$$p(\text{Data}, \theta) = p(\text{Data}|\theta)\,\boxed{p(\theta)}$$

$$p_j(\text{Data}, \theta) \qquad j \in \{1, \ldots, V\}$$

$$\quad\quad p(\text{Data}|\theta)\, p_j(\theta)$$

$$j^* = \arg\max_j\; p_j(\text{Data}) = \arg\max_j \boxed{\int p(\text{Data}|\theta)\, p_j(\theta)\, d\theta} \quad \checkmark$$

$$p(\theta|\text{Data}) = \frac{p(\text{Data}|\theta)\, p(\theta)}{\displaystyle\int p(\text{Data}|\theta)\, p(\theta)\, d\theta}$$

$$p(\text{Data})$$

Evidence

$$j \in \mathbb{R}$$

$$\log p(\text{Data}) = \int q(\theta) \log p(\text{Data})\, d\theta = \left\{ p(\text{Data}) = \frac{p(\text{Data}, \theta)}{p(\theta|\text{Data})} \right\} =$$

$$= \int q(\theta) \log \frac{p(\text{Data}, \theta)\; q(\theta)}{p(\theta|\text{Data})\; q(\theta)}\, d\theta =$$

$$\boxed{\square + \square} \qquad = \int q(\theta) \log \frac{p(\text{Data}, \theta)}{q(\theta)}\, d\theta + \int q(\theta) \log \frac{q(\theta)}{p(\theta|\text{Data})}\, d\theta \; \geq\; \mathcal{L}(q, \theta)$$

$$\|$$

$$\mathcal{L}(q, \theta)$$

$$\log p_j(\text{Data}) \geq \int q(\theta) \log \frac{p_j(\text{Data}, \theta)}{q(\theta)} d\theta \longrightarrow \underset{q, j}{\overset{\forall q}{\max}}$$

1) Ideal case: $q^*(\theta) = p(\theta | \text{Data}) \iff p(\text{Data}) - \text{tractable}$

2) Mean-field: $q(\theta) = \prod_{j=1}^{m} q_j(\theta_j)$

3) Parametric VI: $q(\theta) = q(\theta | \eta)$

$$\mathcal{L}(q_i j) = \mathcal{L}(\eta, j) = \int q(\theta | \eta) \log \frac{p_j(\text{Data}, \theta)}{q(\theta | \eta)} d\theta \longrightarrow \underset{\eta, j}{\max}$$

$\hookrightarrow$ we don't need to compute it

| Function | Stochastic gradient | Randomness |
|---|---|---|
| $\sum_{i=1}^{n} f_i(x)$ | $\nabla f_j(x)$ | $j \sim U\{1, \dots, n\}$ |
| $\int p(y) f(x,y)\, dy$ | $\frac{\partial}{\partial x} f(x, \hat{y})$ | $\hat{y} \sim p(y)$ |
| $\int p(y)\left[\sum_{i=1}^{n} f_i(x,y)\right] dy$ | $\frac{\partial}{\partial x} f_j(x, \hat{y})$ | $j \sim U\{1, \dots, n\}$ $\hat{y} \sim p(y)$ |
| $\int p(y\mid x) f(x,y)\, dy$ | ~~$\frac{\partial}{\partial x} f(x, \hat{y})$~~ | $\hat{y} \sim p(y\mid x)$ |

$\sum_{i=1}^{n} \int p(y\mid x) f_i(y)\, dy$

$$p(t, \theta \mid x) = p(t \mid \theta, x) \, p(\theta) = \frac{1}{1 + \exp(-t\theta^T x)} \, \mathcal{N}(\theta \mid 0, \Lambda)$$

$$t \in \{-1, +1\} \qquad x, \theta \in \mathbb{R}^d$$

$$\left( X_{tr}, T_{tr} \right) = \{ (x_i, t_i) \}_{i=1}^{n}$$

$$\Lambda = \begin{pmatrix} \lambda_1^2 & \cdots & 0 \\ 0 & \cdots & \lambda_d^2 \end{pmatrix}$$

$$\theta_{MP} = \arg\max_{\theta} \; p(\theta \mid X_{tr}, T_{tr}) = \arg\max_{\theta} \; p(T_{tr} \mid X_{tr}, \theta) \, p(\theta)$$

1) $\lim_{\Lambda \to 0} \theta_{MP} = 0$

2) $\lim_{\Lambda \to +\infty} \theta_{MP} = \theta_{ML}$

$\cancel{n >> d}$

Фиксируем все $\lambda_j = \lambda_0$, выписали байеса для бедных $\implies$

$\implies$ $\ell 2$-регуляр. логистич. регрессии

Классический RUM имеет сложность $\underline{O(d^3)}$

$$p(t, \theta | x) = \underline{p(t | x, \theta)}\, p(\theta | \Lambda)$$

$$\underline{\frac{1}{1 + \exp(-t\theta^T x)}}$$

$$\mathcal{N}(\theta | 0, \Lambda)$$

В идеале $\Lambda^* = \arg\max_{\Lambda} p(T_{tr} | X_{tr}, \Lambda) =$

$$= \arg\max_{\Lambda} \int \underline{p(T_{tr} | X_{tr}, \theta)}\, \underline{p(\theta | \Lambda)}\, d\theta$$

i) $u, d \sim 10 - 10^3$ — Classical RUM

2) $n \gg d$ — Logistic regression with no regularization

3) Пусть $u, d \gg 1$

$$\log p(T_{tr} | X_{tr}, \Lambda) \geq \int q(\theta | \eta) \log \frac{p(T_{tr}, \theta | X_{tr}, \Lambda)}{q(\theta | \eta)}\, d\theta =$$

$$= \int q(\theta | \eta) \log \frac{p(T_{tr} | X_{tr}, \theta)\, p(\theta | \Lambda)}{q(\theta | \eta)}\, d\theta \longrightarrow \max_{\Lambda, \eta}$$

$$\underline{p(\theta | X_{tr}, T_{tr}, \Lambda)} - \text{пол- больше}$$

$$q(\theta | \eta) \approx p(\theta | X_{tr}, T_{tr})$$

$$q(\theta | \eta) = \prod_{j=1}^{d} \mathcal{N}(\theta_j | \mu_j, \sigma_j^2) \qquad p(\theta | \Lambda) = \prod_{j=1}^{d} \mathcal{N}(\theta | 0, \lambda_j^2)$$
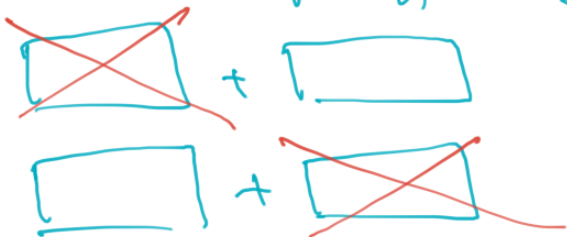
$$\mathcal{L}(\eta, \Lambda) = \int q(\theta | \eta) \log \frac{p(T_{tr} | X_{tr}, \theta) p(\theta | \Lambda)}{q(\theta | \eta)} d\theta =$$

$$= \int q(\theta | \eta) \log p(T_{tr} | X_{tr}, \theta) d\theta + \int q(\theta | \eta) \log \frac{p(\theta | \Lambda)}{q(\theta | \eta)} d\theta =$$

$$= \underbrace{\int q(\theta | \eta) \log p(T_{tr} | X_{tr}, \theta) d\theta}_{\text{Date term}} - \underbrace{KL(q(\theta | \eta) \| p(\theta | \Lambda))}_{\text{Regularisator}} \longrightarrow \max_{\eta, \Lambda}$$
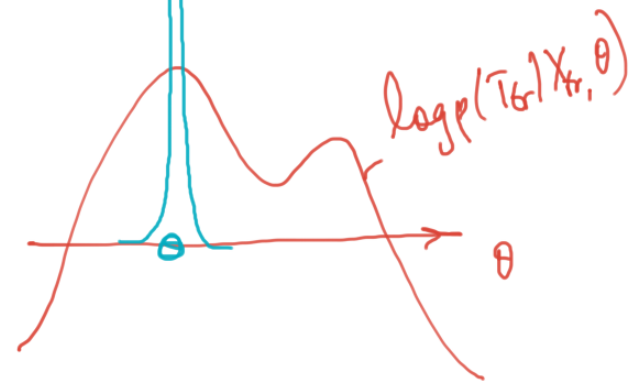
Пусть i) $p(\theta | \Lambda)$ - фиксировано

ii) $q(\theta | \eta)$ - содержит все возм. распр-ия

□̸ + ▱

▢ + ▭̸

$$q(\theta | \eta) = p(\theta | \Lambda)$$

$$q(\theta | \eta) = \delta(\theta - \theta_{ML})$$

$\log p(T_{tr} | X_{tr}, \theta)$

$\theta$

$$p(\Theta|\lambda) = \prod_{j=1}^{d} \mathcal{N}(\theta_j|0, \lambda_j^2) \qquad q(\Theta|\eta) = \prod_{j=1}^{d} \mathcal{N}(\theta_j|\mu_j, \sigma_j^2)$$

$$\max_{\eta, \lambda}$$

$$\eta = \{\vec{\mu}, \vec{\sigma}\}$$

$$\mathcal{L}(\eta, \lambda) = \int q(\Theta|\eta) \log p(T_{tr}|X_{tr}, \Theta) \, d\Theta - \left( \int q(\Theta|\eta) \log \frac{\prod_j q(\Theta|\eta)}{\prod_j p(\Theta|\lambda)} \, d\Theta = \right.$$

tractable

$$= \sum_{i=1}^{n} \int q(\Theta|\eta) \log p(t_i|x_i, \Theta) \, d\Theta - \sum_{j=1}^{d} KL\left( \mathcal{N}(\theta_j|\mu_j, \sigma_j^2) \| \mathcal{N}(\theta_j|0, \lambda_j^2) \right)$$

$$KL\left( q(\theta_j|\eta_j) \| p(\theta_j|\lambda_j) \right) = \log \frac{\lambda_j}{\sigma_j} + \frac{\sigma_j^2 + \mu_j^2}{2\lambda_j^2} \quad \Big| \frac{\partial}{\partial \lambda_j}, \, = 0$$

$$\frac{1}{\lambda_j} - \frac{\sigma_j^2 + \mu_j^2}{\lambda_j^3} = 0 \quad \Longleftrightarrow \quad \boxed{\lambda_{j*}^2 = \sigma_j^2 + \mu_j^2}$$

$$\sum_{j=1}^{d} \log \frac{\sigma_j^2}{\sigma_j^2 + \mu_j^2} + Const$$

$$KL\left( q(\theta_j|\eta_j) \| p(\theta_j|\lambda_{j*}) \right) = \log \frac{\sigma_j^2 + \mu_j^2}{\sigma_j^2} + Const$$

$$\frac{\partial}{\partial \mu}, \frac{\partial}{\partial \sigma}$$

Рассм. 1ое слагаемое в $\mathcal{L}(\eta, \Lambda)$

$$\boxed{\sum_{i=1}^{n} \int q(\theta|\eta) \log p(t_i|x_i, \theta)\, d\theta} \quad \text{©}$$

Reparametrization trick

$$q(\theta|\eta) \longrightarrow r(\varepsilon)$$
$$\theta = g(\varepsilon, \eta)$$

Если $q(\theta|\eta) = \prod_{j=1}^{d} \mathcal{N}(\theta_j|\mu_j, \sigma_j^2)$,

тогда $\theta_j = g(\varepsilon_j, \eta_j) = \underline{\mu_j + \Sigma \sigma_j}$, где $\varepsilon \sim \mathcal{N}(\varepsilon|0, 1)$

$$\text{©} = \sum_{i=1}^{n} \int \boxed{r(\varepsilon)} \log p(t_i|x_i, g(\varepsilon, \eta))\, d\varepsilon \quad | \text{ stoch. grad}$$

$\underbrace{\qquad}_{\text{не зависит от } \eta}$

$j \sim U\{1, \ldots, n\} \quad \tilde{\varepsilon} \sim \mathcal{N}(\varepsilon|0, I)$

$$\frac{\partial}{\partial \eta} \int q(\theta|\eta) \log p(T_{tr}|X_{tr}, \theta)\, d\theta \approx \frac{\partial}{\partial \eta} \log p(t_j|x_j, g(\hat{\varepsilon}, \eta))$$

Нужно $\exists \dfrac{\partial g}{\partial \eta}$

$$\mathcal{L}(\eta) = \sum_{i=1}^{n} \int q(\Theta|\eta) \log p(t_i|x_i, \Theta) \, d\Theta + \frac{1}{2} \sum_{j=1}^{d} \log \frac{\sigma_j^2}{\sigma_j^2 + \mu_j^2} \longrightarrow \max_{\eta}$$

$$\lambda_{j*} = \sigma_j^2 + \mu_j^2 \qquad \underline{\text{Doubly}} \text{ Stochastic Variational Inference}$$

↪ i) Mini-batching

ii) Monte-Carlo estimation of intractable integrals

Значит. вост $\lambda_{j*} \to 0$, м.с. $j^{\text{oui}}$ нейроди вомено исключ вость