

# Topics in Natural Language Processing

Shay Cohen

Institute for Language, Cognition and Computation

University of Edinburgh

Lecture 3

# Administrativa

---

Do we need an online class forum? You would be able to:

- Ask your peers questions about the material
- Look for team members for presentations, etc.
- Ask your peers general questions about NLP

# Last class

Maximum likelihood estimation:

$$p(w|\theta) \quad w_1 \dots w_n$$

data

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

$$\theta^* = \arg \max_{\theta} L(\theta, w_1, \dots, w_n)$$

$$\Theta = [a, 1]$$

$$\theta^* = \frac{a}{n} \leftarrow \text{count of "aigh"}$$

$$\Omega = \{1, \dots, d\}$$

$$\theta_i^* = \frac{\text{count}(i \text{ in } w_1, \dots, w_n)}{n}$$

Multinomial

Distribution

MLE

# 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?  $p(X=x, Y=y)$  given

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

proof:

$$\begin{aligned} p(X=x|Y=y) p(Y=y) &= \\ &= p(Y=y|X=x) p(X=x) \\ &= p(X=x, Y=y) \end{aligned}$$

chain rule

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:  $p(w, \theta) = p(\theta) p(w|\theta)$

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

↑ posterior      ↑ Bayes' rule

Note that  $\Theta$  is continuous, therefore we need  $\int_{\Theta} p(\theta) d\theta = 1$ .

This replaces sum-to-1 constraint



# Posterior inference

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

basic posterior inference

$$p(w) = \int_{\theta} p(w | \theta) p(\theta) d\theta$$

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

$$\Rightarrow p(w) = \int_{\theta} p(w | \theta) p(\theta) d\theta$$

① choose  $p(\theta)$

What do we need to do next?

② compute  $p(w)$

# Priors

---

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

Is that a good or bad thing? *Neither?*

# 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  $\{\text{argh}, \text{blah}\}$ .

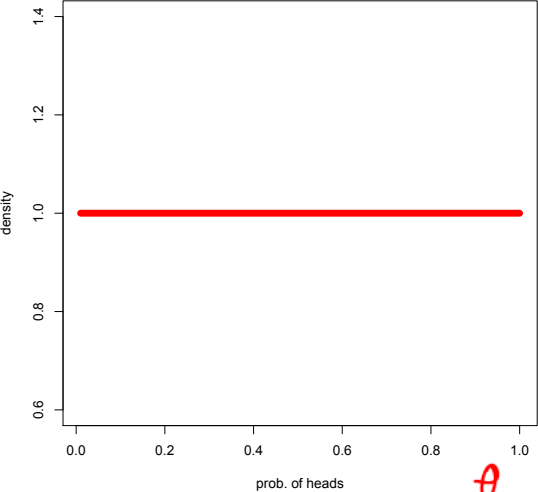
Our  $\Theta$  is  $[0, 1]$ .

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

Then,  $p(w|\theta) = \theta^{I(w)}(1 - \theta)^{(1-I(w))}$ .

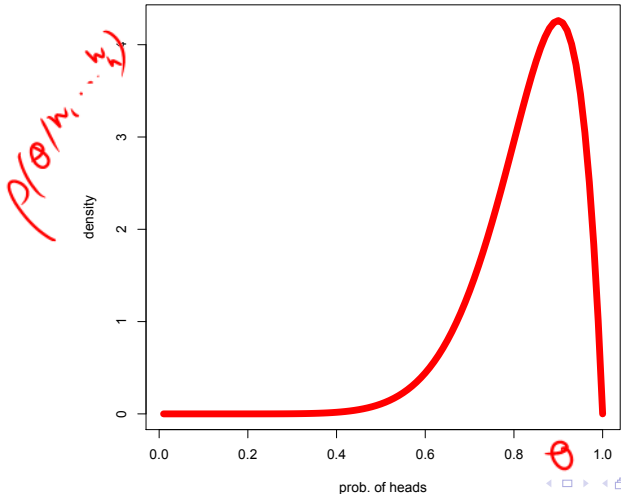
# Uniform prior, 0.7 prob. for argh

$p(\theta)$



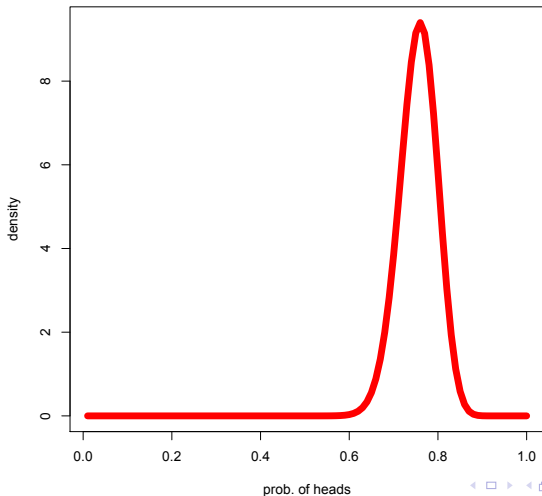
$\theta$

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



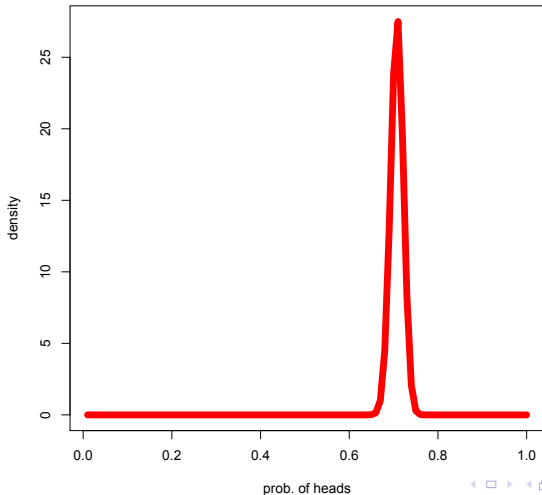
# Posterior with 100 datapoints, truth is 0.7 prob. for arg

---



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

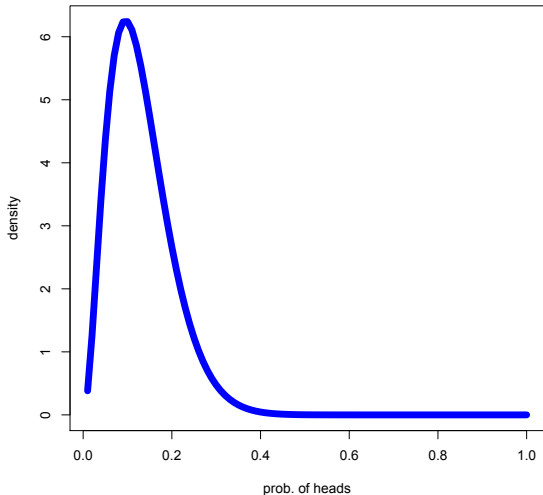
---





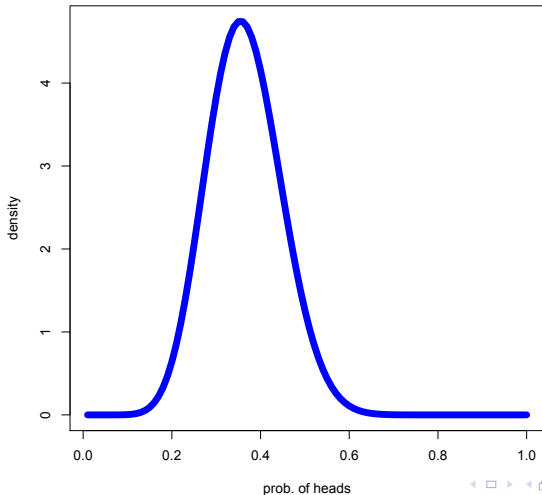
# Non-uniform prior, truth is 0.7 prob. for argh

---



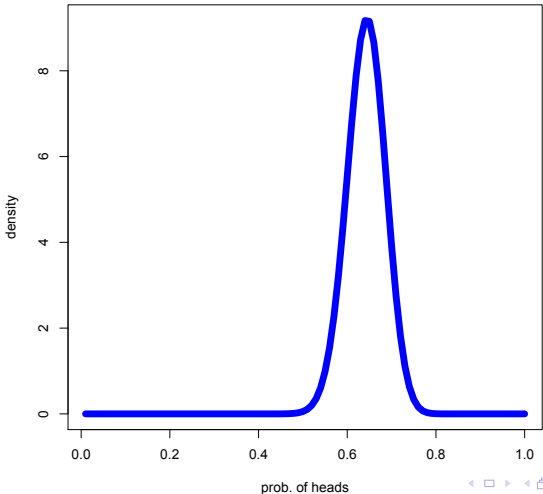
# Posterior with 10 datapoints, truth is 0.7 prob. for arg

---



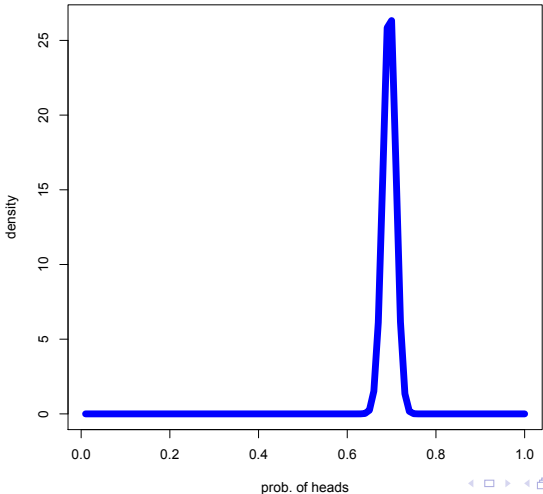
# Posterior with 100 datapoints, truth is 0.7 prob. for argh

---



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

---



# Priors for binary outcomes

$$p(\theta) \propto \theta^\alpha (1 - \theta)^\beta$$

$$p(w|\theta) = \theta^{I(w)} (1 - \theta)^{(1-I(w))}$$

What is the posterior?  $\mathcal{D} = \{w_1, \dots, w_n\}$

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

$$\left[ \prod_{i=1}^n \theta^{I(w_i)} (1 - \theta)^{1 - I(w_i)} \right] \times \theta^\alpha (1 - \theta)^\beta \Big/ p(w_1, \dots, w_n)$$
$$= \theta^{\sum_{i=1}^n I(w_i) + \alpha} (1 - \theta)^{(n - \sum_{i=1}^n I(w_i)) + \beta} \Big/ p(w_1, \dots, w_n)$$

$$= \left( \theta^{a + \alpha} (1 - \theta)^{b + \beta} \right) \Big/ p(w_1, \dots, w_n)$$

# Maximum a posteriori estimate (MAP)

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

$$\theta^* = \underset{\theta}{\operatorname{argmax}} p(\theta | w_1, \dots, w_n) = \underset{\theta}{\operatorname{argmax}} \frac{\theta^{a+\alpha} (1-\theta)^{b+\beta}}{p(w_1, \dots, w_n)}$$

$$= \underset{\theta}{\operatorname{argmax}} \theta^{a+\alpha} (1-\theta)^{b+\beta} = \underset{\theta}{\operatorname{argmax}} (a+\alpha) \log \theta + (b+\beta) \log(1-\theta)$$

↑  
take log

$$\theta^* = \frac{a+\alpha}{a+\alpha+b+\beta}$$

MAP estimate

① smoothing!

$$\textcircled{2} \theta^* \approx \frac{a}{a+b} \quad \text{as } n \rightarrow \infty$$

# 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