

Sampling for estimation

Idea: Parameter estimation via Simulation

$$y = \beta_0 + \beta_1 x + \epsilon \quad \epsilon \sim N(0, \sigma^2)$$

$$y \sim N(\beta_0 + \beta_1 x, \sigma^2) \rightarrow \text{Assume known.}$$

$$\beta \sim N(\vec{\phi}, \vec{\sigma}^2 \mathbf{I}) \quad \text{Prior distribution over } \beta$$

Suppose we observe data $D = \{x, y\}^n$

$$\begin{aligned} P(\beta | D) &= ? \\ &= \frac{P(D | \beta) P(\beta)}{P(D)} \quad \text{Bayes!} \end{aligned}$$

Note that D is observed

$$\begin{aligned} P(D | \beta, \sigma^2) &= \prod_i P(y_i | x_i, \beta, \sigma^2) \\ &= \prod_i N(y_i | \beta_0 + \beta_1 x_i, \sigma^2) \end{aligned}$$

$$P(\beta) = N(\beta | \vec{\phi}, \vec{\sigma}^2 \mathbf{I})$$

Here is a dumb algorithm to estimate β .

For T steps

Sample a $\hat{\beta}_t$ somehow

$w_t \leftarrow P(\hat{\beta}_t | D)$

$Z = \sum_{\tau} w_{\tau}$

$\hat{\beta}^* \leftarrow \sum_{\tau} \frac{w_{\tau} \hat{\beta}_{\tau}}{Z}$

weighted average over $\hat{\beta}$'s

A type of Importance Sampling.

We use weighted samples to approx. a posterior over β . Weights here are unnormalized: hence Z .

Note: Can use same strategy to Predict \hat{y}_i . How?

More generally: Assume a data generating distribution f w. parameters θ , and a prior π .

$$\underbrace{\int f(\theta) \pi(\theta) d\theta}_{\text{might be intractable}} \approx \underbrace{\frac{1}{S} \sum_{s=1}^S f(\theta^s)}_{\text{approximate by sampling}} \quad \theta^s \sim \pi$$

"Monte Carlo" Integration

$$\frac{1}{S} \sum_{s=1}^S \frac{w^s}{Z} \theta^s \quad \text{where} \quad Z = \sum_{s=1}^S w^s$$

and w^s is an unnormalized joint $P(D, \theta^s)$.

Obviously, π is critical here.

The key trick:

$$P(\theta|D) = \frac{P(D|\theta) P(\theta)}{P(D)} \propto P(D, \theta)$$

[See Jupyter Notebook]

For more complex cases, simple importance weighting is not going to fly. It would take forever to find good θ .

Can we be smarter about picking θ^s ?

Markov Chain Monte Carlo (MCMC) is a method that tries to simulate draws from a dist. of interest.

$$P(\theta^{(s+1)} | \theta^1 \dots \theta^s) = T(\theta^{(s+1)} \leftarrow \theta^s)$$

→ transition probability from current parameters.

Metropolis-Hastings is a particular version of MCMC. Basically, start somewhere θ^0 . Then: make a proposal θ^{t+1} that you accept or reject with some probability.

The **accept** p should be high if it is a better fit.

Gibbs Sampling is a simple recipe where we update a particular parameter θ_j conditioned on all others.

Gibbs - Sample

Initialize $\theta \sim q()$

for T steps

$$\theta_1^T \sim P(\theta_1 \mid \theta_2^{(T-1)}, \theta_3^{(T-1)}, \dots, \theta_m^{(T-1)})$$

$$\theta_2^T \sim P(\theta_2 \mid \theta_1^{(T-1)}, \theta_3^{(T-1)}, \dots, \theta_m^{(T-1)})$$

...

$$\theta_m^T \sim P(\theta_m \mid \theta_1^{(T-1)}, \theta_2^{(T-1)}, \dots, \theta_{m-1}^{(T-1)})$$

Return θ^T

(Content below derived from Jordan Boyd-Graber)

Coming back to LDA

$$\beta_k \sim \text{Dirichlet}(\eta)$$

$$\Theta_d \sim \text{Dirichlet}(\alpha)$$

$$z_{d,n} \sim \text{Discrete}(\Theta_d)$$

$$w_{d,n} | z_{d,n} \sim \text{Discrete}(\beta_{z_{d,n}})$$

the
model

For LDA, we will estimate the probability of a specific word's topic assignment, conditioned on all other assignments.

$$P(z_{d,n} = k | \underbrace{\vec{z}_{-d,n}}_{\text{all other word topic assignments}}, \vec{w}, \alpha, \lambda)$$

$$= \frac{P(z_{d,n} = k, \vec{z}_{-d,n} | \vec{w}, \alpha, \lambda)}{P(\vec{z}_{-d,n} | \vec{w}, \alpha, \lambda)}$$

Requires Integrating out Θ and β . This is a bit hairy, but we end with:

$$= \frac{n_{d,k} + \alpha_k}{\sum_i^k n_{d,i} + \alpha_i} \cdot \frac{V_{k,w_{d,n}} + \eta_{w_{d,n}}}{\sum_{w'} V_{k,w'} + \eta_{w'}}$$

How much This doc "likes" This topic How much Topic k "likes" This word

Count of topic k in doc d

Count of topic k using word $w_{d,n}$

Note! Given $P(z_{d,n})$ for all $z_{d,n}$, we can derive Θ, β