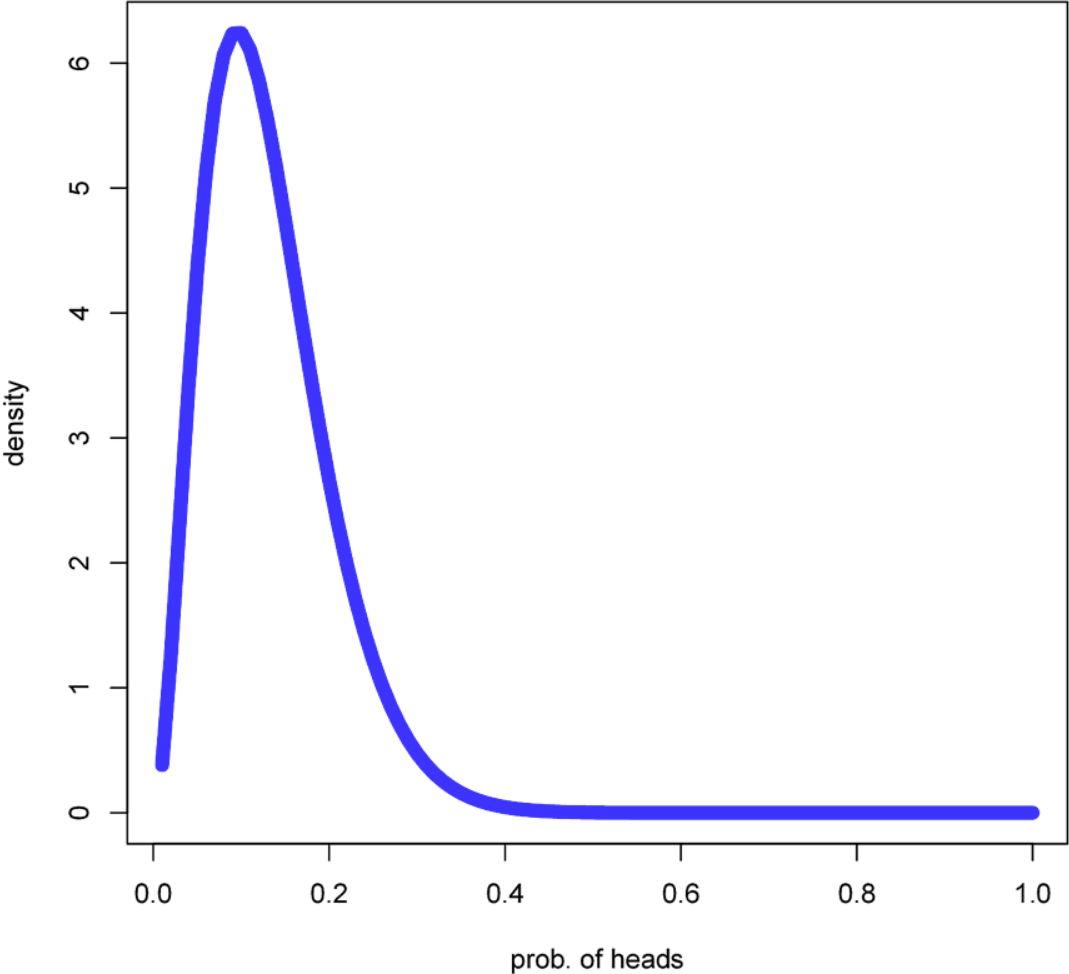
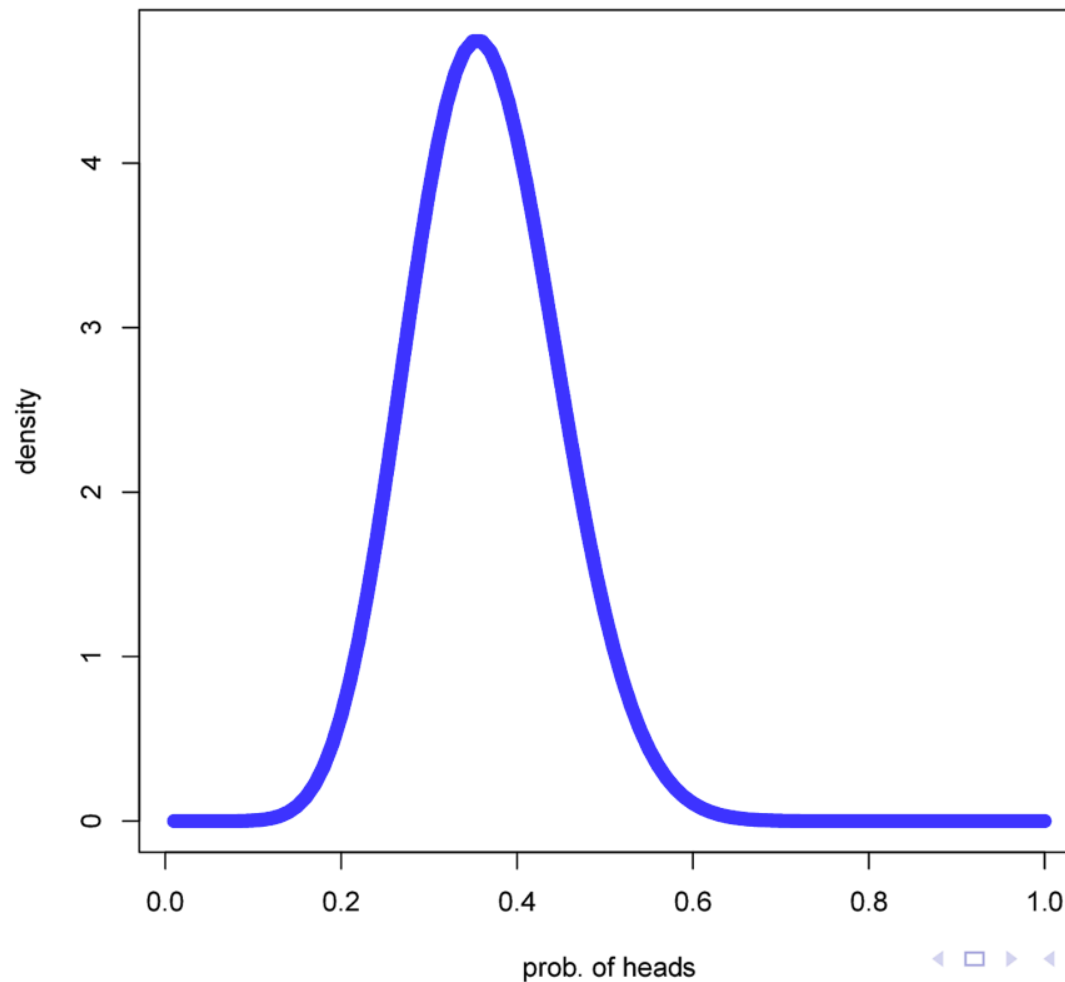# Posterior with 100 datapoints, truth is 0.7 prob. for argh

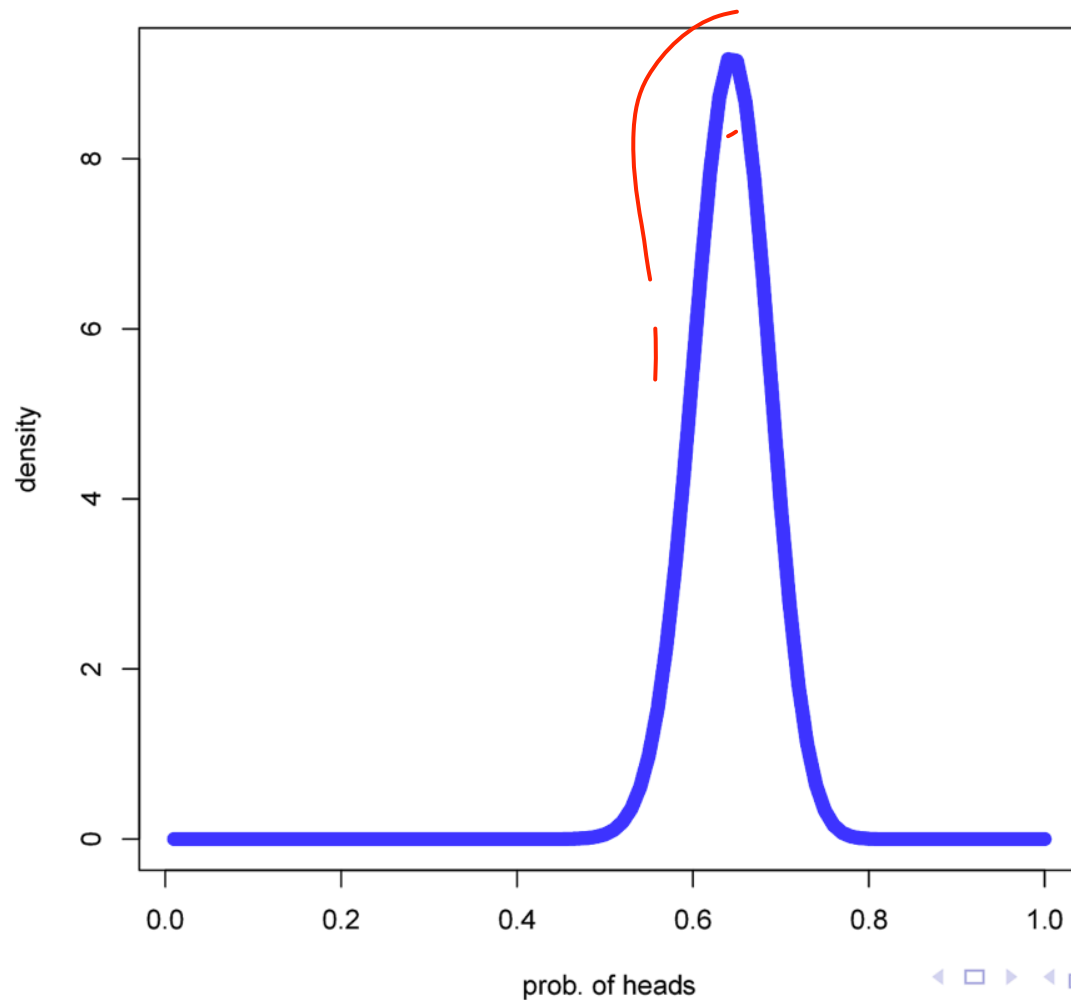# Posterior with 1000 datapoints, truth is 0.7 prob. for argh

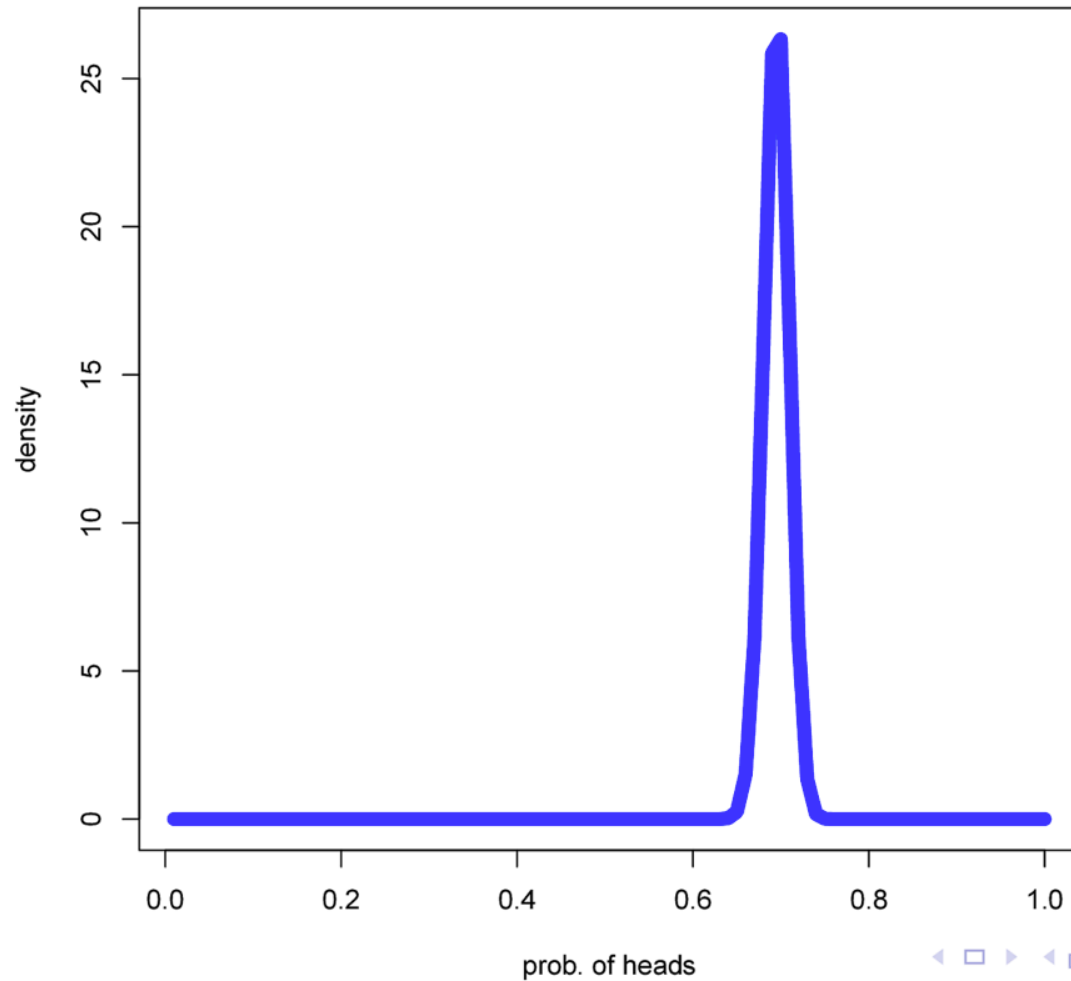# Non-uniform prior, truth is 0.7 prob. for argh

# Posterior with 10 datapoints, truth is 0.7 prob. for argh

# Posterior with 100 datapoints, truth is 0.7 prob. for argh

# Posterior with 1000 datapoints, truth is 0.7 prob. for argh

# Priors for binary outcomes

Beta dist

$\hookrightarrow p(\theta) \propto \theta^{\alpha}(1-\theta)^{\beta}$ $\qquad\qquad p(w|\theta) = \theta^{I(w)}(1-\theta)^{(1-I(w))}$

if $\alpha = \beta = 1$ $\quad p(\theta) \propto 1$, which means uniform

## What is the posterior?

$D = \{w_1 \cdots w_n\}$

$$\underline{p(\theta|w_1 \cdots w_n)} = \frac{\overline{p(w_1 \cdots w_n|\theta) \, p(\theta)}}{p(w_1 \cdots w_n)} = \frac{\left( \prod_{i=1}^{n} p(w_i|\theta) \right) p(\theta)}{p(w_1 \cdots w_n)}$$

likelihood

$$= \left( \prod_{i=1}^{n} \theta^{I(w_i)} (1-\theta)^{1-I(w_i)} \right) \times \theta^{\alpha}(1-\theta)^{\beta} \Big/ p(w_1 \cdots w_n)$$

$$= \theta^{\sum_{i=1}^{n} I(w_1)} (1-\theta)^{\sum_{i=1}^{n} 1-I(w_i)} \times \theta^{\alpha}(1-\theta)^{\beta} \Big/ p(w_1 \cdots w_n)$$

$$= \theta^{\alpha + \sum I(w_i)} (1-\theta)^{\beta + \sum 1-I(w_i)} \Big/ p(w_1 \cdots w_n)$$

# Maximum a posteriori estimate (MAP)

"Bayesian estimation": find $\theta^*$ that maximises the posterior:

$$\theta^* = \underset{\theta}{\arg\max}\ p(\theta \mid w_1 \dots w_n) = \underset{\theta}{\arg\max}\ \theta^{a+\alpha}(1-\theta)^{b+\beta}$$

$$= \underset{\theta}{\arg\max}\ (a+\alpha)\log\theta + (b+\beta)\log(1-\theta)$$

$$\theta^* = \frac{a+\alpha}{\underbrace{a+b}_{n}+\alpha+\beta} = \frac{a+\alpha}{n+(\alpha+\beta)}$$

$$a = \sum \mathbb{I}(w_i) \qquad b = n - \sum \mathbb{I}(w_i)$$

# MAP and posteriors

In general,

- Priors are especially important when the amount of data is small

- As there is more data, the prior becomes less influential on the posterior

- Under some mild conditions, the posterior is a distribution concentrated around the MLE

# Next class

- Conjugacy of Bayesian priors to the likelihood

- Structure in NLP