

Machine learning problem

Input data
↓
 $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$

Decision theory : Minimize $E[L(f(x), y)]$

Prediction function
↓

↳ \hat{y} output

Optimize decision when expected difference between prediction function $f(x)$ and actual value y is minimized

Key quantity $p(\hat{y}|x)$ vs. y for x

Generative model assumption :

Data is generated from a process

⇒ Model the data generating process itself → Richer model

Joint distribution : $p(x, y)$

Factorization :

$$p(x, y) = p(y|x)p(x) \Rightarrow$$

$$p(y|x) = \frac{p(x, y)}{p(x)}$$

Predictive distribution

or

$$p(x, y) = p(x|y)p(y)$$

Generative distribution

② Sample x given y ① Sample y

Marginal distribution (w.r.t θ) for predictive distribution

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

Posterior dist. GIVEN a parameter

Assuming parameters independent of input x

$$= \frac{\int p(y|\theta, x, D) p(\theta|D) d\theta}{\text{Easy}} \quad \text{Hard} \quad \text{Posterior dist. of (all) parameters}$$

$\int \dots d\theta$ intractable

\Rightarrow Approx. integral w. "Deterministic approximation" approach

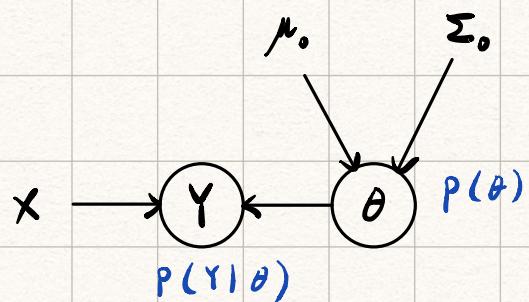
- Laplace approx. (Distr. of $\theta \Rightarrow$ Gaussian approx.)
- Variational methods (Find approx. using optimization)

Estimating parameter distribution using Bayesian inference

$$p(\theta|D) = \frac{p(\theta, D)}{p(D)} = \frac{p(D|\theta)p(\theta)}{p(D)} \quad \text{i.e. Data } \{x_1, \dots, x_n\} \text{ generated by modeled generative process}$$

Likelihood of data: $p(D|\theta) = p(x_1, \dots, x_n|\theta) = \prod_{i=1}^n [p(x_i|\theta)]$

$$p(\theta|D) = \frac{\prod_{i=1}^n [p(x_i|\theta)] p(\theta)}{p(D)} \quad \text{Prior parameter distribution}$$



Bayesian Logistic Regression

- Want the posterior distribution for θ which maximizes the likelihood of the data $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$

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

"Brave" step to skip the normalization term

(Not important as $p(\theta|D)p(D)$
eventually approx. as Gaussian)

- Goal: Get an expression for $p(\theta|D) \propto \underline{p(D|\theta)p(\theta)}$?

- Assume a prior parameter distribution that is Gaussian

Parameter Vector \downarrow Hyperparameters ("Known") \downarrow
 $p(\theta) = N(\theta | \mu_0, \Sigma_0)$

- Assume logistic regression formulation for likelihood function

$$\begin{aligned}
 p(y|\theta) &= \prod_{i=1}^n p(L_i=1|x_i, \theta)^{y_i} p(L_i=0|x_i, \theta)^{1-y_i} \\
 &= \prod_{i=1}^n f(x_i, \theta)^{y_i} (1 - f(x_i, \theta))^{1-y_i} \\
 &= \prod_{i=1}^n (\hat{y}_i)^{y_i} (1 - \hat{y}_i)^{1-y_i} \quad \text{Data 'x' is subsumed} \\
 &\quad \text{by output '}\hat{y}\text{'}
 \end{aligned}$$

$\hookrightarrow \hat{y}_i = \sigma(\theta^\top x_i)$

$$p(\theta|D) \propto \prod_{i=1}^n (\hat{y}_i)^{y_i} (1 - \hat{y}_i)^{1-y_i} \cdot N(\theta | \mu_0, \Sigma_0)$$

- Approximation approach: Laplace approximation (i.e. Gaussian approx.)

- Gaussian w. mean = mode

Variance = maximize likelihood
at data

- Mode of $P(\theta|D) \rightarrow$ MAP solution θ_{MAP}

$$\log p(\theta|D) \propto \sum_{i=1}^n [y_i \log \hat{y}_i + (1 - y_i) \log (1 - \hat{y}_i)]$$

$$-\frac{1}{2}(\theta - \mu_0)^T \Sigma_0^{-1} (\theta - \mu_0) + \text{const.}$$

$$\Rightarrow \theta_{\text{MAP}} = f(\mu_0, \Sigma_0, D)$$

2. Variance of $P(\theta | D) \rightarrow$ Inverse of negative log likelihood
Hessian matrix (second derivative)

$$\Rightarrow \Sigma_n = -\nabla^2 \log P(\theta | D) = \Sigma_0^{-1} + \sum_{i=1}^n \hat{y}_i(1 - \hat{y}_i) x_i x_i^T$$

• Finish: Gaussian approximation to posterior distr. $p(\theta | D)$

*

$$f(\theta) = \mathcal{N}(\theta | \theta_{\text{MAP}}, \Sigma_n)$$

Predictive distribution for Bayesian Logistic Regression

$$p(Y = C_1 | x^{\text{new}}, D) = \int p(C_1, \theta | x^{\text{new}}, D) d\theta$$

$$= \int p(C_1 | x^{\text{new}}, \theta) p(\theta | D) d\theta$$

$$\approx \int \sigma(\theta^T x^{\text{new}}) f(\theta) d\theta$$

Integrating this expression using Dirac function substitution

$$\sigma(\theta^T x^{\text{new}}) = \int \delta(a - \theta^T x^{\text{new}}) \sigma(a) da \quad \begin{matrix} \langle \text{Contusing} \\ \text{part ...} \end{matrix}$$

$$\int \sigma(\theta^T x^{new}) q(\theta) d\theta = \int \sigma(a) p(a) da$$

$$p(a) = \int \delta(a - \theta^T x^{new}) q(\theta) d\theta \equiv \int q(\theta) d\theta \text{ when } a = \theta^T x^{new}$$

$$\mathbb{E}[a] = \int p(a) a da = \int q(\theta) \theta^T x^{new} d\theta = \boxed{\theta_{MAP}^T x^{new}} = \mu_a$$

$$\text{var}[a] = \int p(a) (a^2 - \mathbb{E}[a]^2) da = \int q(\theta) [(\theta^T x^{new})^2 - (\theta_{MAP}^T x^{new})^2] d\theta$$

$$= \boxed{(x^{new})^T \sum_n x^{new} = \sigma_a^2}$$

Variational approximation to the predictive distribution:

$$p(Y=c_1 | x^{new}, D) \approx \int \sigma(a) N(a | \mu_a, \sigma_a^2) da$$

Sigmoid function CANNOT be integrated

analytically w. Gaussian

\Rightarrow Approximate $\sigma(a)$ by Probit function $\Phi(\lambda \cdot a)$
(which CAN be integrated analytically!)

- What is λ ? Coefficient for adjusting approximation

Requiring $\sigma(a)$, $\Phi(\lambda a)$ to have same derivative
(i.e. slope) at origin

$$\Rightarrow \lambda^2 = \frac{1}{8} \pi$$

Inegration

$$\int \underline{\Phi(\lambda \cdot a)} \mathcal{N}(a | \mu_a, \sigma_a^2) da = \underline{\Phi\left(\frac{\mu_a}{(\lambda^2 + \sigma_a^2)^{1/2}}\right)}$$

$\simeq \sigma(a)$

$$\lambda \cdot \frac{\mu_a}{(\lambda^2/\lambda^2 + \lambda^2 \sigma_a^2)^{1/2}}$$

↙

Again approximating

$$\sigma(a) \simeq \Phi(\lambda a)$$

$\simeq \sigma(a^*)$

$$\Rightarrow \int \sigma(a) \mathcal{N}(a | \mu_a, \sigma_a^2) da = \sigma(a^*)$$

$$a^* = \frac{\mu_a}{(1 + \lambda^2 \sigma_a^2)^{1/2}} = \mu_a \left(1 + \frac{1}{8} \pi \sigma_a^2\right)^{-1/2}, \quad \lambda^2 = \frac{1}{8} \pi$$

K(σ_a^2)

\Rightarrow Final approx. predictive distribution

*

$$p(Y=c_1 | X^{new}, D) = \sigma(K(\sigma_a^2) \cdot \mu_a)$$

Variational Logistic Regression w. set priors

- Improvement on regular Bayesian Log. Reg:

 - Better Gaussian approx. of posterior dist. $p(\theta)$ than Laplace approx.

 - Optimize lower bound of approx. and true dist.

 - (NOT restricting mode = mean, etc. !)

\Rightarrow Greater flexibility, greater accuracy

(param. dist. can become more optimal)

\Rightarrow Approx. posterior normal param. dist. $q(\theta)$ incorporating prior $p(\theta)$ and evidence $p(D|\theta)$

Goal: Approx. logistic regression likelihood func. by

"Exponential or quadratic form" \sim Gaussian-like func.

$$\begin{aligned}
 e^{[ax^2 + bx + c]} &= e^{[a(x^2 + \frac{b}{a}x) + c]} \\
 &= e^{[a(x^2 + 2\frac{b}{2a}x + \frac{b^2}{4a^2}) - \frac{b^2}{4a^2} + c]} \Rightarrow \sigma(x) \simeq e^{(ax^2 + bx + c)} \underset{\text{* Known analytical integration}}{\equiv} N(\mu, \sigma^2) \\
 &= e^{[a(x + \frac{b}{2a})^2 - \frac{b^2}{4a^2} + c]} \\
 &= e^{(c - \frac{b^2}{4a^2})} e^{a(x + \frac{b}{2a})^2} \\
 &\equiv \frac{(2\pi\sigma^2)^{-\frac{1}{2}}}{\sigma} e^{-\frac{1}{2\sigma^2}(x - \mu)^2} \quad \left| \begin{array}{l} \mu = e^{(c - \frac{b^2}{4a^2})} \\ -\frac{1}{2\sigma^2} = a \\ \mu = -\frac{b}{2a} \end{array} \right. \Rightarrow \text{Lower bound to optimize}
 \end{aligned}$$

Marginal likelihood of Bayesian logistic regression model:

* Likelihood of data over all models

Likelihood of data conditioned on a model θ

$$p(D) = \int p(D, \theta) d\theta = \int p(D|\theta) p(\theta) d\theta = \int \left[\prod_{i=1}^N p(\hat{y}_i | \theta) \right] p(\theta) d\theta$$

$$\text{Rewriting } p(\hat{y} | \theta) = \sigma \left(\frac{\theta^T x}{a} \right)^y \left(1 - \sigma \left(\frac{\theta^T x}{a} \right) \right)^{(1-y)}$$

Known "family of models"
 \sim Parametric function

Substitution

$$a = \theta^T x$$

$$= \sigma(a)^y (1 - \sigma(a))^{(1-y)}$$

$$= \left(\frac{1}{1 + e^{-a}} \right)' \left(1 - \frac{1}{1 + e^{-a}} \right)$$

⋮ ?

$$p(\hat{y} | \theta) = e^{ay} \sigma(-a)$$

- Obtaining a lower bound on $p(D)$ by using known lower bound on $\sigma(\cdot)$

*

$$\sigma(z) \geq \sigma(\xi) \exp \left[\frac{z-\xi}{2} - \lambda(\xi)(z^2 - \xi^2) \right]$$

Optimal variational parameter λ

which gives optimal lower bound $\lambda(\xi) = \frac{1}{2\xi} \left[\sigma(\xi) - \frac{1}{2} \right]$
approximation $q(x, \xi)$ for $\sigma(x)$

Expressing likelihood of data $p(\hat{y} | \theta)$ in terms of the lower bound $q(\cdot)$ for $\sigma(\cdot)$

$$p(\hat{y}_i | \theta) = e^{a_i y_i} \sigma(-a_i) \geq e^{a_i y_i} \sigma(\xi_i) \exp \left(-\frac{a_i - \xi_i}{2} - \lambda(\xi_i)(a_i^2 - \xi_i^2) \right)$$

Sample-specific parameters: $x_i \rightarrow a_i, y_i, \xi_i$

Expressing the joint distribution (of all data) using the lower bound $q(\cdot)$ for $\sigma(\cdot)$

$$p(D, \theta) = p(D | \theta) p(\theta) = \left[\prod_{i=1}^N p(\hat{y}_i | \theta) \right] p(\theta)$$

$$\geq \left[\prod_{i=1}^N \exp(\theta^T x_i \cdot y) \sigma(\xi_i) \exp\left(-\frac{1}{2}(\theta^T x_i - \xi_i)^2 - \lambda(\xi_i)([\theta^T x_i]^2 - \xi_i^2)\right) \right] p(\theta)$$

"Label y " term " $x = \xi$ " term "Variational param. λ " term
 $= \left[\prod_{i=1}^N \sigma(\xi_i) \exp\left(\theta^T x_i \cdot y_i - \frac{1}{2}(\theta^T x_i - \xi_i)^2 - \lambda(\xi_i)([\theta^T x_i]^2 - \xi_i^2)\right) \right] p(\theta)$

$h(\theta, \xi)$

Parameters that vary in the functional?

Log of the joint distribution $p(D, \theta)$

Why? - Because it is easier to manipulate the expressions

$$\log p(D, \theta) \geq$$

$$\log(a \cdot b) = \log a + \log b$$

$$\log \left[\prod_{i=1}^N \sigma(\xi_i) \exp\left(\theta^T x_i \cdot y_i - \frac{1}{2}(\theta^T x_i - \xi_i)^2 - \lambda(\xi_i)([\theta^T x_i]^2 - \xi_i^2)\right) \right] p(\theta)$$

$$= \log p(\theta) + \sum_{i=1}^N \left[\log \sigma(\xi_i) + \theta^T x_i \cdot y_i - \frac{1}{2}(\theta^T x_i - \xi_i)^2 - \lambda(\xi_i)([\theta^T x_i]^2 - \xi_i^2) \right]$$

Assuming a Gaussian prior for parameters (like in Bayesian log. reg.)

$$p(\theta) = \mathcal{N}(\Sigma_0) e^{-\frac{1}{2}(\theta - \mu_0)^T \Sigma_0^{-1} (\theta - \mu_0)}$$

$$\log p(\theta) = \log \mathcal{N}(\Sigma_0) - \frac{1}{2}(\theta - \mu_0)^T \Sigma_0^{-1} (\theta - \mu_0)$$

Inscribing into the joint prob. expression

$$\begin{aligned} \log p(D, \theta) &\geq \log \mathcal{N}(\Sigma_0) - \frac{1}{2}(\theta - \mu_0)^T \Sigma_0^{-1} (\theta - \mu_0) \\ &+ \sum_{i=1}^N \left[\log \sigma(\xi_i) + \theta^T x_i \cdot y_i - \frac{1}{2}(\theta^T x_i - \xi_i)^2 - \lambda(\xi_i)([\theta^T x_i]^2 - \xi_i^2) \right] \end{aligned}$$

Optimizing $\theta \Rightarrow$ Discard terms NOT affecting θ

$$\log p(D, \theta)$$

$$\geq -\frac{1}{2}(\theta - \mu_0)^T \Sigma_0^{-1} (\theta - \mu_0) + \sum_{i=1}^N \left[\theta^T x_i y_i - \frac{1}{2}(\theta^T x_i - \xi_i) - \lambda(\xi_i) \theta^T (x_i x_i^T) \theta \right] + \text{const}$$

POINT : The above joint distribution $p(D, \theta)$ encompass the prior $p(\theta)$ and evidence $p(D|\theta)$

\Rightarrow Approx. posterior $q(\theta)$ incorporating both

$$\exp(\log p(D, \theta)) \geq q(\theta) = \exp(\dots)$$

*

$$\Rightarrow q(\theta) = N(\theta | \mu_N, \Sigma_N)$$

(

$$\mu_N = \Sigma_N \left(\Sigma_0^{-1} \mu_0 + \sum_{i=1}^N \left(y_i - \frac{1}{2} \right) x_i \right)$$

$$\Sigma_N = \Sigma_0^{-1} + 2 \sum_{i=1}^N \lambda(\xi_i) x_i x_i^T$$

$$e^{(c - \frac{b^2}{4\sigma^2})} e^{a(x + \frac{b}{2\sigma})^2} \\ (2\pi\sigma^2)^{-\frac{1}{2}} e^{-\frac{1}{2\sigma^2}(x - \mu)^2}$$

* Need to read original paper to see how to connect with normal form ... ?

POINT : Because $q(\theta)$ has a Gaussian form \Rightarrow Already normalized

Evidence prior

$$\text{posterior} \quad \text{Normalization}$$

$$p(\theta | D) = \frac{p(D|\theta)p(\theta)}{p(D)}$$

$$< \Rightarrow p(D) p(\theta | D) = p(D|\theta) p(\theta)$$

"negligible" $\Rightarrow q(\theta) = N(\cdot)$

$$\Rightarrow q(\theta) \leq p(D|\theta)p(\theta)$$

Determine variational parameters ξ_i

How? Maximize the marginal log likelihood lower bound

$$\log p(D) = \log \int p(D|\theta)p(\theta) d\theta \geq \log \int h(\theta, \xi) p(\theta) d\theta = \underline{\mathcal{L}(\xi)}$$

↑
Use the EM algorithm to optimize
- Latent variable is θ

① Initialize $\xi^{\text{old}} = \{\xi_1, \dots, \xi_n\}$ to some value

② E-step : Use ξ^{old} to find "best-estimate" posterior distribution $q(\theta)$ (i.e. the latent params.)

⇒ Compute μ_N, Σ_N using $(D, \mu_0, \Sigma_0, \xi^{\text{old}})$

⇒ Current "best-estimate" $q(\theta) = N(\theta | \mu_N, \Sigma_N)$

③ M-step : Maximize expected complete-data log likelihood
w.r.t. current "best estimate" $q(\theta)$

"Standard variational method approach" REPEAT!

$$Q(\xi, \xi^{\text{old}}) = \mathbb{E}_{\xi} [\mathcal{L}(\xi)] = \mathbb{E}_{\xi} [\log h(\theta, \xi) p(\theta)]$$

$$\begin{aligned} &= \mathbb{E}_{\xi} [\log h(\theta, \xi) + \underline{\log p(\theta)}] \\ &\quad \rightarrow \text{const. as NOT depending on } \xi \\ &= \mathbb{E}_{\xi} [\log h(\theta, \xi)] + \text{const.} \end{aligned}$$

$$= \mathbb{E}_{\xi} \sum_{i=1}^n \left[\log \sigma(\xi_i) + \underline{\theta^T x_i y_i - \frac{1}{2} (\theta^T x_i - \xi_i)^2 - \lambda(\xi_i) ([\theta^T x_i]^2 - \xi_i^2)} \right] + \text{const.}$$

→ const

→ const

$$= E_g \sum_{i=1}^n \left[\log \sigma(\xi_i) - \frac{1}{2} \xi_i^2 - \lambda(\xi_i) (x_i^\top E[\theta \theta^\top] x_i - \xi_i^2) \right] + \text{const.}$$

Taking the derivative to find maximum of expression

Applying dd. of $\sigma(s)$, $\lambda(s)$

$$\frac{d}{d\xi_i} [Q(\xi, g^{dd})] = \dots = \lambda'(\xi_i) (x_i^\top E[\theta \theta^\top] x_i - \xi_i^2) = 0$$

- As $\lambda'(s) \neq 0$ as it is
 - A monotonic function
 - Symmetric around $s=0$

$$\Rightarrow x_i^\top E[\theta \theta^\top] x_i - \xi_i^2 = 0$$

Expectation identity:

$$E[x^2] = E[x]^2 + \text{var}[x]$$

Re-estimation equation:

*

$$(\xi_i^{\text{new}})^2 = x_i^\top E[\theta \theta^\top] x_i = x_i^\top (\Sigma_n + \mu_n \mu_n^\top) x_i$$

Evaluating the lower bound $\mathcal{L}(\xi)$:

- Possible to do analytically because

How?

① $p(\theta)$ is Gaussian

Details in original paper?

② $h(\theta, \xi)$ is "exponential of a quadratic function of θ "

* $\Rightarrow L$ can be formulated as a Gaussian

$$\mathcal{L}(\xi) = \int h(\theta, \xi) p(\theta) d\theta$$

$$= \frac{1}{2} \log \frac{|\Sigma_N|}{|\Sigma_0|} - \frac{1}{2} \mu_N^\top \Sigma_N^{-1} \mu_N + \frac{1}{2} \mu_0^\top \Sigma_0^{-1} \mu_0 + \sum_{i=1}^n [\log \sigma(\xi_i) - \frac{1}{2} \xi_i^2 - \lambda(\xi_i) \xi_i^2]$$

Excuse: Inferring hyperparameters from data

- Approach: Combining 'Global' and 'Local' variational approximation

Extendable to prior with nonzero mean?

Parameter prior assumption : $p(\theta) = N(\theta | 0, \tilde{\alpha}^{-1} I)$
 $\mu_0 = 0, \Sigma_0 = \tilde{\alpha}^{-1} I$

Hyperprior assumption : $p(\alpha) = \text{Gam}(\alpha | a_0, b_0)$

Joint distribution : $p(y, \theta, \alpha) = p(y|\theta)p(\theta|\alpha)p(\alpha)$

Marginal likelihood : $p(D) = \iint p(y, \theta, \alpha) d\theta d\alpha$

Approximate using both 'Global' and
 'Local' variational approaches in same model

Introducing variational distribution $q(\theta, \alpha)$

- General formulation :

Lower bound
 \downarrow Ditt. between lower bound
 and true dist. \sim Remainder

$$\log p(D) = \mathcal{L}(q) + KL(q \parallel p)$$

Repcat...

$$\mathcal{L}(q) = \iint q(\theta, \alpha) \log \left[\frac{p(\theta, \alpha, y)}{q(\theta, \alpha)} \right] d\theta d\alpha$$

$$KL(q \parallel p) = - \iint q(\theta, \alpha) \log \left[\frac{p(\theta, \alpha | y)}{q(\theta, \alpha)} \right] d\theta d\alpha$$

POINT : Intractability of $\mathcal{L}(q)$ stems from $p(y|\theta)$
 w. the logistic sigmoid ($p(\theta|\alpha)p(\alpha)$ analytically ok)

\Rightarrow Apply the Local variational bound

$$\begin{aligned} p(y, \theta, \alpha) &= p(y|\theta)p(\theta|\alpha)p(\alpha) \\ &\geq h(\theta, \xi)p(\theta|\alpha)p(\alpha) \end{aligned}$$

$$\Rightarrow \log p(D) \geq \mathcal{L}(q) \geq \tilde{\mathcal{L}}(q, \xi) \quad \text{w. local approx. at sigmoid}$$

$$= \iint q(\theta, \alpha) \log \left[\frac{h(\theta, \xi)p(\theta|\alpha)p(\alpha)}{q(\theta, \alpha)} \right] d\theta d\alpha$$

Mean Field Approximation

$$q(\theta, \alpha) = q(\theta) q(\alpha)$$

$$\log q_j^*(z_j) = \mathbb{E}_{\gamma_j} [\log p(x, z)] + \text{const}$$

Applying the "general result"

\Rightarrow Optimal factors for $q(\theta)$, $q(\alpha)$

① $q(\theta)$:

$$\log q(\theta) = \mathbb{E}_\alpha [\log p(y, \theta, \alpha)] + \text{const}$$

$$\simeq \mathbb{E}_\alpha [\log (h(\theta, \xi)p(\theta|\alpha)p(\alpha))] + \text{const}$$

$$= \mathbb{E}_\alpha [\log h(\theta, \xi) + \log p(\theta|\alpha) + \log p(\alpha)] + \text{const}$$

Independent of $\theta \rightarrow \text{const}$

$$= \log h(\theta, \xi) + E_{\alpha} [\log p(\theta | \alpha)] + \text{const}$$

Known

Known

$$= \sum_{i=1}^n \left[\underbrace{\log \sigma(\xi_i)}_{\text{Independent of } \theta \rightarrow \text{const}} + \theta^T x_i y_i - \frac{1}{2} (\theta^T x_i - \xi_i)^2 - \lambda(\xi_i) ([\theta^T x_i]^2 - \xi_i^2) \right]$$

$$+ E_{\alpha} \left[\log \Gamma(\alpha) + \log \exp \left(-\frac{1}{2} (\theta - \bar{\theta})^T \Sigma^{-1} I (\theta - \bar{\theta}) \right) \right] + \text{const.}$$

$$= \sum_{i=1}^n \left[\theta^T x_i y_i - \frac{1}{2} \theta^T x_i - \lambda(\xi_i) \theta^T x_i x_i^T \theta \right] + E_{\alpha} \left[-\frac{1}{2} \alpha \theta^T I \theta \right] + \text{const}$$

$$= \sum_{i=1}^n \left[(y_i - \frac{1}{2}) \theta^T x_i - \lambda(\xi_i) \theta^T x_i x_i^T \theta \right] - \frac{1}{2} E[\alpha] \theta^T \theta + \text{const}$$

POINT : Quadratic function of $\theta \Rightarrow$ Express $g(\theta)$ in Gaussian form

*

$$g(\theta) = N(\theta | \mu_N, \Sigma_N)$$

?

$$\mu_N = \sum_N \left(\sum_{i=1}^n (y_i - \frac{1}{2}) x_i \right)$$

$$\Sigma_N^{-1} = E[\alpha] I + 2 \sum_{i=1}^n \lambda(\xi_i) x_i x_i^T$$

Expectation over Gamma dist. param.

$$E[\alpha] = \frac{a_N}{b_N}$$

② $g(\alpha) :$

$$\log g(\alpha) = E_{\theta} [\log p(y, \theta, \alpha)] + \text{const}$$

$$= E_{\theta} \left[\log h(\theta, \alpha) + \log p(\theta | \alpha) + \log p(\alpha) \right] + \text{const}$$

Independent of $\alpha \rightarrow \text{const}$

$$= E_{\theta} [\log p(\theta | \alpha)] + \log p(\alpha) + \text{const}$$

$$= E_{\theta} \left[\log \frac{1}{(2\pi)^D} \frac{1}{|\alpha^T I|^{1/2}} \exp(-\frac{1}{2}(\theta - \bar{\theta})^T \alpha I (\theta - \bar{\theta})) \right]$$

$$+ \log \frac{1}{T(a_0)} b_0^{a_0} \alpha^{a_0-1} e^{-b_0 \alpha} + \text{const}$$

$$|\alpha^T I| = \alpha^T \cdot \alpha^T \dots = \alpha^D$$

$$= E_{\theta} \left[\log (2\pi)^D + \log |\alpha^T I|^{-1/2} - \frac{1}{2} \alpha^T \theta^T \theta \right] + \log T(a_0) b_0^{a_0} + \log \alpha^{a_0-1} + \log e^{-b_0 \alpha} + \text{const}$$

$\rightarrow \text{const}$ Dim. of covariance matrix $\alpha^T I$ $\rightarrow \text{const}$

$$= E_{\theta} \left[\frac{D}{2} \log \alpha - \frac{1}{2} \alpha^T \theta^T \theta \right] + (a_0-1) \log \alpha - b_0 \alpha + \text{const}$$

$$= \frac{D}{2} \log \alpha - \frac{1}{2} \alpha^T E_{\theta} [\theta^T \theta] + (a_0-1) \log \alpha - b_0 \alpha + \text{const}$$

POINT: With ingenuity, see that this expression has the form of a Gamma distribution

*

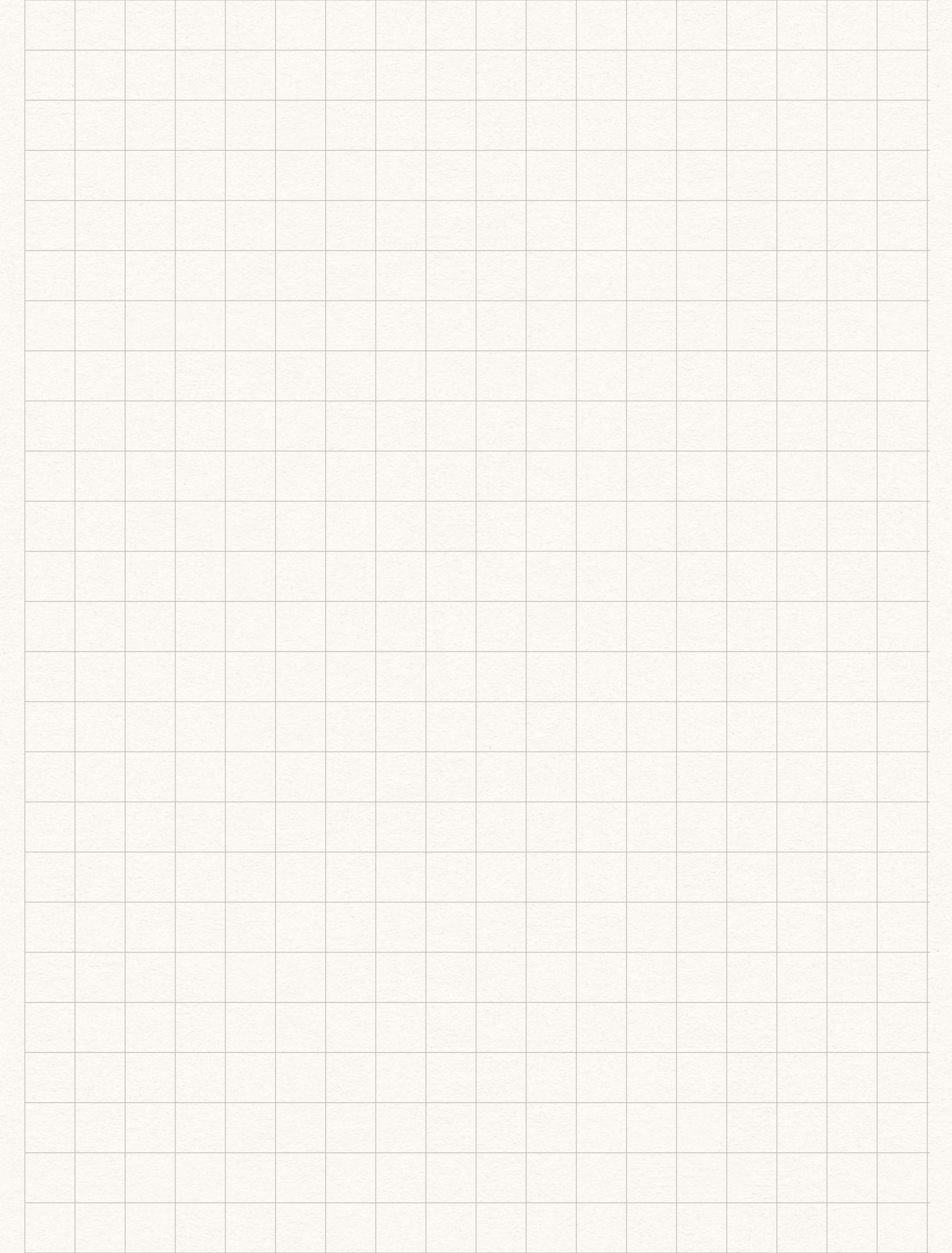
$$f(\alpha) = \text{Gam}(\alpha | a_N, b_N) = \frac{1}{T(a_N)} a_N^{b_N} \alpha^{a_N-1} e^{-b_N \alpha}$$

↑

$$a_N = a_0 + \frac{D}{2}$$

$$b_N = b_0 + \frac{1}{2} E_{\theta} [\theta^T \theta]$$

Optimizing ξ has same form as previously



$$\Rightarrow f(\alpha), f(w)$$

$$(\xi_i^*)^2 = \phi_i^\top (\Sigma_N + \mu_N \mu_N^\top) \phi_i$$

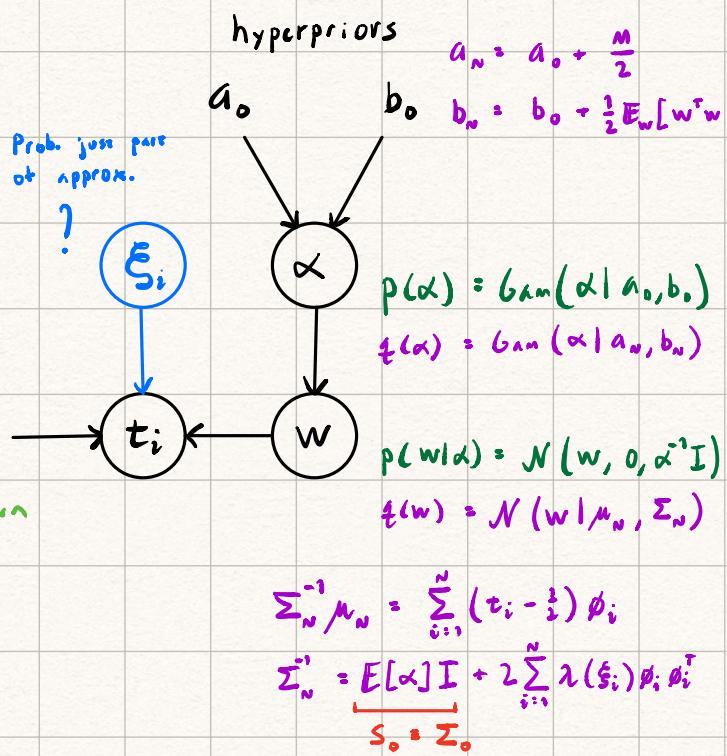
$$E[\alpha] = \frac{a_0}{b_0} \text{ Gamma}$$

$$E[w^\top w] = \Sigma_N + \mu_N^\top \mu_N \text{ Gaussian}$$

$$(E[ww^\top] = \Sigma_N + \mu_N \mu_N^\top) \text{ Matrix}$$

Initialize parameters

- a_0, b_0, ξ_i



$$\lambda(\xi) = \frac{1}{2\epsilon} \left(\sigma(\xi) - \frac{1}{2} \right)$$

$\epsilon \rightarrow 0 \text{ as } \xi \rightarrow 0$
* add numerical ϵ

Expectation

$$E[\alpha] = \frac{a_0}{b_0}$$

Expectation step

$$1. E[\alpha]$$

$$2. E[w^\top w]$$

$$E_w[w^\top w] = (\Sigma_N + \mu_N^\top \mu_N)$$

→ Removed?

Maximization (Re-estimation equations)

$$f(w) = N(w | \mu_N, \Sigma_N)$$

$$\cdot \mu_N = \sum_N \sum_{i=1}^n (t_i - \frac{1}{2}) \phi_i$$

Maximization step

$$(1. a_N)$$

$$2. b_N$$

$$3. \xi$$

$$\cdot \Sigma_n^{-1} = E[\alpha] I + 2 \sum_{i=1}^n \lambda(\xi_i) \phi_i \phi_i^\top \quad 4. \quad \Sigma_n^{-1}$$

$$5. \quad \mu_n$$

$$g(\alpha) = \text{Gam}(\alpha | a_n, b_n)$$

$$\cdot a_n = a_0 + \frac{m}{2} \leftarrow \text{Feature dim.}$$

$$\cdot b_n = b_0 + \frac{1}{2} E_w [w^\top w]$$

Remember!

$$(\xi_i^{new})^2 = \phi_i^\top (\Sigma_n + \mu_n \mu_n^\top) \phi_i$$

Predictive dist. $p(t | \phi, D)$

- Assumption: "Same term"

\Rightarrow Same approximate predictive distribution?

$$p(t | \phi, D) = \sigma(K(\sigma_n^2) \cdot \mu_n)$$

$$K(\sigma_n^2) = \left(1 + \frac{\pi}{8} \sigma_n^2\right)^{-1/2}$$

$$\mu_n = \mu_n^\top \phi$$

$$\Sigma_n = S_n \text{ in } \mathcal{N}(w_{new}, \Sigma_n)$$

$$\sigma_n^2 = \phi^\top \Sigma_n \phi$$