

LeBron James : $\{0, 1, 0, 0, 1, 1, 0, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, \dots\}$

exchangeable data model (EDM)

binary data

$$x_i \in \{0, 1\}$$

dataset: $x_{1:n}$

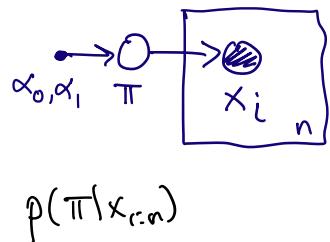
assume:

1) each shot is indep drawn from $p(x|\pi)$

2) the parameter comes from a prior-

$$p(\pi, x) = p(\pi) \prod_{i=1}^n p(x_i | \pi)$$

↑ ↑
prior likelihood



Bernoulli:

$$p(x|\pi) = \pi^x (1-\pi)^{1-x} \quad \pi \in [0, 1]$$

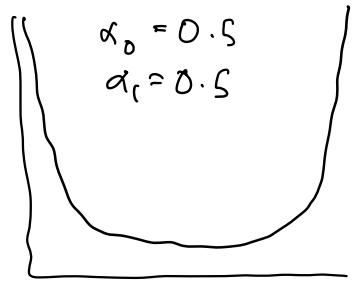
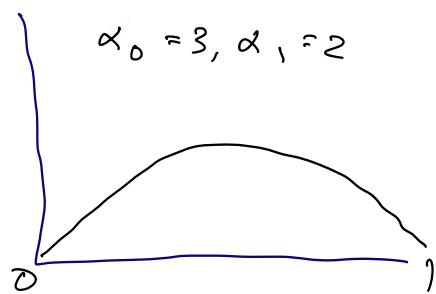
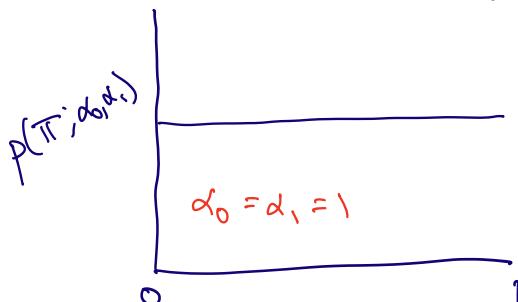
$$\mathbb{E}[X|\pi] = \pi$$

$$\text{Var}(X|\pi) = \pi(1-\pi)$$

Beta

↖ normalizing constant

$$p(\pi; \alpha_0, \alpha_1) = \frac{\Gamma(\alpha_1 + \alpha_0)}{\Gamma(\alpha_1)\Gamma(\alpha_0)} \pi^{\alpha_1 - 1} (1-\pi)^{\alpha_0 - 1} \quad \alpha_0, \alpha_1 > 0 \\ \pi \in (0, 1)$$



$$\mathbb{E}[\pi; \alpha_1, \alpha_0] = \frac{\alpha_1}{\alpha_1 + \alpha_0}$$

$$\mathbb{E}[1-\pi; \alpha_1, \alpha_0] = \frac{\alpha_0}{\alpha_1 + \alpha_0}$$

$$\text{Var}(\pi; \alpha_1, \alpha_0) = \frac{\alpha_1 \alpha_0}{(\alpha_1 + \alpha_0)^2 (\alpha_1 + \alpha_0 + 1)}$$

$x_{1:n}$: binary dataset

$$\pi \sim \text{Beta}(\alpha_0, \alpha_1)$$

$$p(\pi) \prod_{i=1}^n p(x_i | \pi)$$

↑ ↑
Beta Bernoulli

$$x_i \sim \text{Bern}(\pi) \quad i = 1 \dots n$$

$$p(a|b) \& p(a, b)$$

$$\begin{aligned} p(\pi | x_{1:n}; \alpha_0, \alpha_1) &\propto p(\pi; \alpha_0, \alpha_1) \prod_{i=1}^n p(x_i | \pi) \\ &= \pi^{\alpha_1 - 1} (1-\pi)^{\alpha_0 - 1} \prod_{i=1}^n \pi^{x_i} (1-\pi)^{1-x_i} \\ &= \frac{(\alpha_1 + \sum_{i=1}^n x_i - 1)}{\pi} \frac{(\alpha_0 + \sum_{i=1}^n (1-x_i) - 1)}{(1-\pi)} \end{aligned}$$

This is a beta distribution.

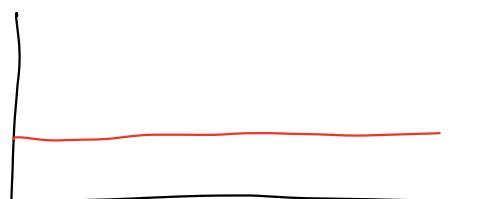
$$\begin{aligned} \hat{\alpha}_1 &= \alpha_1 + \sum_{i=1}^n x_i & n_1 \\ \hat{\alpha}_0 &= \alpha_0 + \sum_{i=1}^n (1-x_i) & n_0 \end{aligned}$$

posterior beta

Conjugacy - posterior is in the same family as the prior.

Beta / Bernoulli are called conjugate pair

$$\mathbb{E}[\pi | x_{1:n}] = \frac{\alpha_1 + n_1}{\alpha_1 + \alpha_0 + n_1 + n_0} = \frac{\alpha_1 + n_1}{\alpha_1 + \alpha_0 + n}$$



$$\text{joint} = p(\pi, x_{1:n})$$

$$\text{posterior} : p(\pi | x_{1:n})$$

$$\begin{aligned}\text{predictive} : p(x' | x_{1:n}) &= \int p(\pi, x' | x_{1:n}) d\pi \\ &= \int p(x' | \pi, x_{1:n}) p(\pi | x_{1:n}) d\pi \\ &= \int \pi \cdot p(\pi | x_{1:n}) d\pi \\ &= \mathbb{E}[\pi | x_{1:n}]\end{aligned}$$

Categorical data

Count data

Real-valued data

x_i takes on one of K "categories" e.g., terms in a vocabulary.

Example: Text.

"To be or not to be, that is the question."

$x_1 \ x_2 \ x_3 \ x_4 \ x_5$

indicator vector, "one-hot" vector $x \in \mathbb{I}^K$

$$\begin{array}{l} \text{a} \\ \text{of} \\ \text{the} \\ \text{Juliet} \\ \text{Romeo} \\ \text{be} \end{array} \left(\begin{array}{c} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{array} \right) \mid K$$

Categorical distribution

parameter θ - prob of each value

$$\theta = \left(\begin{array}{c} 0.01 \\ 0.02 \\ 0.0 \\ 0.005 \\ \vdots \end{array} \right) \mid \sum_{k=1}^K \theta_k = 1 \quad \text{1-Simplex}, \quad \theta \in \Delta^{K-1}$$

$$\theta_k \geq 0$$

$$p(x|\theta) = \prod_{k=1}^K \theta_k^{x^{(k)}} \quad x^{(k)} \sim \text{Bern}(\pi) \quad \pi \sim \text{Dir}(1-\pi)$$

$$\mathbb{E}[x^{(k)}|\theta] = \theta_k$$

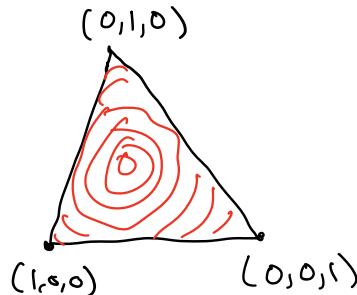
$$\text{Var}[x^{(k)}|\theta] = \theta_k(1-\theta_k)$$

$$\text{Cov}[x^{(j)}, x^{(k)}|\theta] = -\theta_j \theta_k$$

$$p(\theta, x_{1:n}) = p(\theta) \prod_{i=1}^n p(x_i|\theta)$$

$$p(\theta; \alpha) = \left(\frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \right) \prod_{k=1}^K \theta_k^{\alpha_k - 1}$$

$$\alpha = (\alpha_1, \dots, \alpha_K), \alpha_k > 0$$



$$\mathbb{E}[\theta_j; \alpha] = \frac{\alpha_j}{\sum_{k=1}^K \alpha_k} \quad \text{Var}[\theta_j; \alpha] = \frac{\alpha_j (\sum_{k \neq j} \alpha_k)}{(\sum_k \alpha_k)^2 (\sum_k \alpha_k + 1)}$$

$$\text{Cov}[\theta_j, \theta_\ell; \alpha] = \frac{-\alpha_j \alpha_\ell}{\sum_{k=1}^K \alpha_k}$$

Exchangeable Dirichlet

$$p(\theta; \alpha_0) : \text{Dir}(\underbrace{\alpha_0, \alpha_0, \dots, \alpha_0}_{K \text{ times}})$$

$$\theta \sim \text{Dir}_K(\alpha_0)$$

$$x_i \sim \text{Cat}(\theta)$$

$$p(\theta | x_{1:n}) : \text{Dirichlet}(\hat{\alpha}_1, \dots, \hat{\alpha}_K)$$

$$\hat{\alpha}_k = \alpha_0 + \sum_{i=1}^n x_i^{(k)} \quad \sum_i x_i^{(k)}$$

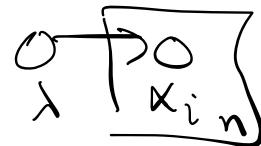
$$x: \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} \quad \alpha: \begin{pmatrix} 3 \\ 2 \\ 9.6 \\ 0.5 \end{pmatrix}$$

$$\theta: \begin{pmatrix} p(x=1) \\ p(x=2) \\ p(x=3) \\ p(x=4) \end{pmatrix}$$

$$\mathbb{E}[\theta_k | x_{1:n}; \alpha] = \frac{\alpha_k + n_k}{\sum_{j=1}^k \alpha_j + n} \quad \text{Laplace smoothing.}$$

daily clicks: $\{0, 5, 3, 1, 9, 6, \dots\}$ count data.

$$p(\lambda; x_{1:n}) = p(\lambda) \prod_{i=1}^n p(x_i | \lambda)$$



Poisson distribution

$$p(x | \lambda) = \frac{\lambda^x e^{-\lambda}}{x!} \quad \lambda \geq 0$$

$$\mathbb{E}[x | \lambda] = \lambda \quad \text{Var}[x | \lambda] = \lambda$$

$$p(\lambda; \alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \exp\{-\beta\lambda\} \lambda^{\alpha-1} \quad \text{Gamma prior}$$

$$\alpha > 0, \beta > 0$$

$\alpha = 1$: exponential dist.

$$\mathbb{E}[\lambda] = \alpha / \beta \quad \text{Var}[\lambda] = \alpha / \beta^2$$

$$p(\lambda | x_{1:n}; \alpha, \beta) : \text{Gamma}(\hat{\alpha}, \hat{\beta})$$

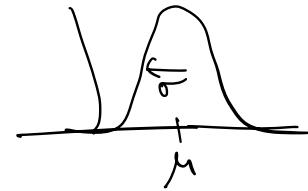
$$\begin{aligned}\hat{\alpha} &= \alpha + \sum_{i=1}^n x_i \\ \hat{\beta} &= \beta + n\end{aligned}$$

$$\mathbb{E}[\lambda | x_{1:n}; \alpha, \beta] = \frac{\alpha + \sum_{i=1}^n x_i}{\beta + n}$$

$$\mathbb{E}[\lambda; \hat{\alpha}, \hat{\beta}]$$

Gaussian.

$$p(x | \mu, \sigma^2) \quad x \in \mathbb{R} \quad \begin{array}{l} \sigma^2 > 0 \\ \mu \in \mathbb{R} \end{array}$$



$$p(\mu, x_{\text{obs}}) = p(\mu) \prod_{i=1}^n p(x_i | \mu)$$

$$p(x_i | \mu) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{(x_i - \mu)^2}{2\sigma^2} \right\}$$

$$\mu \sim N(\mu_0, \lambda^2)$$

$$x_i | \mu \sim N(\mu, \sigma^2)$$

$$\mu | x_{\text{obs}} \sim N(\hat{\mu}_0, \hat{\lambda}^2) \quad \bar{x} = \sum_i x_i / n$$

$$\hat{\mu}_0 = \left(\frac{\mu_0 / \lambda^2 + n \bar{x} / \sigma^2}{1 / \lambda^2 + n / \sigma^2} \right)$$

hyperparameters. parameters to the priors.

"Bayesian purist"

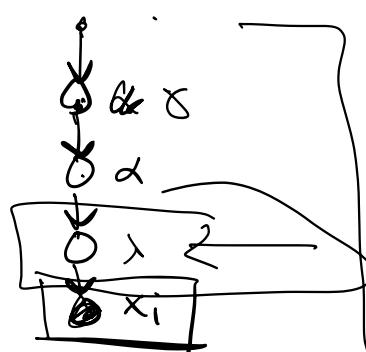
Model Criticism.

"Orthodox Bayesian"

"Crayesians" ~~think~~

Uniform prior.

"Vague prior"

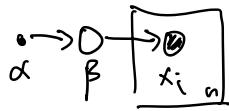


filled out by likelihood.

$$\beta \sim p(\beta; \alpha)$$

$$x_i \sim p(x_i | \beta)$$

generative process



graphical model

$$p(\beta, x_{1:n}) = p(\beta; \alpha) \prod_{i=1}^n p(x_i | \beta)$$

factorized joint

$$p(\beta | x_{1:n}; \alpha)$$

posterior



$$p(x_{\text{new}} | x_{1:n}; \alpha)$$

posterior predictive

Held out log probability

$$p(\theta, x_{1:n}) = p(\theta; \alpha_0) \prod_{i=1}^n p(x_i | \theta)$$

\uparrow
categorical dist
parameters
over terms

\uparrow
categorical
dist w/
param θ

$$p(x_{\text{new}} | x_{1:n}; \alpha_0) = \int p(x_{\text{new}} | \theta) p(\theta | x_{1:n}; \alpha_0) d\theta$$

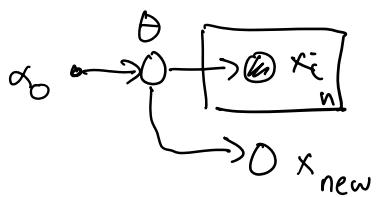
x_{new} : a word

θ : dist. over words, Δ^{V-1}

$$= \mathbb{E}[\theta_v | x_{1:n}; \alpha_0]$$

$$\hat{\alpha}_v = \alpha_0 + n_v \leftarrow \# \text{ times Shakespeare used term } v$$

$$\mathbb{E}[\theta_v | x_{1:n}; \alpha_0] = \frac{\alpha_0 + n_v}{\sqrt{\alpha_0 + n}}$$



$$\mathbb{E}[\theta | x_{1:n}] \triangleq \hat{\theta}$$

data splitting

$$x_{1:n}, 1:n_{in}$$

$$x_{out}, 1:n_{out}$$

$$\hat{\Lambda}(\alpha_0) \triangleq \sum_{i=1}^{n_{out}} \log p(x_{out,i} | x_{in}; \alpha_0)$$

Model selection.

$$\hat{\alpha}_0 = \arg \max_{\alpha_0} \log p(x_{i:n}; \alpha_0)$$

prior predictive check (Box)

posterior " " (Rubin...)

$$\int p(x_{i:n} | \theta) p(\theta) d\theta$$

$F(x)$ — true population distribution of a datapoint.

$$\Lambda(\alpha_0) \triangleq \mathbb{E}_{F} [\log p(X_{\text{new}} | x_{i:n}; \alpha_0)]$$

"proper scoring function"

$$\hat{\Lambda}(\alpha_0) \triangleq \frac{1}{n_{\text{out}}} \sum_{i=1}^{n_{\text{out}}} \underbrace{\log p(x_{i,\text{out}} | x_{i:n}; \alpha_0)}_{=}$$

Conditional models

Linear regression

Logistic regression

Stochastic optimization for MAP estimation.

$$y_i | x_i, \beta \sim N(\beta \cdot x_i, \sigma^2)$$

\uparrow \uparrow
 p-vector p-vector
 of of
 features coefficients

$$p(\beta, y_{1:n} | x_{i:n}; \lambda) = p(\beta; \lambda) \prod_{i=1}^n p(y_i | x_i, \beta)$$

$$\beta_k \sim N(0, \lambda^2)$$

$$p(\beta | y_{1:n}, x_{1:n}; \lambda)$$

$$\hat{\beta}_{MAP} = \arg \max_{\beta} p(\beta | y_{1:n}, x_{1:n}; \lambda)$$

$$\hat{\beta}_{MAP} = \arg \max_{\beta} \left(\log p(\beta) + \sum_{i=1}^n \log p(y_i | x_i, \beta) \right)$$

Generative process:

For each coefficient k :

$$y_i = \beta \cdot x_i + \epsilon_i$$

$$\epsilon_i \sim N(0, \sigma^2)$$

$$p(\beta | \mathcal{D}; \lambda) = \beta_k \sim N(0, \lambda^2) \text{ prior}$$

$$\frac{p(\beta) \prod p(y_i | x_i, \beta)}{p(y_{1:n} | x_{1:n})} \text{ For each data point } i:$$

$$y_i | x_i, \beta \sim N(\beta \cdot x_i, \sigma^2)$$

(x_i, y_i)

$i = 1, \dots, n$

$$p(\beta | y_{1:n}, x_{1:n}; \lambda)$$

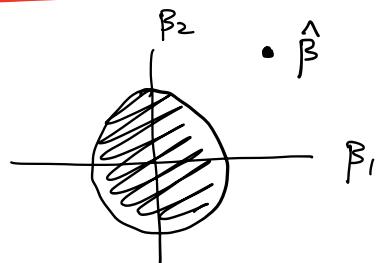
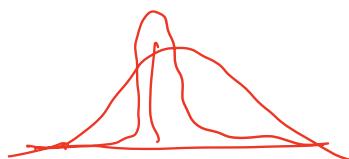
$$\hat{\beta}_{MAP} = \arg \max_{\beta} \log p(\beta) + \sum_{i=1}^n \log p(y_i | x_i, \beta) - \log p(y_{1:n} | x_{1:n})$$

$$= \arg \max_{\beta} \underbrace{\log p(\beta) + \sum_{i=1}^n \log p(y_i | x_i, \beta)}_{\text{log joint}} + \text{const.}$$

log joint

Ridge regression.

$$= \arg \max_{\beta} -\frac{1}{2\lambda^2} \sum_{k=1}^p \beta_k^2 - \sum_{i=1}^n \frac{(y_i - \beta \cdot x_i)^2}{2\sigma^2} + \text{const.}$$



Laplace prior or Lasso.

$$\beta_k \sim \text{Laplace}(\lambda) \quad k = 1 \dots p$$

$$y_i | x_i, \beta \sim N(\beta \cdot x_i, \sigma^2) \quad i = 1 \dots n$$

$$p(\beta_k; \lambda) = \frac{1}{2\lambda} \exp \left\{ -\frac{|\beta_k|}{\lambda} \right\}$$

$$\mathcal{L}_{\text{MAP}}(\beta) \triangleq -\frac{1}{\lambda} \sum_{k=1}^p |\beta_k| + \sum_{i=1}^n \log p(y_i | x_i, \beta)$$

MAP estimate tends to be sparse

Movie review data $(x_i, y_i) \quad i = 1 \dots n$

$x_i : (x_i^{(1)}, \dots, x_i^{(p)})$ word counts across a vocabulary.

$$y_i : \{0, 1\} \quad 0 = \text{bad} \\ 1 = \text{good}$$

$$\beta_k \sim N(0, \lambda^2) \quad k = 1 \dots p$$

$$y_i | x_i, \beta \sim \text{Bern}(\sigma(\beta \cdot x_i))$$

$\sigma(w)$: sigmoid function $\mathbb{R} \rightarrow (0, 1)$

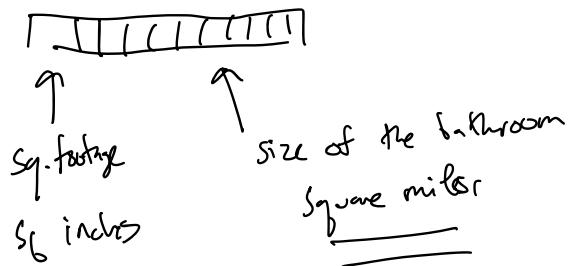
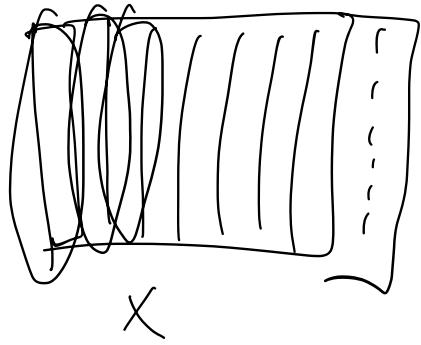
$$\sigma(\beta \cdot x_i) = \frac{1}{1 + \exp\{-\beta \cdot x_i\}} - \frac{1}{2\lambda^2} \alpha^2$$

$$\hat{\beta}_{\text{MAP}} = \arg \max_{\beta} -\frac{1}{2\lambda^2} \sum_{k=1}^p \beta_k^2 + \sum_{i=1}^n y_i \log \sigma(\beta \cdot x_i) + (1-y_i) \log \sigma(-\beta \cdot x_i)$$

$$p(z) = \pi^2 (1-\pi)^{1-z} \quad 1 - \sigma(w) = \sigma(-w)$$

$$p(\beta | y_{1:n}, x_{1:n}; \lambda)$$

- ① Standardizing the covariates
 $x_i, y_i \quad i=1 \dots n$



$$\textcircled{2} \quad \text{The intercept. } y_i \sim N\left(\frac{\beta \cdot x_i + \alpha}{n}, \sigma^2\right)$$

$$\beta_k \sim N(0, \lambda^2)$$

$$\alpha \sim N(0, \gamma^2) \quad \leftarrow \gamma^2 \text{ shall be big.}$$

Setting the regularizer.

$$\mathbb{E}[Y | x^*, \hat{\beta}_{MAP}] = \hat{\beta}_{MAP} \cdot x^* \quad \text{linear regression}$$

$$\mathbb{E}[Y | x^*, \hat{\beta}_{MAP}] = \sigma(\hat{\beta}_{MAP} \cdot x^*) \quad \text{logistic regression}$$

$$(x_{in}, y_{in}) \quad (x_{out}, y_{out})$$

$$\hat{\Lambda}(\lambda) = \frac{1}{m} \sum_{i=1}^m \log p(y_{i,out} | x_{i,out}, \hat{\beta}(\lambda))$$

$\hat{\beta}(\lambda)$ = MAP estimate with hyperparameter λ .

$$\beta_k \sim N(0, \lambda^2) \quad k=1 \dots p$$

$$y_i | x_i, \beta \sim p(y_i | x_i, \beta) \quad i=1 \dots n$$

$$\hat{\beta}_{MAP} = \arg \max_{\beta} \log p(\beta; \lambda) + \sum_{i=1}^n \log p(y_i | x_i, \beta)$$

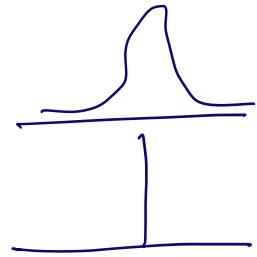
$$\Lambda(\lambda) \triangleq \underbrace{E_F [\log p(Y|X, \theta)]}_{\text{posterior predictive}} \quad \begin{aligned} x', y' &\sim F \\ \theta &= \{(x_i, y_i)\}_{i=1}^n - \text{observed} \end{aligned}$$

$$\hat{\Lambda}(\lambda) \approx \frac{1}{n_{out}} \sum_{i=1}^{n_{out}} \log p(y_{i,out} | x_{i,out}, \theta_{in}) \quad \theta_{in} = \{(x_{in,i}, y_{in,i})\}_{i=1}^{n_{in}}$$

$$p(y | x, \theta) = \int p(y | x, \beta) p(\beta | \theta) d\beta$$

$$\delta(\beta_{MAP})$$

$$\approx p(y | x, \underline{\beta_{MAP}(\theta)})$$



$$\hat{\Lambda}(\lambda) \approx \frac{1}{n_{out}} \sum_{i=1}^{n_{out}} \log p(y_{i,out} | x_{i,out}, \hat{\beta}_{MAP}^\lambda(\theta_{in}))$$

$$\equiv \hat{\Lambda}(\lambda) = \sum_{i=1}^{n_{out}} - \underbrace{(y_{i,out} - \hat{\beta}_{MAP}^\lambda \cdot x_{i,out})^2}_{2\sigma^2} / 2\sigma^2$$

$$\equiv \hat{\Lambda}(\lambda) = \sum_{i=1}^{n_{out}} y_{i,out} \log \sigma(\hat{\beta}_{MAP}^\lambda \cdot x_{i,out}) + (1 - y_{i,out}) \log \sigma(-\hat{\beta}_{MAP}^\lambda \cdot x_{i,out})$$

How do we find β_{MAP} ?

In : θ Out : β_{MAP}

$$\mathcal{L}(\beta) = \log p(\beta; \lambda) + \sum_{i=1}^n \log p(y_i | x_i, \beta)$$

gradient ascent.

$$\nabla_{\beta} \mathcal{L}(\beta) = \left(\frac{\partial}{\partial \beta_1} \mathcal{L}(\beta), \dots, \frac{\partial}{\partial \beta_p} \mathcal{L}(\beta) \right)$$

initialize $\beta[0]$

$$\beta[t+1] = \beta[t] + \rho t \nabla_{\beta} \mathcal{L}(\beta[t])$$

\uparrow "step size", "learning rate", 0.001

$$\mathcal{L}(\beta) = -\frac{1}{2\lambda^2} \sum_{k=1}^p \beta_k^2 - \sum_{i=1}^n \frac{(y_i - \beta \cdot x_i)^2}{2\sigma^2}$$

$$\frac{\partial}{\partial \beta_k} \mathcal{L}(\beta) = -\beta_k / \lambda^2 + \sum_{i=1}^n \frac{(y_i - \beta \cdot x_i)}{\sigma^2} x_{i,k}$$

$$\boxed{\nabla_{\beta} \mathcal{L}(\beta) = -\frac{1}{\lambda^2} \beta + \sum_{i=1}^n \frac{(y_i - \beta \cdot x_i)}{\sigma^2} x_i}$$

$$\begin{aligned} \text{Signed residual} &\triangleq (y_i - \mathbb{E}[Y|x_i, \beta]) \\ &= (y_i - \beta \cdot x_i) \end{aligned}$$

Logistic regression

$$\mathcal{L}(\beta) = -\frac{1}{2\lambda^2} \sum_{k=1}^p \beta_k^2 + \sum_{i=1}^n [y_i \log \sigma(\beta \cdot x_i) + (1-y_i) \log \sigma(-\beta \cdot x_i)]$$

$$\frac{d}{d\eta} \sigma(\eta) = \sigma(\eta)(1-\sigma(\eta))$$

$$\nabla \mathcal{L}(\beta) = -\frac{\beta}{\lambda^2} + \sum_{i=1}^n$$

$$\frac{d}{d\eta_i} \log p(y_i | x_i, \beta) = \frac{y_i \sigma(\eta_i)(1-\sigma(\eta_i))}{\sigma(\eta_i)} + \frac{(1-y_i) \sigma(-\eta_i)(1-\sigma(-\eta_i))}{\sigma(-\eta_i)}$$

$$\begin{aligned} \eta_i &\triangleq \beta \cdot x_i \\ &= y_i (1-\sigma(\eta_i)) - (1-y_i)(1-\sigma(-\eta_i)) \\ &= y_i (1-\sigma(\eta_i)) - (1-y_i)\sigma(\eta_i) \end{aligned}$$

$$= y_i - y_i \sigma(\eta_i) - \sigma(\eta_i) + y_i \sigma(\eta_i)$$

$$= y_i - \sigma(\eta_i)$$

$$\boxed{\nabla_{\beta} \mathcal{L}(\beta) = -\frac{1}{N} \beta + \sum_{i=1}^n (y_i - \sigma(\beta \cdot x_i)) x_i}$$

$$\epsilon_i = y_i - \mathbb{E}[Y | x_i, \beta] = \sigma(\beta \cdot x_i)$$

$$\boxed{\nabla_{\beta} \mathcal{L}(\beta) = -\frac{1}{N} \beta + \sum_{i=1}^n (y_i - \mathbb{E}[Y | x_i, \beta]) x_i}$$

Stochastic optimization

10M more reviews + sentiment.

large dataset $y_{1:n}$

$$\mathcal{L}(\beta; y_{1:n}) = h(\beta) + \sum_{i=1}^n f(\beta; y_i)$$

$$\hat{\beta} = \arg \max_{\beta} \mathcal{L}(\beta; y_{1:n})$$

$$\nabla_{\beta} \mathcal{L} = \nabla_{\beta} h(\beta) + \sum_{i=1}^n \nabla_{\beta} f(\beta; y_i)$$

stochastic gradient ascent.

Robbins + Monro (1951), Bottou (1998)

stochastic gradient:

$$g \sim p(g; \beta)$$

$$\mathbb{E}[G] = \nabla_{\beta} \mathcal{L}(\beta; y_{1:n})$$

unbiased stochastic gradient

repeat ~~for~~ until converge

$$g_t \sim p(g; \beta[t])$$

$$\beta[t+1] = \beta[t] + \rho_t g_t$$

ρ_t step size schedule

Robbins Monro conditions

$$\sum_t \rho_t = \infty \quad \sum_t \rho_t^2 < \infty$$

\hat{F} = empirical dist of data $\forall n$ on each data point

$$y \sim \hat{F}$$

$$g = \nabla_{\beta} h(\beta) + n \nabla_{\beta} f(\beta[t]; y)$$

$$\mathbb{E}_{\hat{F}}[g] = \nabla L(\beta) \quad - \text{confirm.}$$

$$\underline{\text{minibatch}} \quad g = \nabla_{\beta} h(\beta) + \frac{n}{B} \sum_{b=1}^B \nabla_{\beta} f(\beta[t]; y_b)$$

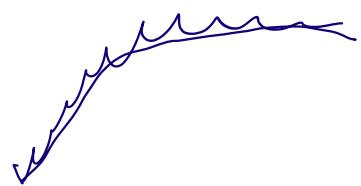
B batch size

$$y_1 \dots y_B \sim \hat{F}$$

Practice

① Convergence.

- norm of the noisy gradient
- objective plateaued
- validation objective



② Minibatches + data sweeps
w/o replacement

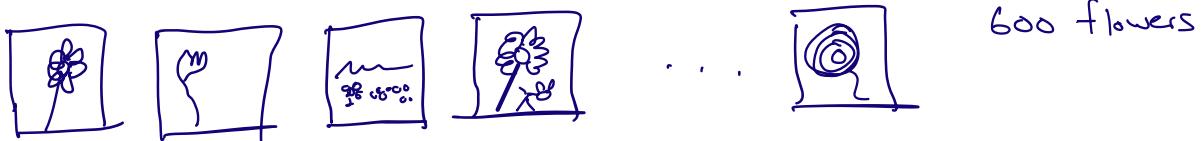
③ Step size Schedule

- Robbins Monro
- constant step size (η)
- Adaptive step sizes - AdaGrad (η)
ADAM (η)
RMS Prop (η)
AdaDelta (η)

$$p > n$$

- Mixture models, mixed membership models, factor models
- Gibbs sampling (MCMC), Variational inference.

Mixture model



group the data

- ① we don't know which data belongs to which group
- ② we don't know what the clusters are.

$$x_i : (x_{i,1}, \dots, x_{i,d}) - d \text{ dimensions, } d = 512$$

Core assumption: each x_i belongs to a cluster,
drawn from a distribution associated with it.

generative process, etc.

- observed
- n data points x_1, \dots, x_n
 - K - # mixture components.
- latent
- n assignments z_1, \dots, z_n
 - z_i is a K -categorical variable, $z_i \in \mathbb{I}^K$
 - K mixture components, β_1, \dots, β_K
parametrizes a dist over x
 - proportions $\theta \in \Delta^K$

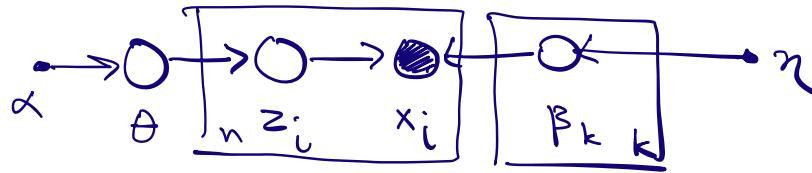
- ① generative process.

- ② factorized joint

$$p(\theta, \beta_{1:k}, z_{1:n}, x_{1:n}) = p(\theta) \prod_{k=1}^K p(\beta_k) \prod_{i=1}^n \left(p(z_i | \theta) p(x_i | z_i, \beta_{1:k}) \right)$$

$$\begin{cases}
 p(\theta) : \text{Dirichlet} \\
 p(\beta_k) : \text{Gaussians} \\
 p(z_i | \theta) : \text{Categorical}(\theta) \\
 p(x_i | z_i, \beta_{1:k}) = \prod_{k=1}^K f(x_i; \beta_k)^{z_i^{(k)}} \\
 f(x_i; \beta) = \prod_{j=1}^d \ell(x_{id}; \beta^{(d)})
 \end{cases}
 \quad f_{\beta_k}(x_i)$$

③



$$p(\theta, \beta_{1:k}, z_{1:n} | x_{1:n})$$

GOAL, Assumptions. DATA

\downarrow \downarrow
 MODEL \rightarrow POSTERIOR \rightarrow USE IT

10/5

Goal: Cluster data

Mixture model:

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

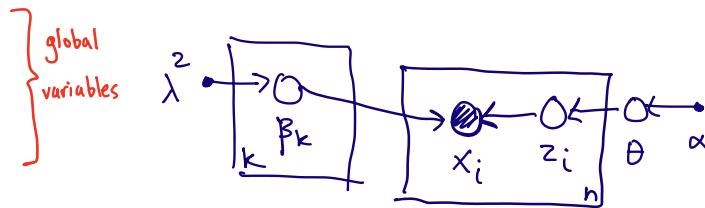
For each component k :

$$\beta_k \sim N_d(0, \lambda^2)$$

For each data point i :

$$z_i | \theta \sim \text{Cat}(\theta)$$

$$x_i | z_i, \beta_{1:k} \sim N(\beta_{z_i}, \sigma^2)$$



local variables

$$p(\underbrace{\theta, \beta_{1:k}, z_{1:n}, x_{1:n}}_{\text{hidden}}, \underbrace{x_{1:n}}_{\text{observed}}; \lambda, \alpha) = p(\theta) \prod_k p(\beta_k) \prod_i p(z_i | \theta) p(x_i | z_i, \beta_{1:k})$$

local latent variable: z_i

posterior mixture

$$p(\theta, \beta_{1:k}, z_{1:n} | x_{1:n})$$

$$\log p(\theta, \beta_{1:k}, z_{1:n} | x_{1:n}) = \underbrace{\log p(\theta, \beta_{1:k}, z_{1:n}, x_{1:n})}_{\text{posterior}} - \log p(x_{1:n})$$

$$\log p(-) = \underbrace{\log p(\theta)}_{\theta = 1/k} + \underbrace{\sum_{k=1}^K \log p(\beta_k)}_{\text{global}} + \underbrace{\sum_{i=1}^n \log p(z_i | \theta)}_{\text{local}} + \underbrace{\sum_{i=1}^n \log p(x_i | z_i, \beta_{z_i})}_{\text{observed}}$$

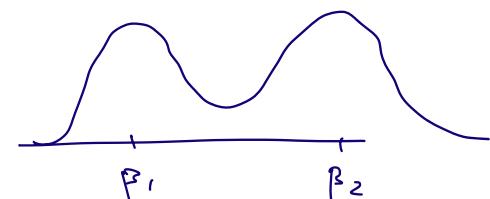
$$\theta = 1/k$$

n is huge

$$\sum_{i=1}^n \log p(x_i | z_i, \beta_{z_i}) = \sum_{i=1}^n \log f(x_i; \beta_{z_i})$$

$$\log f(x_i; \beta_k) = -\frac{1}{2\sigma^2} \sum_{j=1}^d (x_i^{(j)} - \beta_k^{(j)})^2 + \text{const.}$$

$$p(x | \beta_{1:k}) = \sum_{k=1}^K p(x | \beta_k) \theta_k$$



$$\log p(z_i | \theta) = -\log K$$

approximate posterior inference — Gibbs sampling
 Variational inference
 MAP estimation

Gibbs sampler for a model

input: data

output: samples of latent variables from the posterior.

$$p(\theta, \beta_{1:k}, z_{1:n} | x_{1:n}) = \frac{p(\cdot)}{\int \int \sum_{z'_{1:n}} p(\theta', \beta'_{1:k}, z'_{1:n}, x_{1:n}) d\beta'_{1:k} d\theta}$$

$\underbrace{\qquad\qquad\qquad}_{p(x_{1:n})}$

- initialize θ, β, z to something
- iteratively sample from complete conditionals

$$\theta \sim p(\theta | \beta, z, x)$$

$$z_i \sim p(z_i | \theta, \beta, z_{-i}, x) \quad i=1..n$$

$$\beta_k \sim p(\beta_k | \theta, \beta_{-k}, z, x) \quad k=1..K$$

General model $p(h, x)$ $h = h_{1:d}$

Goal $p(h | x)$

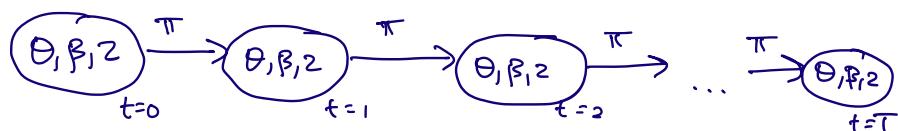
Gibbs $p(h_j | h_{-j}, x) \quad j=1..d$

Why does this recipe work?

MCMC - Markov chain Monte Carlo

$$\{\theta, \beta, z\}_0 \sim \pi_0(\theta, \beta, z) \quad \swarrow \text{transition distribution.}$$

$$\{\theta, \beta, z\}_t \mid \{\theta, \beta, z\}_{t-1} \sim \pi(\theta, \beta, z \mid \{\theta, \beta, z\}_{t-1})$$



$\pi_t(\theta, \beta, z) \triangleq$ margin at time t $\pi_\infty(\theta, \beta, z) = p(\theta, \beta, z | x)$

Complete conditional of θ

$$p(\theta | \beta_{r:k}, z_{1:n}, x_{1:n}) \propto p(\theta, \beta_{r:k}, z_{1:n}, x_{1:n})$$

posterior
 $p(\theta | x_{1:n})$ $\propto p(\theta) \prod_{i=1}^n p(z_i | \theta)$
 Dirichlet Categorical

$$\begin{aligned} p(a, b) &= p(a)p(b) \\ p(a|b) &= p(a) \end{aligned}$$

$$\theta \perp\!\!\!\perp \{\beta_{r:k}, x_{1:n}\} | z_{r:n}$$

$$n_k(z) \triangleq \sum_{i=1}^n z_i^{(k)} = \text{Dirichlet}(\hat{\alpha}(z)) \geq (n_1(z), \dots, n_K(z))$$

$$\hat{\alpha}(z) = \alpha + n(z)$$

Complete conditional of the assignments

$$p(z_i | \beta_{r:k}, \theta, z_{-i}, x_{1:n}) \in \text{joint distribution}$$

$$\propto p(z_i | \theta) p(x_i | z_i, \beta_{r:k})$$

$$= \theta_{z_i} f(x_i; \beta_{z_i})$$

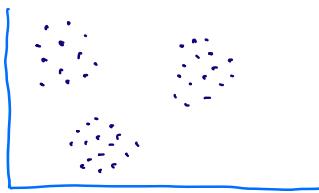
$$p(z_i = k | \dots) = \frac{\theta_k f(x_i; \beta_k)}{\sum_{k'=1}^K \theta_{k'} f(x_i; \beta_{k'})}$$

Complete conditional of the components

We know: how the data are partitioned z_i

$p(\beta)$: Gaussian

$f(x; \beta)$: Gaussian w/ mean β (fixed variance)



$p(\beta_k | x_{1:n}, z_{r:n})$: Gaussian with "posterior" mean and variance from the EMM

=

- the practice of Gibbs sampling
- the history of " "
- the landscape of approximate inference.

- Geman and Geman (1984) : Gibbs for the Ising model

- Gelfand and Smith (1990) : Gibbs for hierarchical Bayesian models

• burn-in, initialization, lag, convergence

$$\pi_0 \rightarrow h_0 \xrightarrow{\pi} h_1 \xrightarrow{\pi} h_2 \rightarrow \dots \xrightarrow{\pi} h_T$$

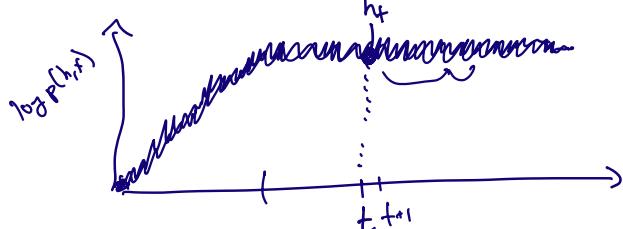
π : transition

π_0 : initial.

$\pi_T \triangleq \max_t \pi(h_t)$

$$\lim_{t \rightarrow \infty} \pi_t(h_t) = p(h | x)$$

monitor $\log p(h, x)$ \in deviance



$$p(h | x) \approx \sum_{b=1}^B \delta(h_b)$$

$$h_b \sim p(h | x)$$