# TALENT COURSE I I LEARNING FROM DATA: BAYESIAN METHODS AND MACHINE LEARNING

Lecture 10: Choosing a prior

Daniel Phillips
Ohio University
TU Darmstadt
ExtreMe Matter Institute









TALENT Course II is possible thanks to funding from the STFC

Thomas Bayes (1701?-1761)



$$\operatorname{pr}(A \mid B, I) = \frac{\operatorname{pr}(B \mid A, I)\operatorname{pr}(A \mid I)}{\operatorname{pr}(B \mid I)}$$

Thomas Bayes (1701?-1761)



http://www.bayesian-inference.com

$$\operatorname{pr}(A \mid B, I) = \frac{\operatorname{pr}(B \mid A, I)\operatorname{pr}(A \mid I)}{\operatorname{pr}(B \mid I)}$$

Thomas Bayes (1701?-1761)



$$\operatorname{pr}(A \mid B, I) = \frac{\operatorname{pr}(B \mid A, I)\operatorname{pr}(A \mid I)}{\operatorname{pr}(B \mid I)}$$

http://www.bayesian-inference.com

$$pr(\text{model} | \text{data}, I) = \frac{pr(\text{data} | \text{model}, I)pr(\text{model} | I)}{pr(\text{data} | I)}$$

Thomas Bayes (1701?-1761)



$$\operatorname{pr}(A \mid B, I) = \frac{\operatorname{pr}(B \mid A, I)\operatorname{pr}(A \mid I)}{\operatorname{pr}(B \mid I)}$$

http://www.bayesian-inference.com

$$pr(\text{model} | \text{data}, I) = \frac{pr(\text{data} | \text{model}, I)pr(\text{model} | I)}{pr(\text{data} | I)}$$

Posterior

Thomas Bayes (1701?-1761)



http://www.bayesian-inference.com

$$\operatorname{pr}(A \mid B, I) = \frac{\operatorname{pr}(B \mid A, I)\operatorname{pr}(A \mid I)}{\operatorname{pr}(B \mid I)}$$

Likelihood

$$pr(\text{model} | \text{data}, I) = \frac{pr(\text{data} | \text{model}, I)pr(\text{model} | I)}{pr(\text{data} | I)}$$

Posterior

Thomas Bayes (1701?-1761)



http://www.bayesian-inference.com

$$\operatorname{pr}(A \mid B, I) = \frac{\operatorname{pr}(B \mid A, I)\operatorname{pr}(A \mid I)}{\operatorname{pr}(B \mid I)}$$

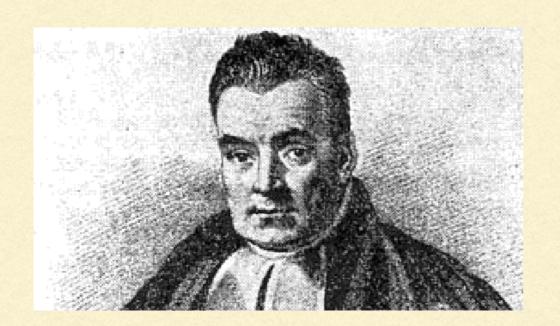
Likelihood

Prior

$$pr(\text{model} | \text{data}, I) = \frac{pr(\text{data} | \text{model}, I)pr(\text{model} | I)}{pr(\text{data} | I)}$$

Posterior

Thomas Bayes (1701?-1761)



http://www.bayesian-inference.com

$$\operatorname{pr}(A \mid B, I) = \frac{\operatorname{pr}(B \mid A, I)\operatorname{pr}(A \mid I)}{\operatorname{pr}(B \mid I)}$$

Likelihood

Prior

$$pr(\text{model} | \text{data}, I) = \frac{pr(\text{data} | \text{model}, I)pr(\text{model} | I)}{pr(\text{data} | I)}$$

Posterior

Normalization

Thomas Bayes (1701?-1761)



http://www.bayesian-inference.com

$$\operatorname{pr}(A \mid B, I) = \frac{\operatorname{pr}(B \mid A, I)\operatorname{pr}(A \mid I)}{\operatorname{pr}(B \mid I)}$$

Likelihood

Prior

$$pr(model | data, I) = \frac{pr(data | model, I)pr(model | I)}{pr(data | I)}$$

$$pr(data | I)$$

Probability as degree of belief cf. frequentist view

Consider a six-sided dice

- Consider a six-sided dice
- How do we assign  $pr(X_i|I)$ ? (i=1, 2, 3, 4, 5, 6)

- Consider a six-sided dice
- How do we assign  $pr(X_i|I)$ ? (i=1, 2, 3, 4, 5, 6)
- We do know

$$\sum_{i} \operatorname{pr}(X_i | I) = 1$$

- Consider a six-sided dice
- How do we assign  $pr(X_i|I)$ ? (i=1, 2, 3, 4, 5, 6)
- We do know

$$\sum_{i} \operatorname{pr}(X_i | I) = 1$$

■ Invariance under labeling $\Rightarrow$ pr( $X_i|I$ )=1/6.

provided I says nothing that breaks the symmetry

#### INDIFFERENCE PRIORS: CONTINUOUS

#### INDIFFERENCE PRIORS: CONTINUOUS

#### Location

$$\operatorname{pr}(x \mid I) dx = \operatorname{pr}(x + x_0 \mid I) dx$$

• Invariance under origin position⇒pr(x|I)=constant
provided I says nothing that breaks the symmetry

#### INDIFFERENCE PRIORS: CONTINUOUS

#### Location

$$\operatorname{pr}(x \mid I) dx = \operatorname{pr}(x + x_0 \mid I) dx$$

- Invariance under origin position⇒pr(x|I)=constant
   provided I says nothing that breaks the symmetry
  - Scale

$$\operatorname{pr}(\lambda | I) d\lambda = \operatorname{pr}(\beta \lambda | I) d(\beta \lambda)$$

- Invariance under unit choice  $\Rightarrow$  pr( $\lambda \mid I$ )  $\propto 1/\lambda$  provided I says nothing that breaks the symmetry
- "Jeffreys prior": uniform prior on  $log(\lambda)$

 First: rescale to get rid of units, i.e., express plot in terms of "natural units" for the problem

- First: rescale to get rid of units, i.e., express plot in terms of "natural units" for the problem
- Intercept: location, scale, or something else?

- First: rescale to get rid of units, i.e., express plot in terms of "natural units" for the problem
- Intercept: location, scale, or something else?
- Slope: location, scale, or something else?

- First: rescale to get rid of units, i.e., express plot in terms of "natural units" for the problem
- Intercept: location, scale, or something else?
- Slope: location, scale, or something else?

$$pr(m, b | I) \propto \frac{1}{(1 + m^2)^{3/2}}$$
 Jaynes (1967); Jeffreys (1946)

"Symmetric prior"

- First: rescale to get rid of units, i.e., express plot in terms of "natural units" for the problem
- Intercept: location, scale, or something else?
- Slope: location, scale, or something else?

$$pr(m, b | I) \propto \frac{1}{(1 + m^2)^{3/2}}$$
 Jaynes (1967); Jeffreys (1946)

"Symmetric prior"

Standard deviation: location, scale, or something else?

- First: rescale to get rid of units, i.e., express plot in terms of "natural units" for the problem
- Intercept: location, scale, or something else?
- Slope: location, scale, or something else?

$$pr(m, b | I) \propto \frac{1}{(1 + m^2)^{3/2}}$$
 Jaynes (1967); Jeffreys (1946)

"Symmetric prior"

Standard deviation: location, scale, or something else?

$$\operatorname{pr}(\sigma|I) \propto \frac{1}{\sigma}$$

$$\sum_{i=1}^{6} p_i = 1$$

(something like this is always true)

$$\sum_{i=1}^{6} ip_i = 4.5$$

Suppose we have some information on our dice, e.g.

$$\sum_{i=0}^{6} p_i = 1$$
 (something like this is always true)

$$\sum_{i=1}^{6} ip_i = 4.5$$

Suppose we have some information on our dice, e.g.

$$\sum_{i=1}^{6} p_i = 1$$
 (something like this is always true)

$$\sum_{i=1}^{6} ip_i = 4.5$$

Principle of maximum entropy seeks to maximize

$$S = -\sum_{i=1}^{6} p_i \log(p_i)$$
 subject to these constraints

Suppose we have some information on our dice, e.g.

$$\sum_{i=1}^{6} p_i = 1$$
 (something like this is always true)

$$\sum_{i=1}^{6} ip_i = 4.5$$

Principle of maximum entropy seeks to maximize

$$S = -\sum_{i=1}^{6} p_i \log(p_i)$$
 subject to these constraints

■ Define 
$$Q(\{p_i\}; \lambda_0, \lambda_1) = -\sum_{i=1}^6 p_i \log(p_i) + \lambda_0 \left(1 - \sum_i p_i\right) + \lambda_1 \left(4.5 - \sum_i i p_i\right)$$

#### HAIR COLOR-HEIGHT EXAMPLE REVISITED

- I/3 of all Australians are blond, and I/4 are tall
- What proportion are blond and tall?

#### HAIR COLOR-HEIGHT EXAMPLE REVISITED

- I/3 of all Australians are blond, and I/4 are tall
- What proportion are blond and tall?

$$Q(\lbrace p_i \rbrace; \lambda_0, \lambda_1, \lambda_2) = -\sum_{i=1}^4 p_i \log(p_i) + \lambda_0 \left( 1 - \sum_{i=1}^4 p_i \right) + \lambda_1 \left( \frac{1}{3} - p_1 - p_2 \right) + \lambda_2 \left( \frac{1}{4} - p_1 - p_3 \right)$$

#### HAIR COLOR-HEIGHT EXAMPLE REVISITED

- I/3 of all Australians are blond, and I/4 are tall
- What proportion are blond and tall?

$$Q(\lbrace p_i \rbrace; \lambda_0, \lambda_1, \lambda_2) = -\sum_{i=1}^4 p_i \log(p_i) + \lambda_0 \left( 1 - \sum_{i=1}^4 p_i \right) + \lambda_1 \left( \frac{1}{3} - p_1 - p_2 \right) + \lambda_2 \left( \frac{1}{4} - p_1 - p_3 \right)$$

	BLOND	BROWN	SUMS
TALL	×	1/4-×	1/4
SHORT	1/3-×	5/12+×	3/4
SUMS	1/3	2/3	

Thomas Shahan - https://www.flickr.com/photos/

 A model for assigning probabilities to M different alternatives: monkeys throwing N balls into M equally sized boxes



- A model for assigning probabilities to M different alternatives: monkeys throwing N balls into M equally sized boxes
- Number of micro-states, W, as a function of {p<sub>i</sub>} is

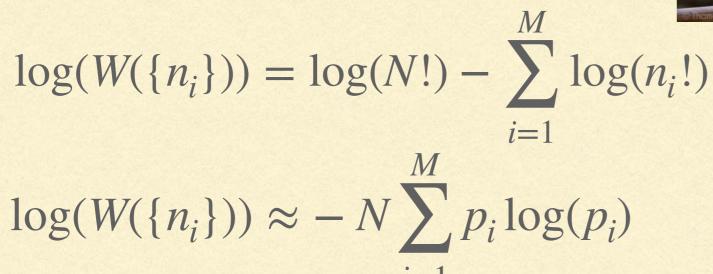


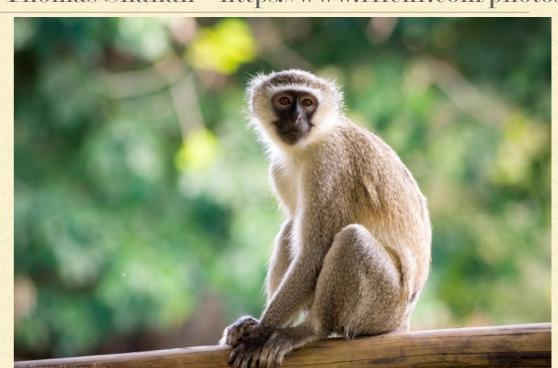
- A model for assigning probabilities to M different alternatives: monkeys throwing N balls into M equally sized boxes
- Number of micro-states, W, as a function of {p<sub>i</sub>} is

$$\log(W(\{n_i\})) = \log(N!) - \sum_{i=1}^{M} \log(n_i!)$$



- A model for assigning probabilities to M different alternatives: monkeys throwing N balls into M equally sized boxes
- Number of micro-states, W, as a function of {p<sub>i</sub>} is



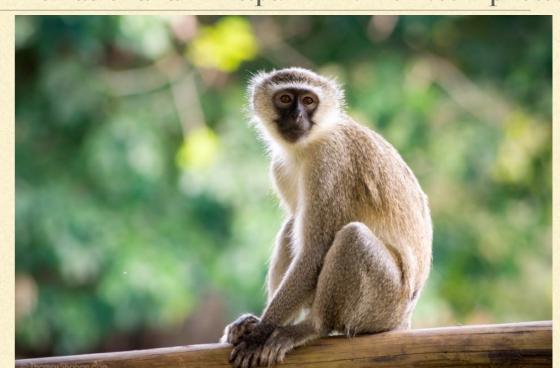


- A model for assigning probabilities to M different alternatives: monkeys throwing N balls into M equally sized boxes
- Number of micro-states, W, as a function of {p<sub>i</sub>} is

$$\log(W(\{n_i\})) = \log(N!) - \sum_{i=1}^{M} \log(n_i!)$$

$$\log(W(\lbrace n_i \rbrace)) \approx -N \sum_{i=1}^{M} p_i \log(p_i)$$

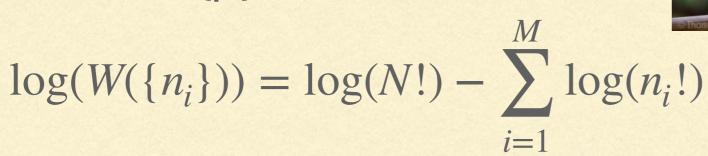
So define 
$$S=-\sum_{i=1}^{M} p_i \log(p_i)$$



### MONKEYS ARGUMENT

Thomas Shahan - https://www.flickr.com/photos/

- A model for assigning probabilities to M different alternatives: monkeys throwing N balls into M equally sized boxes
- Number of micro-states, W, as a function of {p<sub>i</sub>} is



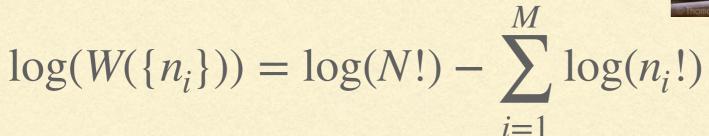
$$\log(W(\lbrace n_i \rbrace)) \approx -N \sum_{i=1}^{M} p_i \log(p_i)$$

So define 
$$S=-\sum_{i=1}^{M} p_i \log(p_i)$$
 Entropy

### MONKEYS ARGUMENT

Thomas Shahan - https://www.flickr.com/photos/

- A model for assigning probabilities to M different alternatives: monkeys throwing N balls into M equally sized boxes
- Number of micro-states, W, as a function of {p<sub>i</sub>} is



$$\log(W(\lbrace n_i \rbrace)) \approx -N \sum_{i=1}^{M} p_i \log(p_i)$$

So define 
$$S=-\sum_{i=0}^{M} p_i \log(p_i)$$
 Entropy

i=1



Jaynes: derive thermo using MaxEnt to assign pdfs subject to macro constraints

### WHY MAXIMIZETHE ENTROPY?

- Information theory: maximum entropy=minimum information (Shannon, 1948)
- Logical consistency (Shore & Johnson, 1960)
- Uncorrelated assignments related monotonically to S (Skilling, 1988)

### WHY MAXIMIZETHE ENTROPY?

- Information theory: maximum entropy=minimum information (Shannon, 1948)
- Logical consistency (Shore & Johnson, 1960)
- Uncorrelated assignments related monotonically to S (Skilling, 1988)

VARIATIONAL FUNCTION	OPTIMAL X	IMPLIED CORRELATION
$-\sum p_i \log(p_i)$	0.0833	None
$-\sum p_i^2$	0.0417	Negative
$\sum \log(p_i)$	0.1060	Positive
$\sum \sqrt{p_i}$	0.0967	Positive

### CONTINUOUS CASE

 Return to monkeys, but now with different probabilities for each bin. Then

$$S = -\sum_{i=1}^{M} p_i \log \left(\frac{p_i}{m_i}\right)$$

Shannon-Jaynes entropy
Kullback number
cross entropy

### CONTINUOUS CASE

 Return to monkeys, but now with different probabilities for each bin. Then

$$S = -\sum_{i=1}^{M} p_i \log \left(\frac{p_i}{m_i}\right)$$

Shannon-Jaynes entropy
Kullback number
cross entropy

Continuous case

### CONTINUOUS CASE

 Return to monkeys, but now with different probabilities for each bin. Then

$$S = -\sum_{i=1}^{M} p_i \log \left(\frac{p_i}{m_i}\right)$$

Shannon-Jaynes entropy
Kullback number
cross entropy

Continuous case

$$S[p] = -\int p(x) \log \left[ \frac{p(x)}{m(x)} \right]$$

m(x) makes the expression invariant under a change of variables

- m(x) makes the expression invariant under a change of variables
- Example I: maximize S subject only to normalization

- m(x) makes the expression invariant under a change of variables
- Example I: maximize S subject only to normalization

$$Q[p;\lambda] = -\int dx \, p(x) \, \log\left[\frac{p(x)}{m(x)}\right] + \lambda \left(1 - \int dx \, p(x)\right)$$

Solution: p(x) proportional to m(x)

- m(x) makes the expression invariant under a change of variables
- Example I: maximize S subject only to normalization

$$Q[p;\lambda] = -\int dx \, p(x) \, \log\left[\frac{p(x)}{m(x)}\right] + \lambda \left(1 - \int dx \, p(x)\right)$$

Solution: p(x) proportional to m(x)

Example 2: maximize S[p] subject to constraint

$$\int dx \, p(x)(x-\mu)^2 = \sigma^2$$

- m(x) makes the expression invariant under a change of variables
- Example I: maximize S subject only to normalization

$$Q[p;\lambda] = -\int dx \, p(x) \, \log\left[\frac{p(x)}{m(x)}\right] + \lambda \left(1 - \int dx \, p(x)\right)$$

Solution: p(x) proportional to m(x)

Example 2: maximize S[p] subject to constraint

$$\int dx \, p(x)(x-\mu)^2 = \sigma^2$$

Solution: p(x) proportional to m(x)