

# Statistical Language Models

## 2019

### Week 5 part 1

Dr. Dave Campbell  
[davecampbell@math.carleton.ca](mailto:davecampbell@math.carleton.ca)

# Approximate Course Outline

- Week 1: ShinyApps and Dashboarding
- Week 2: TidyText & obtaining data, dealing with time events
- Week 3: Regular Expressions; Word co-occurrence explorations
- Week 4: Sentiment Analysis; Stochastic process models
- Week 5: Exponential models for time between events.
- Week 6: Bayesian Basics; Author attribution models; hierarchical models
- Week 7: MCMC Diagnostics
- Week 8: Embeddings and Word2Vec; Cryptography
- Week 9: Clustering; Latent Dirichlet Allocation and topic models.
- Week 10: Variational Inference
- Week 11: Getting Fancier with Language Models
- Week 12: Student projects and presentations

# Cervical Cancer Example

- $N = 5736$
- Observed cancer:  $Y = 1$
- $P(\theta) \propto \theta^{\alpha-1}(1-\theta)^{\beta-1}$ , this is a Beta( $\alpha, \beta$ ) distribution,  $\alpha, \beta$  are equivalent to the prior number of observed cancer and non-cancers
- Using the non-vaccinated group we may choose  $\alpha = 36$ ,  $\beta = 5766$

# Cervical Cancer Example

- $N = 5736$
- Observed cancer:  $Y = 1$
- $P(\theta) \propto \theta^{\alpha-1}(1 - \theta)^{\beta-1}$ ,  $\alpha, \beta$  are equivalent to the prior number of observed cancer and non-cancers. Here  $\alpha = 36, \beta = 5766$
- $P(Y | n, \theta) \propto \theta^y(1 - \theta)^{n-y}$  this is a Binomial( $n, \theta$ )
- $P(\theta | Y, n) \propto P(Y | n, \theta)P(\theta) = \theta^{y+\alpha-1}(1 - \theta)^{n-y+\beta-1}$  this is Beta( $\alpha+y, \beta+n-y$ )

# MCMC

- Create a stochastic process  $X(t)$  such that its limiting probabilities match our target distribution.
- If the chain is irreducible, aperiodic, and not transient then a sample from Metropolis Hastings is guaranteed converge to the appropriate limiting probabilities.

# Metropolis Hastings

- Start with  $X(t-1) = j$
- Propose a value  $Y(t) \mid X(t-1)=j$  from transition probability matrix  $Q$  as a candidate for  $X(t)$
- compute  $\alpha_{ij} = \min\left(\frac{P(Y)P_{ji}}{P(X[t])P_{ij}}, 1\right)$  and sample  $u \sim \text{Unif}(0,1)$
- If  $u < \alpha_{ij}$  then accept the proposal and set  $X_t=Y$  and if not then set  $X(t)=X(t-1)$ .
- Repeat (T times) until you obtain a sufficient sample from the distribution of  $X$

```

Niters = 1000000
samples = c(1/5736, rep(0,1,Niters))
y = 1; N = 5736
a = 36 ; b = 5766
stepvar = .004
for(lp in 1:(Niters-1)){
  x = runif(1,samples[lp]-stepvar,samples[lp]+stepvar)          # sample from arbitrary transition distribution
  if(x>0 & x<1){
    alpha = dbinom(y,N,x) * (x^(a-1)) * (1-x)^(b-1) /
            (dbinom(y,N,samples[lp])*(samples[lp]^(a-1))*(1-samples[lp])^(b-1))
    u = runif(1)
    if(u<alpha){
      samples[lp+1] = x
    }else{
      samples[lp+1] = samples[lp]
    }
  }
  else{
    samples[lp+1] = samples[lp]
  }
}

hist(samples,probability=TRUE,500)
x = seq(0,.01,length=1000)
lines(x,dbeta(x, a+y, b+N-y ),lwd=4,col="red")

```

# Language Model

- Probability Transition Matrix from "AN" to "NA", "NB", ... "N\_"
- Each row is a Dirichlet (sums to 1, 27 valid categories)

- $$P(\theta) = \frac{\Gamma(\alpha_A + \dots + \alpha_-)}{\Gamma(\alpha_A) \dots \Gamma(\alpha_-)} \theta_A^{\alpha_A-1} \dots \theta_-^{\alpha_- - 1}$$



- Update the language model using books.
- Data generating process (likelihood):  $P(\text{book} \mid \theta)$
- $P(\text{observed transitions from "AN"} = [n_A, \dots, n_-] \mid \theta) = \text{multinomial for each transition}$
- $P(\text{observed transitions from AN} = [n_A, \dots, n_-] \mid \theta) \propto \theta_{AN,NA}^{n_A} \cdots \theta_{AN,N_-}^{n_-}$

# Now the cypher

- Cypher  $\phi$ , language model  $\theta$ , and (deterministic) decoded message  $Y(\phi)$
- We can directly calculate  $P(Y(\phi) \mid \theta)$  by deciphering the message and summing over probabilities from transitions:

$$P(Y(\phi) \mid \theta) = \prod_t P(Y_t(\phi) \mid \theta)$$

- Sample directly from  $P(Y(\phi) \mid \theta)$  using MCMC

- End result: a sample from the data generating mechanism  $P(Y(\phi) \mid \theta)$  using MCMC

# Author Attribution Models

- $P(\text{Author} \mid \text{text}) = P(\text{text} \mid \text{Author}) P(\text{Author}) / P(\text{text})$
- Requires  $P(\text{Author})$  to be feasible
- Possible models for  $P(\text{text} \mid \text{Author})$ ?