

$$p(t, \theta | x, \lambda) = p(t | x, \theta) p(\theta | \lambda)$$

43

✓ 1) Model Selection, λ - ?

✓ 2) Bayesian inference

$$1) \lambda^* = \arg \max_{\lambda} p(T_{tr} | X_{tr}, \lambda) = \arg \max_{\lambda} \int p(T_{tr} | X_{tr}, \theta) p(\theta | \lambda) d\theta$$

$$2) p(\theta | X_{tr}, T_{tr}, \lambda) = \frac{p(T_{tr} | X_{tr}, \theta) p(\theta | \lambda)}{\int p(T_{tr} | X_{tr}, \theta) p(\theta | \lambda) d\theta}$$

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

Variational Bayes

$$\log p(T_{tr} | X_{tr}, \lambda) \geq \mathcal{L}(\eta, \lambda) = \int q(\theta | \eta) \log \frac{p(T_{tr}, \theta | X_{tr}, \lambda)}{q(\theta | \eta)} d\theta \rightarrow \max_{\eta, \lambda}$$

$$\begin{aligned} & \int q(\theta | \eta) \log p(T_{tr} | X_{tr}, \theta) d\theta - \text{KL}(q(\theta | \eta) \| p(\theta | \lambda)) = \\ & = \sum_{i=1}^n \int q(\theta | \eta) \log p(t_i | x_i, \theta) d\theta - \text{KL}(q(\theta | \eta) \| p(\theta | \lambda)) \end{aligned}$$

MC estimation

$$\lambda^* = \arg \max_{\lambda} \mathcal{L} \quad q(\theta | \eta) \approx p(\theta | X_{tr}, T_{tr}, \lambda)$$

$$p(t|x, \lambda) = p(t|x, \theta) p(\theta|\lambda)$$

$$p(t|x, \theta) \sim \mathcal{N}(\mu, \sigma^2)$$

$$\checkmark \quad p(\theta|\lambda) = \prod_{i,j,l} \mathcal{N}(\theta_{ijl} | 0, \lambda_{ijl}^2)$$

$$(X_{\mathcal{G}}, T_{\mathcal{G}}) = \{x_i, t_i\}_{i=1}^n$$

$$\checkmark \quad q(\theta|\eta) = \prod_{i,j,l} \mathcal{N}(\theta_{ijl} | \mu_{ijl}, \sigma_{ijl}^2)$$

$$\eta = \{\mu_{ijl}, \sigma_{ijl}^2\}$$

не забыть
от λ

$$\mathcal{L}(\eta, \lambda) = \int q(\theta|\eta) \log p(T_{\mathcal{G}} | X_{\mathcal{G}}, \theta) d\theta - \text{KL}(q(\theta|\eta) \| p(\theta|\lambda)) =$$

$$= \sum_{i=1}^n \int q(\theta|\eta) \log p(t_i | x_i, \theta) d\theta - \sum_{i,j,l} \text{KL}(\mathcal{N}(\theta_{ijl} | \mu_{ijl}, \sigma_{ijl}^2) \| \mathcal{N}(\theta_{ijl} | 0, \lambda_{ijl}^2))$$

$$\text{KL}(\mathcal{N}(\theta | \mu, \sigma^2) \| \mathcal{N}(\theta | 0, \lambda^2)) = \log \frac{\lambda}{\sigma} + \frac{\sigma^2 + \mu^2}{2\lambda^2} = \left\{ \left(\frac{\sigma}{\lambda} \right)^2 + \frac{\mu^2}{\lambda^2} \right\} =$$

$$= \text{Const} + \frac{1}{2} \log \frac{\mu^2 + \sigma^2}{\sigma^2}$$

$$\sum_{i=1}^n \int q(\theta | \eta) \log p(t_i | x_i, \theta) d\theta = \left\{ \begin{array}{l} RT \\ \theta = g(\varepsilon, \eta) \\ \varepsilon \sim \pi(\varepsilon) \end{array} \right. \left. \begin{array}{l} \theta_{ijl} = \mu_{ijl} + \varepsilon_{ijl} \delta_{ijl} \\ \varepsilon_{ijl} \sim N(\varepsilon_{ijl} | 0, 1) \\ g(\varepsilon, \eta) \end{array} \right\} =$$

$$= \sum_{i=1}^n \int \pi(\varepsilon) \log p(t_i | x_i, g(\varepsilon, \eta)) d\varepsilon$$

$$\text{stoch grad}_{\eta} \left[\sum_{i=1}^n \int \pi(\varepsilon) \log p(t_i | x_i, g(\varepsilon, \eta)) d\varepsilon \right] = \left\{ j \sim U\{1, \dots, n\} \right.$$

$$= n \cdot \text{stoch grad}_{\eta} \int \pi(\varepsilon) \log p(t_j | x_j, g(\varepsilon, \eta)) d\varepsilon =$$

$$= n \int \pi(\varepsilon) \text{stoch grad}_{\eta} \log p(t_j | x_j, g(\varepsilon, \eta)) d\varepsilon = \left\{ \hat{\varepsilon} \sim \pi(\varepsilon) \right\} =$$

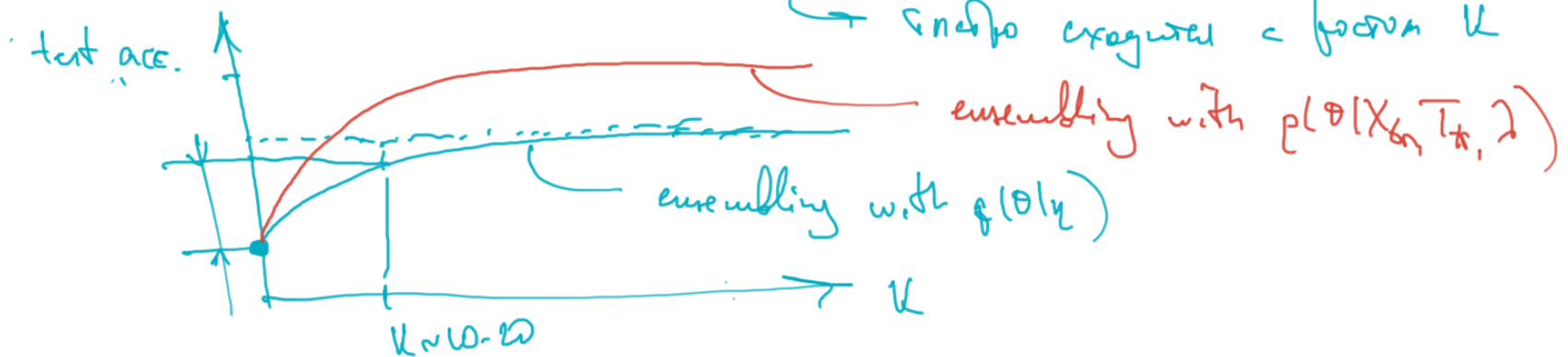
$$= n \cdot \frac{\partial}{\partial \eta} \log p(t_j | x_j, g(\hat{\varepsilon}, \eta)) = n \frac{\partial}{\partial \theta} \log p(t_j | x_j, \theta) \cdot \frac{\partial g(\hat{\varepsilon}, \eta)}{\partial \eta}$$

$$\mathcal{L}(\eta) = \sum_{i=1}^n \int q(\theta|\eta) \log p(T_{tr}|X_{tr}, \theta) d\theta - \frac{1}{2} \sum_{i,j,l} \log \frac{\mu_{ijl}^2 + \sigma_{ijl}^2}{\sigma_{ijl}^2} \rightarrow \max_{\eta}$$

Восмущения $\mu_{ijl} \rightarrow 0$, $\sigma_{ijl} \rightarrow 0$, следовательно $\lambda_{ijl} \rightarrow 0$

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

$$\begin{aligned} 3) \text{ Ensembling } p(t|x, X_{tr}, T_{tr}) &= \int p(t|x, \theta) p(\theta|X_{tr}, T_{tr}) d\theta \approx \\ &\approx \int p(t|x, \theta) q(\theta|\eta) d\theta \approx \frac{1}{K} \sum_{k=1}^K p(t|x, \theta_k), \text{ где } \theta_k \sim \underline{q(\theta|\eta)} \end{aligned}$$





$$y = f\left(\sum_{i=1}^n w_i x_i\right)$$

$$w_i = \tilde{w}_i \cdot z_i$$

$$z_i \sim \text{Ber}(p)$$

- можно
заменить
на нормальное
распределение

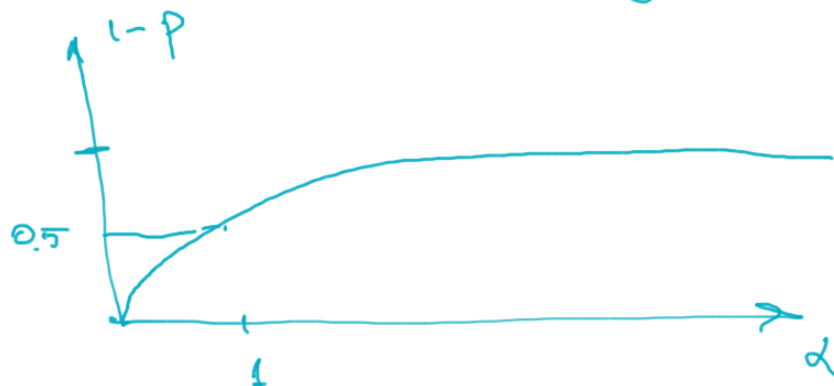
$$v = \sum_{i=1}^m w_i x_i \xrightarrow{m \gg 1} \mathcal{N}(v | \mu, \sigma^2)$$

$$w_i = \tilde{w}_i \cdot \delta$$

$$\delta \sim \mathcal{N}(\delta | 1, \alpha)$$

Существует взаимно-однозначное соответствие между α и p

$$\alpha \approx \frac{1-p}{p}$$



$$L(\eta) = \int q(\theta|\eta) \log p(T_{\text{tr}} | X_{\text{tr}}, \theta) d\theta - \frac{1}{2} \sum_{i,j,l} \log \frac{\mu_{ijl}^2 + \sigma_{ijl}^2}{\sigma_{ijl}^2} =$$

$$= \int \prod_{i,j,l} \mathcal{N}(\theta_{ijl} | \mu_{ijl}, \sigma_{ijl}^2) \log p(T_{\text{tr}} | X_{\text{tr}}, \theta) d\theta - \frac{1}{2} \sum_{i,j,l} \log \frac{\mu_{ijl}^2 + \sigma_{ijl}^2}{\sigma_{ijl}^2}$$

$$= \left\{ \begin{array}{l} \sigma_{ijl}^2 = \alpha_{ijl} \mu_{ijl}^2 \\ \alpha_{ijl} = \frac{\sigma_{ijl}^2}{\mu_{ijl}^2} \end{array} \right\} = \int \prod_{i,j,l} \mathcal{N}(\theta_{ijl} | \mu_{ijl}, \alpha_{ijl} \mu_{ijl}^2) \log p(T_{\text{tr}} | X_{\text{tr}}, \theta) d\theta - \frac{1}{2} \sum_{i,j,l} \log \frac{1 + \alpha_{ijl}}{\alpha_{ijl}}$$

Bycym, bce $\alpha_{ijl} = \alpha$ - functionau

we get μ_{ijl} or μ_{ijl} !

$$\max_{\mu} \left[\int \prod_{i,j,l} \mathcal{N}(\theta_{ijl} | \mu_{ijl}, \alpha \mu_{ijl}^2) \log p(T_{\text{tr}} | X_{\text{tr}}, \theta) d\theta - \frac{1}{2} \sum_{i,j,l} \log \frac{1 + \alpha_{ijl}}{\alpha_{ijl}} \right] =$$

$$= \max_{\mu} \left[\int \prod_{i,j,l} \mathcal{N}(\theta_{ijl} | \mu_{ijl}, \alpha \mu_{ijl}^2) \log p(T_{\text{tr}} | X_{\text{tr}}, \theta) d\theta \right] = \left\{ \begin{array}{l} \theta_{ijl} = \mu_{ijl} \delta_{ijl} \\ \delta_{ijl} \sim \mathcal{N}(\delta_{ijl} | 1, \alpha) \end{array} \right\}$$

$$= \max_{\mu} \left[\int \prod_{i,j,l} \mathcal{N}(\delta_{ijl} | 1, \alpha) \log p(T_{\text{tr}} | X_{\text{tr}}, \mu \cdot \delta) d\delta \right]$$

$$\underline{D(\mu)} = \int \prod_{i,j,l} N(\delta_{ijl} | 1, \alpha) \log p(T_{ijl} | X_{ijl}, \mu, \delta) d\delta \rightarrow \max_{\mu} \quad | \text{stoch. grad}_{\mu}$$

↪ не зависит от μ

$$j \sim U\{1, \dots, n\}$$

$$\hat{\delta}_{ijl} \sim N(\delta_{ijl} | 1, \alpha)$$

$$\text{stoch grad}_{\mu} D(\mu) = n \frac{\partial}{\partial \mu} \log p(t_j | x_j, \underline{\hat{\delta}}) =$$

$$= n \frac{\partial}{\partial \theta} \log p(t_j | x_j, \theta) \cdot \frac{\partial \theta}{\partial \mu} = \text{совпадает со стх. градиентом в гауссовском пространстве}$$

Гауссовский процесс со стх. параметрами рассматриваемый как процесс в μ , где

$$\underline{g}(\theta | \mu) = \prod_{i,j,l} N(\theta_{ijl} | \mu_{ijl}, \alpha \mu_{ijl}^2) \text{ где } \mu = \vec{\mu}, \text{ а } \alpha - \text{дисперсия}$$

нашем параметрическом пространстве μ_{ijl} для каждого θ_{ijl}

$\mathcal{L}(\mu, \alpha) \rightarrow \max_{\mu}$ \Leftrightarrow лагранжевы множители α функции \mathcal{L}

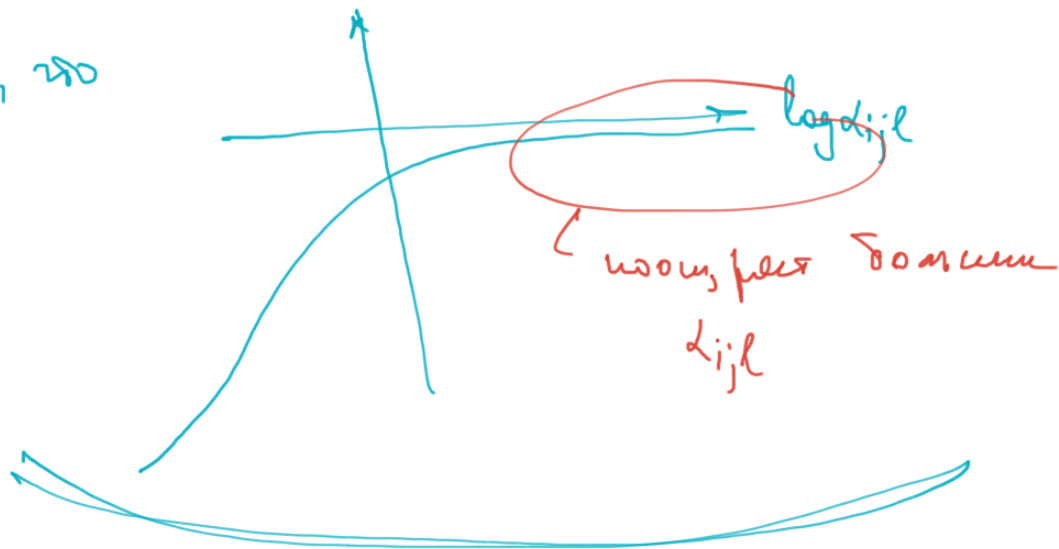
$\mathcal{L}(\mu, \alpha) \rightarrow \max_{\mu, \alpha}$ \Leftrightarrow Sparse Variational Dropout

3 умнож. на $d_{ijl} \rightarrow +\infty$

большая константа $\mathcal{L}(\mu, \alpha)$ умнож. $\log - \frac{1}{2} \sum_{ijl} \log \frac{1 + \alpha_{ijl}}{d_{ijl}}$

\rightarrow если $\mu_{ijl} \rightarrow 0$ так, что

$$\sigma_{ijl}^2 = d_{ijl} \mu_{ijl}^2 \rightarrow 0$$



Новая постановка регуляризации

$$\log p(T_{tr} | X_{tr}, \theta) + \lambda R(\theta) \rightarrow \max_{\theta}$$

$$\theta_{\max} = \arg \max_{\theta} \frac{p(T_{tr} | X_{tr}, \theta) p(\theta)}{\int p(T_{tr} | X_{tr}, \theta) p(\theta) d\theta} = \arg \max_{\theta} p(T_{tr} | X_{tr}, \theta) p(\theta)$$

где $p(\theta) = \frac{1}{A} \exp(\lambda R(\theta))$

Классика

Современное
представление

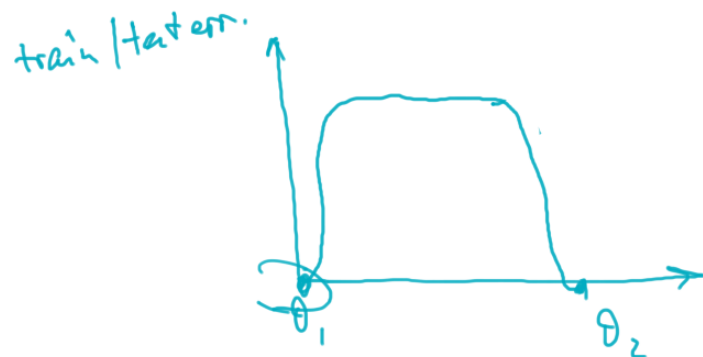
$$\int r(\varepsilon) \log p(T_{tr} | X_{tr}, g(\varepsilon, \theta)) d\varepsilon \rightarrow \max, \text{ где}$$

$g(\varepsilon, \theta)$ — замученная версия θ из θ и ε

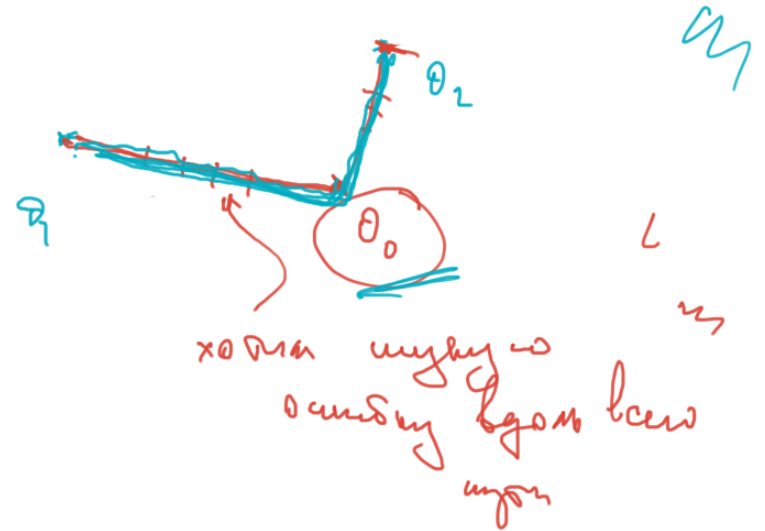
Вуз $\textcircled{r(\varepsilon)}$ — задается пользователем в байесовских моделях

Введение байесовской модели позволяет избежать множества ошибок

θ_1, θ_2 - global maximum minimum $-\log p(T_{tr} | X_{tr}, \theta)$



$$\theta(t; \theta_0) = \begin{cases} (1-t)\theta_1 + t\theta_0, & t \in [0, 1] \\ (2-t)\theta_0 + (t-1)\theta_2, & t \in [1, 2] \end{cases}$$



$$\Phi(\theta_0) = \int_0^2 \underline{r(t)} \log p(T_{tr} | X_{tr}, \theta(t; \theta_0)) dt \rightarrow \max_{\theta_0} \quad r(t) = U[0, 2]$$

$$\text{stoch grad } \underline{\Phi(\theta_0)} = \sim \frac{1}{n} \log p(y_j | x_j, \theta) \cdot \frac{\partial \theta(t; \theta_0)}{\partial \theta_0} \quad \begin{aligned} & \mathbb{I} \sim U[0, 2] \\ & j \sim U\{1, \dots, n\} \end{aligned}$$