

Latent Dirichlet Allocation (LDA)

Data: $w_{d1} \dots w_{dN_d}$ "tokens" in document d

Model:

for topic $k = 1 \dots K$:

$$\beta_k \sim \text{Dir}(\beta_1 \dots \beta_V)$$

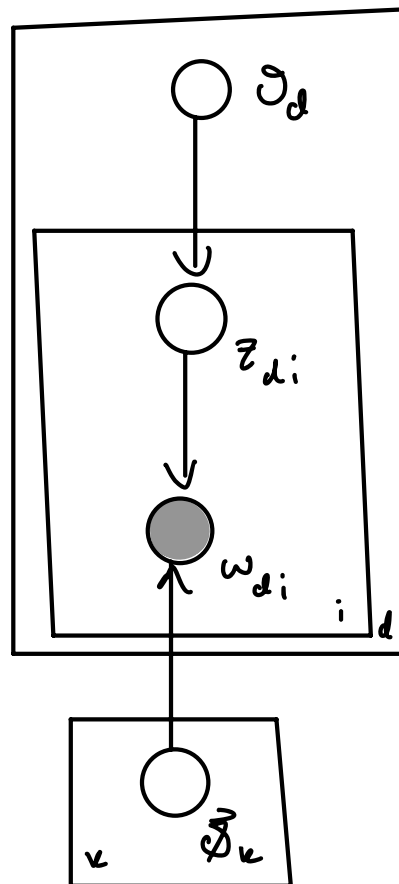
for doc $d = 1 \dots D$:

$$\vec{\theta}_d \sim \text{Dir}(\alpha_1 \dots \alpha_K)$$

for token i in $1 \dots N_d$:

$$z_{di} \sim \text{Lat}(\vec{\theta}_d)$$

$$w_{di} \sim \text{Lat}(\vec{\phi}_{z_{di}})$$



Complete conditionals

$$p(z_{di} = k \mid w_{di} = v, \dots) \propto p(z_{di} = k) p(w_{di} = v \mid z_{di} = k, \dots)$$

$$\propto \theta_{dk} \phi_{kv}$$
$$\rightarrow = \frac{\theta_{dk} \phi_{kv}}{\sum_j \theta_{dj} \phi_{jv}}$$

$$p(\vec{\theta}_d \mid -) \propto \text{Dir}(\vec{\theta}_d; \vec{\alpha}) \prod_i \text{Lat}(z_{di}; \vec{\theta}_d)$$

$$\hookrightarrow = \text{Dir}(\vec{\theta}_d; \tilde{\alpha}_{d1} \dots \tilde{\alpha}_{dK}) = \gamma_{dk}$$

where $\tilde{\alpha}_{dj} = \alpha_j + \sum_{i=1}^{N_d} \mathbb{1}(z_{di} = j)$

$$p(\vec{\phi}_n | -) \propto \text{Dir}(\vec{\phi}_n; \vec{\beta}) \prod_d \prod_i \left[\text{Cat}(w_{di}; \vec{\phi}_n) \right]^{1(z_{di}=k)}$$

$$\hookrightarrow = \text{Dir}(\vec{\phi}_n; \tilde{\beta}_1 \dots \tilde{\beta}_V) \equiv \gamma_{kv}$$

where $\tilde{\beta}_{kv} = \beta_v + \sum_d \sum_i 1(w_{di}=v) 1(z_{di}=k)$

Gibbs sampling

for iter $s = -B \dots -1, 0, 1 \dots ST$:

① for doc d :

for token i :

B "burn-in" iterations
S posterior samples
"thinning" every T samples

$$z_{di}^s \sim \text{Cat} \left(\frac{\theta_{d1} \phi_{1w_{di}}}{\sum_j \theta_{dj} \phi_{jw_{di}}} \dots \frac{\theta_{dK} \phi_{Kw_{di}}}{\sum_j \theta_{dj} \phi_{jw_{di}}} \right)$$

for k, v :

$$\gamma_{dkv}^s = \sum_{i=1}^{N_d} 1(z_{di}^s = k) 1(w_{di} = v)$$

② for doc d :

$$\vec{\theta}_d^s \sim \text{Dir}(\alpha_1 + \sum_v \gamma_{d1v}^s \dots \alpha_K + \sum_v \gamma_{dKv}^s)$$

③ for topic k :

$$\vec{\phi}_k^s \sim \text{Dir}(\beta_1 + \sum_d \gamma_{d1k}^s \dots \beta_V + \sum_d \gamma_{dVk}^s)$$

Notice that the sufficient statistics y_{dv} are invariant to the token ordering $i = 1 \dots Nd$.

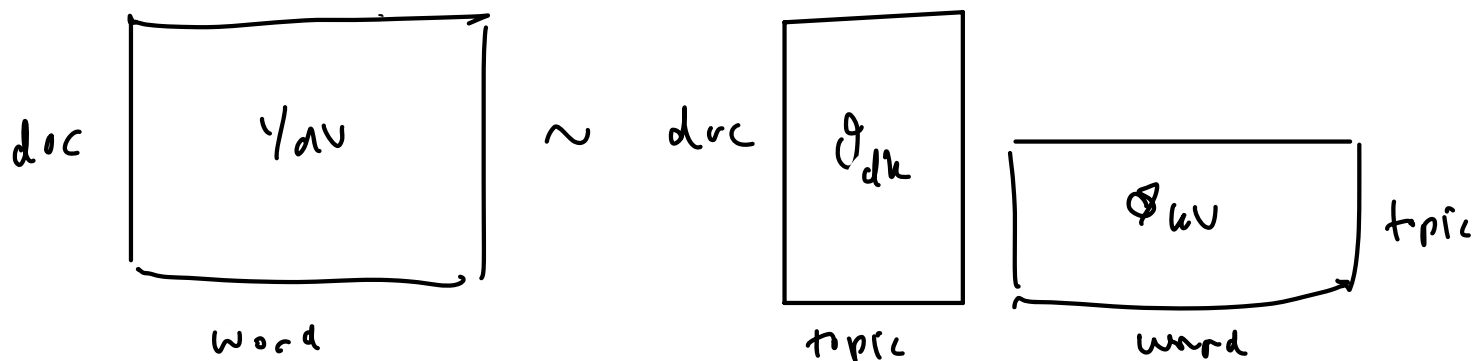
An equivalent way to sample is:

$$\begin{aligned} & \text{for doc } d = 1 \dots D \\ & \quad \text{for word } v = 1 \dots V : \\ & \quad (y_{dv1} \dots y_{dvK}) \sim \text{Mult}(y_{dv} ; \frac{\theta_{d1} \phi_{1v}}{\sum_j \theta_{dj} \phi_{jv}} \dots \frac{\theta_{dK} \phi_{Kv}}{\sum_j \theta_{dj} \phi_{jv}}) \\ & \quad \quad \quad \downarrow \\ & \quad \quad \quad = \sum_i \mathbb{1}(w_{di} = v) \quad (\text{observed}) \end{aligned}$$

Notice that this is exactly the same auxiliary variable sampling step we derived for Poisson matrix factorization

$$y_{dv} \sim \text{Pois}(\sum_j \theta_{dj} \phi_{jv}) \quad \leftarrow$$

LDA is "just" a (very) special case of NMF.



↑ each doc is a "bag of words".

Another way to view LDA is as an admixture or mixed membership model. Each document is part of multiple mixture components ("topics").

If $\vec{\theta}_d$ were a "one-hot" vector, we would recover the mixture model.

Label switching

- In LDA and other forms of AMF, there is no inherent ordering of the topics $1 \dots k$:

$$\sum_k \theta_{dk} \phi_{kv} = \sum_k \theta_{d\Delta(k)} \phi_{\Delta(k)v}$$

where $\Delta(k)$ is a permutation of indices $1 \dots k$

- During MCMC, the "labels" can switch
- True for most (ad)mixtures
- This means you should not average any quantity indexed by k across posterior samples:

$$\text{e.g. } \underbrace{\frac{1}{S} \sum_s \phi_{kv}^s}_{\text{bad}} \quad \text{vs.} \quad \frac{1}{S} \sum_s \sum_k \theta_{dk}^s \phi_{kv}^s$$

good

Variational Inference (VI)

Setting:

data: $x_1, \dots, x_n \equiv x$

latents: z_1, \dots, z_D

posterior: $p(z | x)$
(intractable)

We have seen one way to approximate the posterior using Gibbs sampling:

$$z_d^s \sim p(z_d | z_{-d}, x)$$

VI is another way...

Define a family of approximate densities:

$$q(z) \in \mathcal{Q}$$

find the member that minimizes:

$$q^*(z) = \underset{q \in \mathcal{Q}}{\text{argmin}} \text{KL}(q(z) \parallel p(z | x))$$

How do you minimize if you cannot compute $p(z | x)$?

$$\text{KL}(q \parallel p) = \mathbb{E}_q \left[\log \frac{q(z)}{p(z | x)} \right]$$

$$= \underbrace{\mathbb{E}_q \left[\log \frac{q(z)}{p(z, x)} \right]}_{-\text{ELBO}(q)} + \underbrace{\log p(x)}_{\log \text{ evidence}}$$

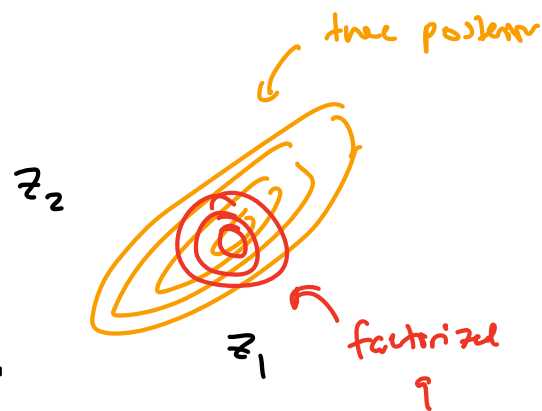
$$q^*(z) = \underset{q}{\text{argmax}} \text{ELBO}(q)$$

Define the family \mathcal{Q} to be "convenient".

Factorized family:

$$q(z) = \prod_{d=1}^D q_d(z_d) \equiv \prod_{d=1}^D q(z_d; \lambda_d)$$

↑
"variational parameter"



Coordinate ascent VI (CAVI)

repeat:

for $d = 1 \dots D$:

$$q^*(z_d) = \underset{q(z_d)}{\text{argmax}} \text{ELBO}(q)$$

fact (Bishop 2006):

$$q^*(z_d) \propto \exp\left(\mathbb{E}_{z_{-d}} [\log p(z_d | z_{-d}, x)]\right)$$

If $p(z_d | z_{-d}, x)$ is exp/con, this will often be "nice".

Proof: $\text{KL}(q(z_d) \parallel q^*(z_d))$

$$= \mathbb{E}_{q(z_d)} [\log q(z_d) - \log q^*(z_d)]$$

$$\propto \mathbb{E}_{q(z_d)} [\log q(z_d)] - \mathbb{E}_{q(z_d)} \left[\mathbb{E}_{q(z_{-d})} [\log p(z_d | z_{-d}, x)] \right]$$

$$\propto -H(q) - \mathbb{E}_q [\log (z, x)]$$

$$\propto -\text{ELBO}$$

Minimizing the KL also maximizes the ELBO.

CAVI for LDA

$$q(\dots) = \left[\prod_d \prod_v q(y_{dv} \dots y_{dvk}) \right] \prod_d q(\vec{\theta}_d) \prod_k q(\vec{\phi}_k)$$

$$q^*(\vec{\theta}_d) \propto \exp \left(\mathbb{E}_{q(\dots)} \left[\log p(\vec{\theta}_d | \dots) \right] \right)$$

$$\propto \exp \left(\mathbb{E} \left[\log \text{Dir}(\vec{\theta}_d; \alpha_1 + \sum_v y_{dv1} \dots \alpha_k + \sum_v y_{dvk}) \right] \right)$$

$$\propto_{\vec{\theta}_d} \exp \left(\mathbb{E} \left[\log \prod_k \theta_{dk}^{\alpha_k + \sum_v y_{dvk} - 1} \right] \right)$$

$$\propto_{\vec{\theta}_d} \prod_k \theta_{dk}^{\alpha_k + \sum_v \mathbb{E}[y_{dvk}] - 1}$$

$$\Downarrow$$
$$q^*(\vec{\theta}_d) = \text{Dir}(\vec{\theta}_d; \lambda_{d1} \dots \lambda_{dk}) \quad (1)$$

$$\text{where } \lambda_{dk} = \alpha_k + \sum_v \mathbb{E}[y_{dvk}]$$

Similarly:

$$q^*(\vec{\phi}_k) = \text{Dir}(\vec{\phi}_k; \lambda_{k1} \dots \lambda_{kv}) \quad (2)$$

$$\text{where } \lambda_{kv} = \beta_v + \sum_d \mathbb{E}[y_{dvk}]$$

$$q^{\alpha}(y_{du1} \dots y_{duk}) \propto \exp(\log \mathbb{E}_Q[\text{Mult}(\dots)])$$

$$\propto \exp\left(\mathbb{E}\left[\log \mathcal{S}(\dots) \frac{y_{du!}}{\prod_k y_{duk}!} \prod_n \pi(\dots)^{y_{duk}}\right]\right)$$

$$\propto \mathcal{S}(\dots) \frac{y_{du!}}{\prod_n y_{duk}!} \exp\left(\mathbb{E}\left[\log \prod_n \left(\frac{\vartheta_{dk} \vartheta_{kv}}{\vartheta_d^T \vartheta_v}\right)^{y_{duk}}\right]\right)$$

$$\propto \mathcal{S}(\dots) \frac{y_{du!}}{\prod_n y_{duk}!} \prod_n \underbrace{\exp(\mathbb{E}[\log(\vartheta_{dk} \vartheta_{kv})])^{y_{duk}}}$$

$$| \equiv \zeta_Q[\vartheta_{dk} \vartheta_{kv}] = \zeta_Q[\vartheta_{dk}] \zeta_Q[\vartheta_{kv}]$$

Def 1

$$\vec{x} \sim \text{Dir}(\alpha_1, \dots, \alpha_n)$$

$$\mathbb{E}[\ln x_v] = \psi(x_v) - \psi\left(\sum_v x_v\right)$$

← digamma function

$$\hookrightarrow q^{\alpha}(y_{du}) = \text{Mult}(y_{du}; \vec{p}_{du})$$

$$\text{where } p_{duk} = \frac{\zeta(\vartheta_{dk}) \zeta(\vartheta_{kv})}{\sum_j \zeta(\vartheta_{dj}) \zeta(\vartheta_{jv})} \quad (3)$$

CAVI for LDA

initialize variational parameters λ_{dh} , λ_{kv} , p_{dvh} and
calculate initial expectations $\mathbb{E}_q[y_{dvh}]$, $G[\theta_{dh}]$, $G[\theta_{kv}]$

repeat until convergence:

① update $q(\vec{y}_{dv}) = \text{Mult}(\vec{y}_{dv}; \vec{y}_{dv}, \vec{p}_{dv})$

$$p_{dvh} = \frac{G_q[\theta_{dh}] G_q[\theta_{kv}]}{\sum_{j=1}^K G_q[\theta_{dj}] G_q[\theta_{jv}]} \quad \leftarrow \text{only need to do this for } y_{dv} > 0$$

$$\mathbb{E}_q[y_{dvh}] = y_{dv} p_{dvh}$$

② update $q(\vec{\theta}_d) = \text{Dir}(\vec{\theta}_d; \vec{\lambda}_d)$

$$\lambda_{dh} = \alpha_h + \sum_v \mathbb{E}_q[y_{dv} h]$$

$$G_q[\theta_{dh}] = \exp(\psi(\lambda_{dh}) - \psi(\sum_j \lambda_{dj}))$$

③ update $q(\vec{\theta}_k) = \text{Dir}(\vec{\theta}_k; \vec{\lambda}_k)$

$$\lambda_{kv} = \beta_v + \sum_d \mathbb{E}_q[y_{dv} k]$$

$$G_q[\theta_{kv}] = \dots$$

④ Calculate $E(\text{Bo}(q))$; assess convergence

ELBO:

$$\mathbb{E}_q \left[\log \frac{p(y, \vec{y}, \theta, \phi)}{q(\vec{y}, \theta, \phi)} \right]$$

$$= \mathbb{E}_q [\log p(y, \vec{y} | \theta, \phi)] + \mathbb{E}_q [\log p(\theta, \phi)]$$

"complete data"
likelihood
+ $H_q[q(\dots)]$
Prior

↑ entropy of variational dist

e.g. $\mathbb{E}_q [\log p(\theta)] = \sum_d \mathbb{E}_q [\log p(\vec{\theta}_d)]$

$$\mathbb{E}_q [\log \text{Pr}(\vec{\theta}_d; \mathbf{z})] = \mathbb{E}_q \left[\log \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_k \theta_{dk}^{\alpha_k - 1} \right]$$

$$= \log \frac{\Gamma(\dots)}{\prod_k \Gamma(\dots)} + \sum_k (\alpha_k - 1) \mathbb{E}_q [\log \theta_{dk}]$$

Recall that $q(\vec{\theta}_d) = \text{Dir}(\vec{\theta}_d; \vec{\lambda}_d)$ and that

$$\mathbb{E}_q [\ln \theta_{dk}] = \psi(\lambda_{dk}) - \psi(\sum_j \lambda_{dj})$$

e.g. $H_q[q(\dots)] = \sum_d H[q(\vec{\theta}_d)] + \dots$

↑ entropy of $\text{Dir}(\vec{\theta}_d; \vec{\lambda}_d)$

The entire ELBO can be (carefully) derived and calculated.

CAV) for conditionally conjugate models

$$p(\eta, z, x) = p(\beta) \prod_i p(z_i, x_i | \eta)$$

- η are global params
- z are local latents
- x are data

complete data likelihood:

$$p(z_i, x_i | \eta) = h_e(z_i, x_i) \exp(\eta^T t_e(z_i, x_i) - A_e(\eta))$$

$$\eta \equiv \eta(\theta)$$

(natural parameterization)

$$p(\eta | \alpha) = F(\eta; \alpha) \propto_{\eta} h_c(\eta) \exp(\alpha^T t_c(\eta))$$

$$t_c(\eta) = [\eta, -A_e(\eta)]$$

$$p(\eta | -) \propto h_c(\eta) \exp(\hat{\alpha}^T t_c(\eta))$$

$$\hat{\alpha} = \begin{bmatrix} \alpha_1 + \sum_{i=1}^n t_e(z_i, x_i) \\ \alpha_2 + n \end{bmatrix}$$

↳

natural parameter of complete conditional

$$= F(\eta; \hat{\alpha}) \quad (\text{same family})$$

Optimal Variational family:

$$q^*(\eta) \propto_{\eta} \exp \left(\mathbb{E}_q [\log p(\eta | -)] \right)$$

$$\propto_{\eta} h_c(\eta) \exp \left(\mathbb{E}_q [\hat{\alpha}]^T t_c(\eta) \right)$$

$$\hookrightarrow = F(\eta; \mathbb{E}_q [\hat{\alpha}])$$

↑
variational natural parameter λ

So $q^*(\eta; \lambda)$ is the same family as the prior $F(\eta; \alpha)$ and complete cond. $F(\eta; \hat{\alpha})$.

Taking gradients:

$$\nabla_{\lambda} \text{ELBO} = \underbrace{A_c''(\lambda)}_{\text{Hessian of the log normalizer}} \left(\mathbb{E}_q [\hat{\alpha}] - \lambda \right)$$

Hessian of the log normalizer

$$\text{i.e. } A_c''(\lambda)_{ij} = \frac{\partial^2 A_c(\lambda)}{\partial \lambda_i \partial \lambda_j}$$

So $\lambda = \mathbb{E}_q [\hat{\alpha}]$ maximizes ELBO.

Stochastic Variational Inference

Setting: large n

Stochastic optimization:

$$\text{|| for iter } t: \\ \lambda_t = \lambda_{t-1} + \zeta_t \widehat{\nabla_{\lambda} \text{ELBO}}$$

if $\widehat{\nabla_{\lambda} \text{ELBO}}$ is an unbiased estimate of the gradient:

$$\mathbb{E}[\widehat{\nabla_{\lambda} \text{ELBO}}] = \nabla_{\lambda} \text{ELBO}$$

and if the steps follow the Robbins-Munroe conditions:

$$\sum_t \zeta_t = \infty ; \sum_t \zeta_t^2 < \infty$$

$$\text{e.g. } \zeta_t = t^{-\kappa}, \kappa \in (1/2, 1)$$

then λ_t will reach a (local) optimum of ELBO.

If the gradient had the following structure:

$$\nabla_{\lambda} \text{ELBO} = \sum_{i=1}^n f(x_i, \dots)$$

then a sub-sampled estimate would be unbiased:

$$i \sim \text{Uniform}(1 \dots n)$$

$$\widehat{\nabla_{\lambda} \text{ELBO}} = n f(x_i, \dots)$$

Unfortunately the Euclidean gradient $\nabla_{\lambda} \text{ELBO}$ does not.

Hoffman et al. (2013) showed:

$$\nabla_{\lambda} \text{ELBO} = A_c''(\lambda) \underbrace{\left(\mathbb{E}_q[\hat{z}] - \lambda \right)}_{\triangleq g(\lambda)}$$

natural gradient of the ELBO

$$g(\lambda) = [A_c''(\lambda)]^{-1} \nabla_{\lambda} \text{ELBO}$$

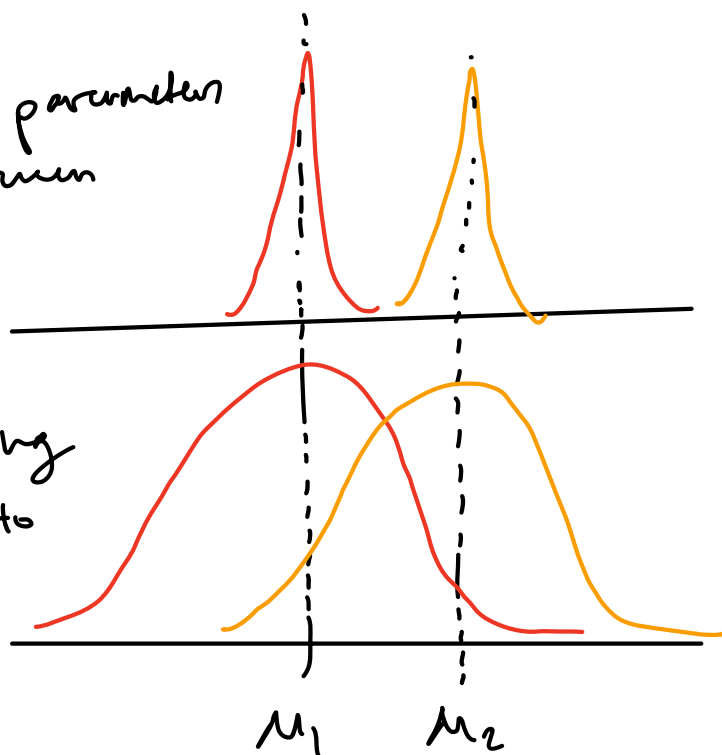
the Euclidean gradient preconditioned by inverse $A_c''(\lambda)$.

$$= \mathbb{E}_q[\hat{z}] - \lambda$$

$$= \begin{bmatrix} \alpha_1 + \sum_i \mathbb{E}_q[t(x_i, z_i)] \\ \alpha_2 + n \end{bmatrix} - \begin{bmatrix} \lambda_1 \\ \lambda_2 \end{bmatrix}$$

Euclidean distance between parameters
 \neq statistical distance between
distributions.

The natural gradient is
defined by a local rescaling
of the Euclidean gradient to
account for this.



SVI

$$\lambda_t = \lambda_{t-1} + \epsilon_t g(\lambda_{t-1})$$

$$= \lambda_{t-1} + \epsilon_t (\mathbb{E}_q[\hat{z}] - \lambda_{t-1})$$

$$\hookrightarrow \lambda_t = (1 - \epsilon_t) \lambda_{t-1} + \epsilon_t \mathbb{E}_q[\hat{z}]$$

replace this with an unbiased estimate

$$i \sim \text{Uniform}(1 \dots n)$$

$$\hat{\mathbb{E}}_q[\hat{z}] = \begin{bmatrix} \alpha_1 \\ \alpha_2 \end{bmatrix} + \begin{bmatrix} n \mathbb{E}_q[t(z_i, x_i)] \\ n \end{bmatrix}$$

Note that $\mathbb{E}_q[t(z_i, x_i)] \equiv \mathbb{E}_{q(z_i)}[t(z_i, x_i)]$

SVI :

for iter t :

$$i \sim \text{Uniform}(1 \dots n)$$

① update local latent

$$q^*(z_i) = \dots$$

② update global param

$$\lambda_t = (1 - \epsilon_t) \lambda_{t-1} + \epsilon_t \left(\alpha + \begin{bmatrix} n \mathbb{E}_q[t(z_i, x_i)] \\ n \end{bmatrix} \right)$$