

Topics in Natural Language Processing

Shay Cohen

Institute for Language, Cognition and Computation

University of Edinburgh

Lecture 4

3

Administrativa

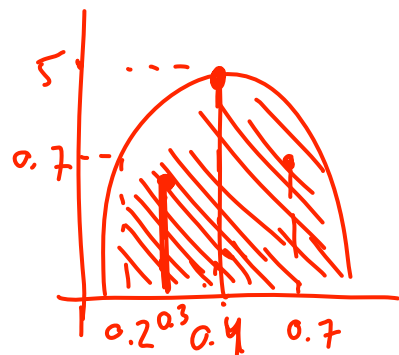
- You should have been added to the Piazza forum
- Today is the last day to send in the papers for presentations. I will send an email over the weekend with assignments.

$p(\theta)$ - density like a prob. dist. over
a continuous space, $[0, 1]$

we want from a density

$$p(\theta) \geq 0$$

$$\int_{\theta} p(\theta) d\theta \cong 1$$



if this is

$$p(\theta | w_1, \dots, w_n)$$

we just
choose

$$\hat{\theta} = 0.4$$

Last class

Bayesian inference: $p(\theta)$ prior $p(w_1, \dots, w_n | \theta) = \prod_{i=1}^n p(w_i | \theta)$

$$p(\theta | w_1, \dots, w_n) \stackrel{\text{Bays}}{=} \frac{p(\theta) p(w_1, \dots, w_n | \theta)}{p(w_1, \dots, w_n)}$$

$$\Rightarrow p(w_1, \dots, w_n) = \int_{\theta} p(w_1, \dots, w_n | \theta) p(\theta) d\theta$$

MAP - maximum a posteriori

$$\theta^* = \underset{\theta}{\operatorname{argmax}} p(\theta | w_1, \dots, w_n)$$

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

Conjugacy of prior and likelihood

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

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

Prior is "hyperparametrised". What is the posterior?

$p(\theta | w_1, \dots, w_n) \propto \theta^{\alpha + a} (1 - \theta)^{\beta + b}$ ← Beta α

$p(\theta | w_1, \dots, w_n) = \frac{\theta^{\alpha + a} (1 - \theta)^{\beta + b}}{\int_{\theta'} \theta'^{\alpha + a} (1 - \theta')^{\beta + b} d\theta'}$

posterior

"new α " $\alpha + a$
 "new β " $\beta + b$

Definition of Conjugacy

(α, β)
↑

(σ, τ)
↑

Let P be a set of priors hyperparametised by a set \mathcal{Q} , for a parameter space Θ . Therefore, each $p \in P$ is a probability distribution $p(\theta | \alpha)$. Let M be a model over Ω such that each $p \in M$ is a probability distribution $p(w | \theta)$. We say, P is conjugate to M , if for any choice of $\alpha \in \mathcal{Q}$ and data w_1, \dots, w_n it holds that $p(\theta | w_1, \dots, w_n, \alpha) \in P$.

Definition of Conjugacy

Let P be a set of priors hyperparametised by a set \mathcal{A} , for a parameter space Θ . Therefore, each $p \in P$ is a probability distribution $p(\theta | \alpha)$. Let M be a model over Ω such that each $p \in M$ is a probability distribution $p(w | \theta)$. We say, P is **conjugate to** M , if for any choice of $\alpha \in \mathcal{A}$ and data w_1, \dots, w_n it holds that $p(\theta | w_1, \dots, w_n, \alpha) \in P$.

Previous example (argh-blah example):

$$M = \{ p(w | \theta) \mid \theta \in [0, 1] \}, \quad p(w | \theta) = \begin{cases} \theta & \text{argh} \\ 1 - \theta & \text{blah} \end{cases}$$

$$P = \{ p(\theta) = \theta^\alpha (1 - \theta)^\beta \mid \alpha \geq 0, \beta \geq 0 \}$$

Posterior new hyperparameters:

$$\theta^{\frac{\alpha + a}{2}} (1 - \theta)^{\frac{\beta + b}{2}}$$

Conjugacy – always useful?

Trivial non-useful example of conjugacy

$P = \{ \text{set } \subset \mathcal{P} \text{ all distributions over } \theta \}$

$M = \text{some model}$

Conjugacy – always useful?

Another trivial non-useful example of conjugacy

choose some θ $\theta = 0.7$

$$P = \left\{ p(\theta) \text{ s.t. } \begin{array}{ll} p(\theta) = 1 & \text{if } \theta = 0.7 \\ 0 & \text{o/w} \end{array} \right\}$$

$$p(\theta/w) = \frac{p(\theta)p(w|\theta)}{p(w)} = \begin{cases} 1 & \theta = 0.7 \\ 0 & \text{o/w} \end{cases}$$

Conjugacy: summary

Conjugacy is useful when:

- The prior is not too poor
- It is easy to calculate the posterior hyperparameters

$$(\alpha, \beta)$$

Minimum Description Length and MAP

What is $-\log_2 p(\theta|w_1, \dots, w_n)$?

Minimum Description Length and MAP

What is $-\log_2 p(\theta|w_1, \dots, w_n)$?

What is $-\log_2 p(\theta)$?

Minimum Description Length and MAP

What is $-\log_2 p(\theta | w_1, \dots, w_n)$?

bits that are required using a "good" code to encode θ given w_1, \dots, w_n

What is $-\log_2 p(\theta)$?

bits that are required using a good code to encode θ a priori

What is $-\log_2 p(w_1, \dots, w_n | \theta)$?

bits that are required to encode the data if we think θ generated it

$$\theta^* = \underset{\theta}{\operatorname{argmax}} p(\theta | w_1, \dots, w_n)$$

Minimum Description Length and MAP

What is $-\log_2 p(\theta|w_1, \dots, w_n)$?

What is $-\log_2 p(\theta)$?

What is $-\log_2 p(w_1, \dots, w_n|\theta)$?

MAP: $\theta^* = \arg \max_{\theta} \log_2 p(\theta) + \log_2 p(w_1, \dots, w_n|\theta)$

$$\begin{aligned}\theta^* &= \arg \max_{\theta} p(\theta|w_1, \dots, w_n) = \\ &= \arg \max_{\theta} \frac{p(\theta) p(w_1, \dots, w_n|\theta)}{p(w_1, \dots, w_n)} \\ &= \arg \max_{\theta} \log_2 p(\theta) + \log_2 p(w_1, \dots, w_n|\theta) \\ &= \arg \min_{\theta} \underbrace{-\log_2 p(\theta) - \log_2 p(w_1, \dots, w_n|\theta)}\end{aligned}$$

Encoding θ^* requires separately:

- Encoding the hypothesis according to the prior
- Encoding the data according to the hypothesis

That's the "minimum description length" criterion

$$\theta^* = \underset{\theta}{\operatorname{arg\,max}} \underbrace{\log p(\theta)} + \sum \log p(\omega_i | \theta)$$

$$\theta_{\text{MLE}} = \underset{\theta}{\operatorname{arg\,max}} \sum_i \log p(\omega_i | \theta)$$

Regularization

$$\theta_{\text{MLE}} = \underset{\theta}{\operatorname{arg\,max}} \sum_i \log p(\omega_i | \theta) - \lambda \sum_i \theta_i^2$$

Summary

Bayesian analysis:

- Only uses Bayes' rule to do inference
- Posterior is a *distribution* over parameters
- Can summarise the posterior, e.g. MAP, to get a point estimate
- Need to be careful about choice of prior
- Especially important with small amounts of data
- MAP has a connection to minimum description length (MDL)