## SL = BAYESIAN LINEAR MODELS

We revisit linear regression & logistic regression from a Bayesian point of view. For background information on Bayesian statistics, see MS: BAYESIAN STATISTICS .

## I. BAYESIAN LINEAR REGRESSION

Consider a learning sample $\mathcal{L}_n = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$, where $(X_i, Y_i)$ are iid, and are assumed to arise from a linear model $Y = X\beta + \mathcal{E}$, where

$$Y = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}_{(n \times 1)}, \quad X = \begin{pmatrix} 1 & X_{11} & \cdots & X_{1d} \\ \vdots & \vdots & & \vdots \\ 1 & X_{n1} & \cdots & X_{nd} \end{pmatrix}_{n \times (d+1)}, \quad \mathcal{E} = \begin{pmatrix} \mathcal{E}_1 \\ \vdots \\ \mathcal{E}_n \end{pmatrix}_{(n \times 1)}, \quad \beta = \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_d \end{pmatrix}$$

see pages 1/2 in SL: LINEAR REGRESSION .

→ The frequentist approach to linear regression assumes that the vector of parameters $\beta$ is <u>fixed</u> and <u>unknown</u>. It is estimated by maximizing the likelihood function ($\equiv$ least squares estimate under normal errors $\mathcal{E}$).

→ The Bayesian treatment considers the vector $\beta$ to be random, with prior distribution $\underline{f(\beta)}$:

$$\boxed{\begin{array}{c} Y = X\beta + \mathcal{E}, \qquad \mathcal{E} \sim \mathcal{N}(\mathcal{E} \mid 0, \gamma^{-1} I_n) \\[2mm] + \text{prior} \quad \beta \sim \mathcal{N}(\beta \mid 0, \alpha^{-1} I_d) = f(\beta) \end{array}}$$

BAYESIAN LINEAR MODEL

For notational convenience, we assume that $\beta \in \mathbb{R}^d$

Remarks : (i) More generally, we can assume that
$$f(\beta) = \mathcal{N}(\beta \mid m_0, S_0).$$ Subsequent calculations can be easily adapted.

(ii) We consider two cases :

↘ $\beta$ unknown, $\gamma$ known (sections I.1 & I.2)

↘ $\beta$ unknown, $\gamma$ unknown (section I.3)

In addition, we consider the case where hyperpriors are introduced for $\alpha$ and $\gamma$; the so-called <u>Empirical Bayes</u> / <u>Type II ML</u> / <u>Evidence approximation</u> approach (section I.4).

In each case, we are interested in

↘ the <u>posterior distribution</u> of $\beta$, having observed $\mathcal{L}_n$ (→ useful for the construction of <u>credible intervals</u>).

↘ the <u>predictive distribution</u> of $Y$ given $\mathcal{L}_n$, and a new input point $x$.

### I.1. Posterior distribution ($\gamma$ known).

• The posterior distribution is $\underbrace{f(\beta \mid \mathcal{L}_n)}_{\text{posterior}} \propto \underbrace{f(\mathcal{L}_n \mid \beta)}_{\text{likelihood}} \underbrace{f(\beta)}_{\text{prior}}$,

"proportional to"

where

• $f(\mathcal{L}_n \mid \beta) = \left(\dfrac{\gamma}{2\pi}\right)^{n/2} \exp\left\{-\dfrac{\gamma}{2}(y - X\beta)^t (y - X\beta)\right\}$,

$f(\beta) = \left(\dfrac{\alpha}{2\pi}\right)^{1/2} \exp\left\{-\dfrac{\alpha}{2}\beta^t\beta\right\}$.

The product of the likelihood by the prior is proportional ③
to:

$$\sim \exp\left\{-\frac{1}{2}\left[\underbrace{\gamma \beta^t X^t X \beta - 2\gamma y^t X\beta + \gamma y^t y}_{\text{from the likelihood}} + \underbrace{\alpha \beta^t \beta}_{\substack{\text{from the} \\ \text{prior}}}\right]\right\}$$

$$= \exp\left\{-\frac{1}{2}\left[\beta^t(\alpha I_d + X^t X)\beta - (\gamma X^t y)^t \beta + \underset{\substack{\text{indpt} \\ \text{of } \beta}}{\text{constant}}\right]\right\}$$

We recognize here the expression of the multivariate normal density.

$$\Rightarrow f(\beta \mid \mathcal{L}_n) = \mathcal{N}(\beta \mid m_n, S_n).$$

To find the expression of $m_n$ and $S_n$, compare the terms in the expression above with

$$(\beta - m_n)^t S_n^{-1}(\beta - m_n) = \beta^t S_n \beta - 2 m_n^t S_n^{-1}\beta + m_n^t S_n^{-1} m_n.$$

We immediately get : $\begin{cases} S_n^{-1} = \alpha I_d + \gamma X^t X \\ m_n = \gamma S_n X^t y \end{cases}$

<div style="border:1px solid red; padding:4px;">

**Summary**: $X = X\beta + \varepsilon$, $\varepsilon \sim \mathcal{N}(\varepsilon \mid 0, \gamma^{-1} I_n)$
$\beta \sim \mathcal{N}(\beta \mid 0, \alpha^{-1} I_d)$

Posterior is $f(\beta \mid \mathcal{L}_n) = \mathcal{N}(\beta \mid m_n, S_n)$, with

$\begin{pmatrix} m_n = \gamma S_n X^t y \\ S_n^{-1} = \alpha I_d + \gamma X^t X \end{pmatrix}$

</div>

(*)

×**Remarks**:(i) Assuming more generally that $f(\beta) = \mathcal{N}(\beta \mid m_0, S_0)$, we easily derive the expression for the posterior mean & covariance:

$$m_n = S_n(S_0^{-1} m_0 + \gamma X^t y) \quad \text{and} \quad S_n^{-1} = S_0^{-1} + \gamma X^t X.$$

(ii) Bayesian linear regression & Ridge Regression (RR) ④

The log of the posterior distribution is :

$$\log f(\beta \mid \mathcal{L}_n) = -\frac{\gamma}{2}\sum_{i=1}^{n}(y_i - x_i^t\beta)^2 - \frac{\alpha}{2}\beta^t\beta$$

$$= -\frac{\gamma}{2}\left\{\underset{=\ RSS_2(\alpha/\gamma)}{\underbrace{\sum_{i=1}^{n}(y_i - x_i^t\beta)^2 + \frac{\alpha}{\gamma}\beta^t\beta}}\right\}$$

see p.5 in SL: RR AND LASSO.

• **Consequences**:

→ $\alpha/\gamma$ is a tuning parameter, and quantifies the trade-off between the goodness-of-fit term ( ≡ likelihood) and the penalty ( ≡ prior).
It can be estimated using cross-validation techniques.

→ Ridge solution = MAP estimator.

(iii) Points arriving sequentially.

Assuming a stream of observations

$$(x_1, y_1) \to (x_2, y_2) \to \cdots \to (x_n, y_n) \to (x_{n+1}, y_{n+1}) \to \cdots;$$

the posterior distribution after $n$ points are collected is $\mathcal{N}(\beta \mid m_n, S_n)$.

A new observation $(x_{n+1}, y_{n+1})$ has density / likelihood

$$f(y_{n+1} \mid x_{n+1}, \beta) = \left(\frac{\gamma}{2\pi}\right)^{1/2}\exp\left[-\frac{\gamma}{2}(y_{n+1} - \beta^t x_{n+1})^2\right],$$

since $Y_{n+1} = \beta^t x_{n+1} + \varepsilon$, where $\varepsilon \sim \mathcal{N}(0, \gamma^{-1})$

The posterior distribution, having observed $(n+1)$ points, is proportional (w.r.t. $\beta$) to:

$$\exp\left\{-\frac{1}{2}\left[(\beta-m_n)^t S_n^{-1}(\beta-m_n) + \gamma(y_{n+1}-\beta^t x_{n+1})^2\right]\right\}$$

$$\|$$

$$\left[\beta^t\left(S_n^{-1} + \gamma\, x_{n+1} x_{n+1}^t\right)\beta\right.$$

$$-2\beta^t\left(S_n^{-1} m_n + \gamma\, x_{n+1}\, y_{n+1}\right)$$

$$\left.+ \text{constant indpt of } \beta\right]$$

Compare this expression with $(\beta-m_{n+1})^t S_{n+1}^{-1}(\beta-m_{n+1})$, appearing in the posterior $f(\beta\mid\mathcal{L}_{n+1}) = \mathcal{N}(\beta\mid m_{n+1}, S_{n+1})$

We see that
$$\begin{cases} S_{n+1}^{-1} = S_n^{-1} + \gamma\, x_{n+1} x_{n+1}^t \\ m_{n+1} = S_{n+1}\left(S_n^{-1} m_n + \gamma\, x_{n+1}\, y_{n+1}\right), \end{cases}$$

with $m_0 = 0$ and $S_0 = \alpha^{-1} I$.

coincides with the formula (*) on page 3 since :

$\bullet$ $S_{n+1}^{-1} = S_{n-1}^{-1} + \gamma\, x_{n+1} x_{n+1}^t + \gamma\, x_n x_n^t$

$$\vdots$$

$$= S_0^{-1} + \gamma \sum_{i=1}^{n+1} x_i x_i^t$$

$$= \alpha I + \gamma\, X^t X$$

$\bullet$ $m_{n+1} = S_{n+1}\left(S_n^{-1}\left[\underbrace{S_n\left(S_{n-1}^{-1} m_{n-1} + \gamma\, x_n y_n\right)}_{m_n}\right] + \gamma\, x_{n+1} y_{n+1}\right)$

$$= S_{n+1}\left(S_{n-1}^{-1} m_{n-1} + \gamma\, x_n y_n + \gamma\, x_{n+1} y_{n+1}\right)$$

$$= S_{n+1}\left(\overset{0}{\cancel{S_0^{-1} m_0}} + \gamma\sum_{i=1}^{n+1} x_i y_i\right)$$

$$= \gamma\, S_{n+1}\, X^t y$$

---

Csq: in a sequential setting,

| posterior distribution $f(\beta\mid\mathcal{L}_n)$ | $=$ | prior distribution on $\beta$ for a new observation $(x_{n+1}, y_{n+1})$. |
|---|---|---|

$$(x_1, y_1), (x_2, y_2), \cdots (x_n, y_n), (x_{n+1}, y_{n+1}), \cdots$$

posterior $\mathcal{N}(\beta\mid m_n, S_n)$ → use as prior for

posterior $\mathcal{N}(\beta\mid m_{n+1}, S_{n+1})$

### I.2. Predictive distribution ($\gamma$ known).

$\bullet$ Given a new input point $x$, the predictive distribution is

$$f(y\mid\mathcal{L}_n, x) = \int f(y, \beta\mid\mathcal{L}_n, x, \alpha, \gamma)\, d\beta$$

$$= \int \underbrace{f(y\mid\beta, x, \gamma)}_{\substack{\text{target var. distrib}\\ \mathcal{N}(y\mid x^t\beta, \gamma^{-1})}}\ \underbrace{f(\beta\mid\mathcal{L}_n, \alpha, \gamma)}_{\substack{\text{posterior distrib}\\ \mathcal{N}(\beta\mid m_n, S_n)}}\, \underbrace{d\beta}_{\substack{\beta\text{ is}\\ \text{integrated}\\ \text{out}}}$$

$$= \text{convolution of two Gaussian distributions}$$
$$\Rightarrow \text{still gaussian.}$$

$$\boxed{f(y\mid\mathcal{L}_n, x) = \mathcal{N}\left(y\mid m_n^t x, \underbrace{\gamma^{-1} + x^t S_n\, x}_{=:\,\sigma_n^2(x)}\right)}$$

More generally, if $p(x) = \mathcal{N}(x\mid\mu, \Lambda^{-1})$
$$p(y\mid x) = \mathcal{N}(y\mid A x + b, L^{-1}),$$

then $p(y) = \mathcal{N}(y\mid A\mu + b, L^{-1} + A\Lambda^{-1}A^t)$, see Bishop p. 93

Remarks: (i) As $n \to +\infty$, $\sigma_n^2(x) \to \gamma^{-1}$ ($\equiv$ noise variance)

Indeed,

$$\sigma_{n+1}^2(x) = \gamma^{-1} + x^t S_{n+1} x, \quad \text{where}$$

$$S_{n+1} = \left( S_n^{-1} + \gamma x_{n+1} x_{n+1}^t \right)^{-1} \quad \text{(page 5)}$$

$$= S_n - \frac{(S_n x_{n+1} \gamma^{1/2})(\gamma^{1/2} x_{n+1}^t S_n)}{1 + \gamma x_{n+1}^t S_n x_{n+1}}$$

$$= S_n - \gamma \frac{S_n x_{n+1} x_{n+1}^t S_n}{1 + \gamma x_{n+1}^t S_n x_{n+1}}$$

Toolbox:

$$(A + vv^t)^{-1}$$

$$\overset{\shortparallel}{A^{-1} - \frac{(A^{-1}v)(v^t A^{-1})}{1 + v^t A^{-1} v}}$$

$$\sigma_{n+1}^2(x) = \gamma^{-1} + x^t \left( \underline{\quad \shortparallel \quad} \right) x \quad \to \text{PSD}$$

$$= \sigma_n^2(x) - \gamma \frac{x^t \boxed{S_n x_{n+1} x_{n+1}^t S_n} x}{1 + \gamma x_{n+1}^t S_n x_{n+1}}$$

non-negative since $S_n$ is positive semi-def. (PSD) $\geqslant 0$

Thus, $\sigma_{n+1}^2(x) \leqslant \sigma_n^2(x)$, as required. ∎

(ii) In a Bayesian linear regression setting, the posterior can be computed analytically. Alternatively, we could sample points $\beta_i \sim \mathcal{N}(m_n, S_n)$ from the posterior distribution, and consider the Monte-Carlo approximation to the predictive distribution:

$$\frac{1}{M} \sum_{i=1}^{M} f(y \mid \beta_i, x, \gamma).$$

will be useful in more complex settings.

---

## I.3. Case when $\beta$ and $\gamma$ are unknown.

We assume now that $Y = X\beta + \varepsilon$, $\varepsilon \sim \mathcal{N}(0, \gamma^{-1} I_n)$, with $\gamma$ unknown.

↳ the conjugate prior on $(\beta, \gamma)$ is the normal-gamma distribution (see MS: BAYESIAN STATISTICS):

$$f(\beta, \gamma) = \mathcal{N}(\beta \mid m_0, \gamma^{-1} S_0) \, \text{Gamma}(\gamma \mid a_0, b_0),$$

where

$$\text{Gamma}(\gamma \mid a_0, b_0) = \frac{b_0^{a_0}}{\Gamma(a_0)} \gamma^{a_0 - 1} e^{-b_0 \gamma}, \quad \gamma > 0$$

↳ the posterior distribution is $f(\beta, \gamma \mid \mathcal{L}_n) \propto f(\mathcal{L}_n \mid \beta, \gamma) f(\beta, \gamma)$

In the log-space,

$$= \log f(\beta, \gamma) + \log f(\mathcal{L}_n \mid \beta, \gamma)$$

$$= \boxed{\log \mathcal{N}(\beta \mid m_0, \gamma^{-1} S_0)} + \boxed{\log \text{Gamma}(\gamma \mid a_0, b_0)}$$

$$+ \boxed{\log \mathcal{N}(y \mid X\beta, \gamma^{-1} I_n)}$$

$$= \boxed{\frac{d}{2} \log \gamma - \frac{1}{2} \log |S_0| - \frac{\gamma}{2} (\beta - m_0)^t S_0^{-1} (\beta - m_0)}$$

$$\boxed{- b_0 \gamma + (a_0 - 1) \log \gamma}$$

$$\boxed{+ \frac{n}{2} \log \gamma - \frac{\gamma}{2} \sum_{i=1}^{n} (y_i - \beta^t x_i)^2}$$

$$+ \text{constant term in } \beta, \gamma.$$

We selected the prior distribution such that the posterior belongs to the same family of distribution.

$$\Rightarrow f(\beta, \gamma \mid \mathcal{L}_n) = \text{normal-gamma}.$$

To find the parameters of the normal-gamma distribution, we write the posterior as a product of two densities (one will correspond to the normal term, the other to the gamma term):

$$f(\beta, \gamma \mid \mathcal{L}_n) = \underbrace{f(\beta \mid \mathcal{L}_n, \gamma)}\, f(\gamma \mid \mathcal{L}_n).$$

We first identify this term, and collect in the expression at the bottom of page 8 all terms involving $\beta$:

$$-\frac{\gamma}{2} \beta^t (X^t X + S_o^{-1}) \beta + \gamma \beta^t (S_o^{-1} m_o + X^t y) + \text{cst in } \beta$$

Compare the terms with those appearing in

$$f(\beta \mid \mathcal{L}_n, \gamma) = \mathcal{N}(\beta \mid m_n, \gamma^{-1} S_n):$$

$$-\frac{1}{2}(\beta - m_n)^t (\gamma^{-1} S_n)^{-1} (\beta - m_n)$$

$$= -\frac{\gamma}{2} \beta^t S_n^{-1} \beta + \gamma \beta^t S_n^{-1} m_n + \text{cst in } \beta$$

we see that

$$\begin{pmatrix} \gamma S_n^{-1} = \gamma (X^t X + S_o^{-1}) \\ \gamma (S_o^{-1} m_o + X^t y) = \gamma S_n^{-1} m_n. \end{pmatrix}$$

$$\Rightarrow \boxed{\begin{aligned} f(\beta \mid \mathcal{L}_n, \gamma) &= \mathcal{N}(\beta \mid m_n, \gamma^{-1} S_n), \text{ with} \\ m_n &= S_n (S_o^{-1} m_o + X^t y) \\ S_n^{-1} &= S_o^{-1} + X^t X \end{aligned}}$$

It remains to identify all remaining terms to identify the parameters of the Gamma distribution $\text{Gamma}(\gamma \mid a_n, b_n)$.
Note that the term $\left(\frac{d}{2} \log \gamma\right)$ appearing in the expression on the bottom of page 8 is incorporated into the expression of $f(\beta \mid \mathcal{L}_n, \gamma)$.

$$\Rightarrow \log f(\gamma \mid \mathcal{L}_n) = \boxed{-b_o \gamma + (a_o - 1) \log \gamma}$$

$$\boxed{+ \frac{n}{2} \log \gamma - \frac{\gamma}{2} \sum_{i=1}^n y_i^2}$$

$$\boxed{- \frac{\gamma}{2} m_o^t S_o^{-1} m_o}$$

$$\boxed{+ \frac{\gamma}{2} m_n^t S_n^{-1} m_n} \longrightarrow \text{from completing the squares}$$

$$= \log \text{ of a gamma distribution, with parameters}$$

$$a_n = a_o + \frac{n}{2}$$

$$b_n = b_o + \frac{1}{2}\left(\sum y_i^2 + m_o^t S_o^{-1} m_o - m_n^t S_n^{-1} m_n\right)$$

**Summary:**
- prior on $(\beta, \gamma)$ is
$$f(\beta, \gamma) = \mathcal{N}(\beta \mid m_o, \gamma^{-1} S_o) \, \text{Gamma}(\gamma \mid a_o, b_o)$$

- posterior is
$$f(\beta, \gamma \mid \mathcal{L}_n) = \mathcal{N}(\beta \mid m_n, \gamma^{-1} S_n) \, \text{Gamma}(\gamma \mid a_n, b_n),$$
with
$$\rightarrow m_n = S_n (S_o^{-1} m_o + X^t y)$$
$$\rightarrow S_n^{-1} = S_o^{-1} + X^t X$$
$$\rightarrow a_n, b_n \text{ given above}$$

Before moving on to the predictive distribution, we review ⑪
a useful result:

"Bayesian view" of     Assume  $X \sim \mathcal{N}(\mu, s\lambda^{-1})$
__Student distribution__ :
                                $\lambda \sim gamma(a, b)$.

Then  $X$ has a Student's t-distribution.

$\hookrightarrow$ Indeed,

$f(x) = \int_0^{+\infty} f(x, \lambda) \, d\lambda$

$= \int_0^{+\infty} f(x \mid \lambda) f(\lambda) \, d\lambda$

$= \int_0^{+\infty} \left(\frac{\lambda}{2\pi s}\right)^{1/2} e^{-\frac{\lambda}{2s}(x-\mu)^2} \frac{b^a}{\Gamma(a)} \lambda^{a-1} e^{-b\lambda} \, d\lambda$

$= \frac{b^a}{\Gamma(a)} \left(\frac{1}{2\pi s}\right)^{1/2} \int_0^{+\infty} \lambda^{a-\frac{1}{2}} e^{\boxed{-\lambda\left(b + \frac{(x-\mu)^2}{2s}\right)}} \, d\lambda$

<span style="color:green">Change of variable $z = \lambda u$,
$u = b + \frac{1}{2s}(x-\mu)^2$</span>

$\underbrace{\qquad\qquad\qquad\qquad}$
$\overset{\shortparallel}{\phantom{a}}$

$u^{-a+\frac{1}{2}} \int_0^{+\infty} e^{-z} u^{-1} z^{a-\frac{1}{2}} \, dz$

$= \frac{b^a}{\Gamma(a)} \left(\frac{1}{2\pi s}\right)^{1/2} u^{-a-1/2} \underbrace{\int_0^{+\infty} z^{a-\frac{1}{2}} e^{-z} \, dz}_{= \Gamma(a+\frac{1}{2})}$

$= \frac{b^a}{\Gamma(a)} \left(\frac{1}{2\pi s}\right)^{1/2} \Gamma(a+\frac{1}{2}) \left(b + \frac{(x-\mu)^2}{2s}\right)^{-a-\frac{1}{2}}$

---

Put $\begin{cases} k = 2a & \longrightarrow & a = k/2 \\ \tau = \frac{a}{b} s^{-1} & \longrightarrow & b = \frac{a}{\tau s} = \frac{k}{2\tau s} \end{cases}$  ⑫

$f(x) = \left(\frac{k}{2\tau s}\right)^{k/2} \frac{1}{\Gamma(k/2)} \left(\frac{1}{2\pi s}\right)^{1/2} \Gamma\left(\frac{k+1}{2}\right) \underbrace{\left(\frac{k}{2\tau s} + \frac{(x-\mu)^2}{2s}\right)^{-\frac{k+1}{2}}}_{\shortparallel}$

$\left(\frac{k}{2\tau s}\right)^{-\frac{k+1}{2}} \left(1 + \frac{\tau}{k}(x-\mu)^2\right)^{-\frac{k+1}{2}}$

$\boxed{f(x) = \frac{\Gamma\left(\frac{k+1}{2}\right)}{\Gamma(k/2)} \left(\frac{\tau}{\pi k}\right)^{1/2} \left(1 + \frac{\tau(x-\mu)^2}{k}\right)^{-\frac{k+1}{2}}}$  <span style="color:green">$k = 2a$<br>$\tau = \frac{a}{b} s^{-1}$</span>

$\uparrow$ Compare with the expression page 39
in <span style="color:green">__PT: POPULAR DISTRIBUTIONS__</span>   $\boxed{X \sim t(k, \mu, \tau)}$

Remark: $\left(1 + \frac{\tau(x-\mu)^2}{k}\right)^{-\frac{k+1}{2}} = \exp\left\{-\frac{k+1}{2} \log\left(1 + \frac{\tau(x-\mu)^2}{k}\right)\right\}$

$\approx \exp\left\{-\frac{k+1}{2}\left(\frac{\tau(x-\mu)^2}{k} + O(k^{-2})\right)\right\}$

$\longrightarrow \exp\left(-\frac{1}{2}\tau(x-\mu)^2\right)$ as $k \to +\infty$

so that
$X \overset{d}{\to} \mathcal{N}(\mu, \tau^{-1})$.

• Back to the __predictive distribution__.

$f(y \mid \mathcal{L}_n, x) = \iint f(y \mid \beta, \gamma, x) f(\beta, \gamma \mid \mathcal{L}_n) \, d\beta \, d\gamma$

$= \iint \mathcal{N}(y \mid x^t\beta, \gamma^{-1}) \, \mathcal{N}(\beta \mid m_n, \gamma^{-1}S_n)$
<span style="color:green">$gamma(\gamma \mid a_n, b_n) \, d\beta \, d\gamma$</span>

× First, compute the integral w.r.t. $\beta$ :

$$\int \mathcal{N}(y \mid x^t \beta, \gamma^{-1}) \, \mathcal{N}(\beta \mid m_n, \gamma^{-1} S_n) \, d\beta$$

↖ Using the general formula at bottom of page 6, we see that this integral is

$$\mathcal{N}\left(y \mid x^t m_n, \; \gamma^{-1}\left(1 + \underbrace{x^t}_{} [S_0^{-1} + \underbrace{\underline{\underline{X}}^t \underline{\underline{X}}}_{}]^{-1} x\right)\right)$$

new point (↑)   matrix of observations (↑)

× It remains to compute the integral

$$\int \mathcal{N}\left(y \mid x^t m_n, \; \gamma^{-1}\left(1 + x^t [S_0^{-1} + \underline{\underline{X}}^t \underline{\underline{X}}]^{-1} x\right)\right) \text{Gamma}(\gamma \mid a_n, b_n) \, d\gamma$$

which is a <u>Student's t distribution</u> $t(y \mid k, \mu, \tau)$, with

- $k = 2 a_n$
- $\mu = x^t m_n$
- $\tau = \dfrac{a_n}{b_n}\left(1 + x^t [S_0^{-1} + \underline{\underline{X}}^t \underline{\underline{X}}]^{-1} x\right)$

× <u>Summary</u> :

→ $Y = X\beta + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \gamma^{-1} I_n)$

→ $\gamma$ known ⟹ predictive distribution is gaussian, under a gaussian prior on $\beta$.

→ $\gamma$ unknown ⟹ predictive distribution is student, under a gaussian-gamma prior on $(\beta, \gamma)$

---

<u>Remarks</u> : (i) Conjugate prior on $(\beta, \sigma^2)$.

Sometimes it is convenient to work with the variance $\sigma^2$ directly, instead of the precision $\gamma = 1/\sigma^2$ :

$$Y = X\beta + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 I_n)$$

The conjugate prior on $(\beta, \sigma^2)$ is

$$f(\beta, \sigma^2) = \mathcal{N}(\beta \mid m_0, \sigma^2 S_0) \, Ig(\sigma^2 \mid a_0, b_0),$$

where

$Ig(x \mid \alpha, \beta)$ denotes the <u>inverse-gamma</u> distribution, with pdf :

$$f(x) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{-(\alpha+1)} e^{-\frac{\beta}{x}}, \quad x > 0$$

↖ with

$$EX = \frac{\beta}{\alpha - 1}, \quad \alpha > 1$$

$$\text{Var } X = \frac{\beta^2}{(\alpha-1)^2 (\alpha-2)}, \quad \alpha > 2$$

$$X \sim \text{Gamma}(x \mid \alpha, \beta)$$
$$\Longleftrightarrow$$
$$\frac{1}{X} \sim Ig(x \mid \alpha, \beta)$$

where we recall that

$$\text{Gamma}(x \mid \alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x}, \quad x > 0$$

with mean $\dfrac{\alpha}{\beta}$ and variance $\dfrac{\alpha}{\beta^2}$

The posterior distribution of $\beta, \sigma^2 \mid \mathcal{L}_n$ is then

$$f(\beta, \sigma^2 \mid \mathcal{L}_n) = \mathcal{N}(\beta \mid m_n, \sigma^2 S_n) \, Ig(\sigma^2 \mid a_n, b_n),$$

where $m_n, S_n, a_n, b_n$ are given on page 10.

(ii) Noninformative prior on $(\beta, \gamma)$ $[/$ on $(\beta, \sigma^2)]$ ⑮

Obtained for $(\beta, \gamma) \sim \frac{1}{\gamma}$ $\left[ (\beta, \sigma^2) \sim \frac{1}{\sigma^2} \right]$

↳ Special case of the conjugate normal-gamma distribution with $m_0$, $S_0^{-1} \to 0$, $a_0 = -\frac{d}{2}$, $b_0 \to 0$.

Indeed, take $\quad m_0 = 0 \qquad a_0 = -\frac{d}{2}$
$\qquad\qquad\qquad S_0^{-1} = \varepsilon I \qquad b_0 = \varepsilon$.

Then

$$f(\beta, \gamma) = \mathcal{N}(\beta \mid 0, \gamma^{-1} \varepsilon^{-1} I) \, gamma(\gamma \mid -\frac{d}{2}, \varepsilon)$$

*conjugate prior*

$$\propto \frac{1}{|\gamma^{-1} \varepsilon^{-1} I|^{1/2}} \exp\left\{ -\frac{1}{2} \gamma \varepsilon \beta^t \beta \right\} \underbrace{\varepsilon^{-\frac{d}{2}} \gamma^{-\frac{d}{2}-1} \underbrace{e^{-\varepsilon \gamma}}_{\downarrow\ 1}}_{\downarrow\ \text{as } \varepsilon \to 0}$$

$$\sim \gamma^{\frac{d}{2}} \varepsilon^{\frac{d}{2}} \varepsilon^{-\frac{d}{2}} \gamma^{-\frac{d}{2}-1} \qquad \text{as } \varepsilon \to 0$$

$$\sim \gamma^{-1}$$

$\Rightarrow$ The posterior distribution of $(\beta, \gamma \mid \mathcal{L}_n)$ given a noninformative prior $(\beta, \gamma) \sim \frac{1}{\gamma}$ is

$$f(\beta, \gamma \mid \mathcal{L}_n) = \mathcal{N}(\beta \mid (X^t X)^{-1} X^t y, \gamma^{-1} (X^t X)^{-1})$$
$$\times \, gamma\left(\gamma \mid \frac{n-d}{2}, \frac{n-d}{2} s^2\right),$$

with
$$s^2 = \frac{1}{n-d} (y - X\hat{\beta})^t (y - X\hat{\beta}) = \frac{1}{n-d} y^t (I-H) y$$
$$\hat{\beta} = (X^t X)^{-1} X^t y = LS \text{ estimate}$$
$$H = X(X^t X)^{-1} X^t = \text{projection matrix}$$

---

since for this choice of $m_0$, $S_0$, $a_0$, $b_0$, we get ⑯

$$S_n = (X^t X)^{-1}$$
$$m_n = (X^t X)^{-1} X^t y = \hat{\beta}$$
$$a_n = \frac{1}{2}(n-d)$$
$$b_n = \frac{1}{2}(y^t y - \hat{\beta}^t (X^t X)\hat{\beta}) = \frac{1}{2} y^t (I-H) y.$$

• The **posterior predictive distribution** is multivariate $t$:

$$y \mid \mathcal{L}_n, x_0 \sim t\left(y \mid n-d, \underset{\text{location}}{x_0^t \hat{\beta}}, \underset{\text{scale}}{s^2(1 + x_0^t (X^t X)^{-1} x_0)}\right)$$

one observation $\in \mathbb{R}^d$

$$y \mid \mathcal{L}_n, X_0 \sim t\left(y \mid n-d, X_0 \hat{\beta}, s^2(I + X_0 (X^t X)^{-1} X_0^t)\right)$$

$m$-new observations $\in \mathbb{R}^{m \times d}$

*see page 12, where the scale parameter above corresponds to $1/\tau$.*

• Likewise, a **non-informative prior** on $(\beta, \sigma^2) \sim \frac{1}{\sigma^2}$ yields the **posterior**

$$(\beta, \sigma^2 \mid \mathcal{L}_n) \sim \mathcal{N}(\beta \mid \hat{\beta}, \sigma^2 (X^t X)^{-1}) \, \text{Inv-}\chi^2(\sigma^2 \mid n-d, s^2),$$

where

$\text{Inv-}\chi^2(x \mid \upsilon, s^2)$ denotes the scaled-inverse $\chi^2$ distr; with pdf

$$f(x) = \frac{(\upsilon/2)^{\upsilon/2}}{\Gamma(\upsilon/2)} s^{\upsilon} x^{-\left(\frac{\upsilon}{2}+1\right)} e^{-\frac{\upsilon s^2}{2x}}, \quad x > 0$$

with $EX = \frac{\upsilon}{\upsilon-2} s^2$

$$Var\, X = \frac{2\upsilon^2}{(\upsilon-2)^2(\upsilon-4)} s^4$$

Indeed, the posterior distribution $\gamma | \mathcal{L}_n$ given on
page 15 can be rewritten,

$$\gamma | \mathcal{L}_n \sim \text{gamma}\left(\gamma \mid \frac{n-d}{2}, \frac{n-d}{2} s^2\right)$$

$\Longleftrightarrow$

$$\sigma^2 = \frac{1}{\gamma} \mid \mathcal{L}_n \sim Ig\left(\sigma^2 \mid \frac{n-d}{2}, \frac{n-d}{2} s^2\right),$$

whose density is precisely $\text{Inv-}\chi^2(\sigma^2 \mid n-d, s^2)$, and
is given by

$$\left|\; \frac{\left(\frac{n-d}{2}\right)^{\frac{n-d}{2}}}{\Gamma\left(\frac{n-d}{2}\right)} \; s^{n-d} \; (\sigma^2)^{-\left(\frac{n-d}{2}+1\right)} \; e^{-\frac{(n-d)s^2}{2\sigma^2}} \right.$$

(iii) <u>Sampling from the posterior distribution.</u>

To draw a sample $y$ from its posterior predictive distribution,
either use its analytical expression page 16, or
- first draw $\beta, \gamma | \mathcal{L}_n$ from its posterior distribution
- then draw $y \sim \mathcal{N}(\beta X, \gamma^{-1} I)$.

$\left(\text{since } f(y | \mathcal{L}_n, x) = \iint \underbrace{f(y, \beta, \gamma | \mathcal{L}_n, x)}\, d\beta\, d\gamma \right.$

$$= f(y | \beta, \gamma, x) \underbrace{f(\beta, \gamma | \mathcal{L}_n)}_{\text{posterior}}$$

Obtain $B$ samples $\{\beta^{(b)}, \gamma^{(b)}\}$, $b=1,..,B$ from the posterior
to get a sample $\{y^{(b)}\}$ from the posterior, which amounts to
obtaining a sample $\{\beta^{(b)}, \gamma^{(b)}, y^{(b)}\}$ & discarding the params thus
marginalizing.

## I.4. <u>Evidence Approximation</u>

- Assume that $Y = X\beta + \mathcal{E}$, $\quad \mathcal{E} \sim \mathcal{N}(0, \gamma^{-1} I_n)$
  $$\beta \sim \mathcal{N}(0, \alpha^{-1} I_d)$$

- Hyperparameters $\alpha$ and $\gamma$ are now treated as random, with
  joint prior $f(\alpha, \gamma)$. This approach is known as <u>Empirical
  Bayes</u> (EB), <u>Evidence Approximation</u>, or type $\text{II}$ <u>Maximum
  Likelihood</u>.

- Predictive distribution is
  $$f(y | \mathcal{L}_n, x) = \iiint f(y, \beta, \alpha, \gamma | \mathcal{L}_n, x)\, d\beta\, d\alpha\, d\gamma$$

  $$= \iiint \underbrace{f(y | \beta, \gamma) f(\beta | \mathcal{L}_n, \alpha, \gamma) f(\alpha, \gamma | \mathcal{L}_n)}\; d\beta\, d\alpha\, d\gamma$$

Compare with the expression on page 6: the product of the
first two terms is precisely $\mathcal{N}(y | x^t \beta, \gamma^{-1}) \mathcal{N}(\beta | m_n, S_n)$,
where $m_n, S_n$ are given on page 3, while $f(\alpha, \gamma | \mathcal{L}_n)$
represents the posterior distribution on the hyperparameters:

$$f(\alpha, \gamma | \mathcal{L}_n) \propto \underbrace{f(\mathcal{L}_n | \alpha, \gamma)}_{\substack{\text{"marginal} \\ \text{likelihood"}}} \underbrace{f(\alpha, \gamma)}_{\text{prior}}$$

As $n$ gets larger, the posterior is more and more peaked
around $(\hat{\alpha}, \hat{\gamma}) = \underset{(\alpha, \gamma)}{\text{argmax}}\, f(\mathcal{L}_n | \alpha, \gamma)$. We may
substitute $(\hat{\alpha}, \hat{\gamma})$ back into the expression of the predictive distribution,
and obtain the approximation

$$f(y | \mathcal{L}_n, x) \approx \int f(y | \beta, \hat{\gamma}) f(\beta | \mathcal{L}_n, \hat{\alpha}, \hat{\gamma})\, d\beta$$

$\Rightarrow \boxed{\begin{array}{l} f(y \mid \mathcal{L}_n, x) \approx \mathcal{N}(y \mid m_n^t x, \hat{\gamma}^{-1} + x^t S_n x), \\[2mm] \text{with} \quad m_n = \hat{\gamma} S_n X^t y, \quad S_n = \hat{\alpha} I_d + \hat{\gamma} X^t X \end{array}}$

The hyperparameters $\alpha$ and $\gamma$ are selected from $\mathcal{L}_n$ directly. No need for cross-validation here; hence the name type II ML : $(\hat{\alpha}, \hat{\gamma})$ maximize the marginal likelihood $f(\mathcal{L}_n \mid \alpha, \gamma)$. It remains to find their value.

$$f(\mathcal{L}_n \mid \alpha, \gamma) = \int \underbrace{f(\mathcal{L}_n \mid \beta, \alpha, \gamma)}_{} \underbrace{f(\beta \mid \alpha, \gamma)}_{} d\beta$$

$$= \prod_{i=1}^n \mathcal{N}(y_i \mid x_i^t \beta, \gamma^{-1}) \qquad = \mathcal{N}(\beta \mid 0, \alpha^{-1} I_d)$$

The integral can be evaluate using the same formula as given at the bottom of page 6. After calculations, we get

$$f(\mathcal{L}_n \mid \alpha, \gamma) = \mathcal{N}(\mathcal{L}_n \mid 0, \gamma^{-1} I_n + \alpha^{-1} X X^t).$$

$$\log f(\mathcal{L}_n \mid \alpha, \gamma) = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \underbrace{\log \mid \gamma^{-1} I_n + \alpha^{-1} X X^t \mid}_{}$$
$$\qquad\qquad\qquad - \frac{1}{2} \underbrace{y^t (\gamma^{-1} I_n + \alpha^{-1} X X^t)^{-1} y}_{}.$$

$\downarrow$ (C.7)

$= \gamma^{-n} \mid I_n + \frac{\gamma}{\alpha} X X^t \mid$   (C.14) page 697 in Bishop

$y^t (\gamma I_n - \gamma X [\alpha I_d + \gamma X^t X]^{-1} X^t \gamma) y$

$\parallel$

$\gamma^{-n} \alpha^{-d} \mid \alpha I_d + \gamma X^t X \mid$

$\parallel$

$\gamma y^t y - \gamma^2 y^t X A^{-1} X^t y$

$\gamma^{-n} \alpha^{-d} \mid A \mid$

$\parallel$

where $A := \alpha I_d + \gamma X^t X$

$\gamma y^t y - \hat{\beta}^t A \hat{\beta}$ $\quad \hat{\beta} = \gamma A^{-1} X^t y$

---

Note that $\hat{\beta} = (X^t X + \frac{\alpha}{\gamma} I_d)^{-1} X^t y$; the ridge solution.

$$\Rightarrow \log f(\mathcal{L}_n \mid \alpha, \gamma) = -\frac{n}{2} \log(2\pi) + \frac{n}{2} \log \gamma + \frac{d}{2} \log \alpha - \frac{1}{2} \log \mid A \mid$$
$$\boxed{- \frac{\gamma}{2} y^t y + \frac{1}{2} \hat{\beta}^t A \hat{\beta}}.$$

$$= \frac{1}{2} (\gamma y^t y - \hat{\beta}^t A \beta)$$

$$= \frac{1}{2} (\gamma y^t y - 2 \hat{\beta}^t A \hat{\beta} + \hat{\beta}^t A \hat{\beta})$$

$$= \frac{1}{2} (\gamma y^t y - 2 \hat{\beta}^t A A^{-1} X^t y \gamma + \hat{\beta}^t (\alpha I + \gamma X^t X) \hat{\beta})$$

$$= \frac{1}{2} (\gamma y^t y - 2 \hat{\beta}^t X^t y \gamma + \alpha \hat{\beta}^t \hat{\beta} + \gamma \hat{\beta}^t X^t X \hat{\beta})$$

$$= \frac{1}{2} (\gamma \| y - X \hat{\beta} \|^2 + \alpha \| \hat{\beta} \|^2)$$

$$= \frac{\gamma}{2} (\| y - X \hat{\beta} \|^2 + \frac{\alpha}{\gamma} \| \hat{\beta} \|^2)$$

$$= \frac{\gamma}{2} \times \text{Penalized Residual Sum of Squares of the Ridge Solution } \mu_n.$$

$$\boxed{\begin{array}{l} \log f(\mathcal{L}_n \mid \alpha, \gamma) = -\frac{n}{2} \log(2\pi) + \frac{n}{2} \log \gamma + \frac{d}{2} \log \alpha - \frac{1}{2} \log \mid A \mid \\[2mm] \qquad\qquad\qquad - \frac{\gamma}{2} \| y - X \hat{\beta} \|^2 - \frac{\alpha}{2} \| \hat{\beta} \|^2, \\[2mm] \text{where } \hat{\beta} = (X^t X + \frac{\alpha}{\gamma} I_d)^{-1} X^t y \end{array}}$$

**MARGINAL LIKELIHOOD**

We are looking for $(\hat{\alpha}, \hat{\gamma})$ maximizing the marginal likelihood
$\hookrightarrow \mathcal{L}_n$ is used to estimate <u>hyperparameters</u> $\to$ Type II ML.

(VS) $[\mathcal{L}_n$ being used to estimate <u>parameter</u> $\beta \to$ ML $]$

↳ First, consider the maximization with respect to $\alpha$.

Let $(l_i, v_i)$ = eigenvalue – eigenvector pairs of $\gamma X^t X$.

Since $A = \alpha I_d + \gamma X^t X$, $A$ has eigenvalue – eigenvector pairs $(\alpha + l_i, v_i)$.

Moreover, $|A| = \prod_{i=1}^{d} (l_i + \alpha) \Rightarrow \log|A| = \sum_{i=1}^{d} \log(\alpha + l_i)$,

and $\dfrac{d}{d\alpha} \log|A| = \sum_{i=1}^{d} \dfrac{1}{\alpha + l_i}$.

$\Rightarrow \dfrac{\partial}{\partial\alpha} \log f(\mathcal{L}_n \mid \alpha, \gamma) = \dfrac{d}{2\hat\alpha} - \dfrac{1}{2}\hat\beta^t\hat\beta - \dfrac{1}{2}\sum_{i=1}^{d}\dfrac{1}{l_i + \hat\alpha} = 0$

$\hat\beta$ depends on $\alpha$.
[we neglected the derivative of $\hat\beta$ with respect to $\alpha$]

$\boxed{\hat\alpha\, \hat\beta^t\hat\beta = d - \hat\alpha \sum_{i=1}^{d}\dfrac{1}{l_i + \hat\alpha} = \sum_{i=1}^{d}\dfrac{l_i}{l_i + \hat\alpha}}$ (∗)

$\hat\alpha$ satisfies this equation. It can be solved iteratively.

↳ Next, consider maximization w.r.t. $\gamma$.

Since $X^t X v_i = \dfrac{l_i}{\gamma} v_i =: \lambda_i v_i$, we see that

$\dfrac{d l_i}{d\gamma} = \lambda_i = \dfrac{l_i}{\gamma} \Rightarrow \dfrac{d}{d\gamma}\log|A| = \dfrac{1}{\gamma}\sum_{i=1}^{d}\dfrac{l_i}{l_i + \alpha}$

$\Rightarrow \dfrac{\partial}{\partial\gamma}\log f(\mathcal{L}_n \mid \alpha, \gamma) = \dfrac{n}{2\gamma} - \dfrac{1}{2}\sum_{i=1}^{n}(y_i - \hat\beta^t x_i)^2 - \dfrac{1}{2\gamma}\sum\dfrac{l_i}{l_i + \alpha}$

Equating to zero

$\boxed{\dfrac{1}{\gamma} = \dfrac{1}{n - \gamma}\sum_{i=1}^{n}(y_i - \hat\beta^t x_i)^2}$ ← solve recursively as well.

---

**Remark** = Evidence Approximation & EM algorithm.

There are close ties between the recursions derived on page 17, and the EM algorithm. Recall that (see UL: CLUSTERING p. 23)

**Goal:** maximize the log-likelihood $\ell(Q) = \log f(x \mid Q)$

**E-step:** compute $Q(Q, Q^{(m)}) = \mathbb{E}_{f(z \mid x, Q^{(m)})}\{\log f(x, z \mid Q)\}$

**M-step:** $Q^{(m+1)} = \underset{Q}{\text{argmax}}\ Q(Q, Q^{(m)})$.

where $Q$ = parameter of interest
$x$ = observed variable
$z$ = latent variable.

**Fact:** $\ell(Q^{(m)}) \leq \ell(Q^{(m+1)})$.

↳ In our context, we want to maximize the marginal (log)-likelihood

$f(\mathcal{L}_n \mid \alpha, \gamma) = \int f(\mathcal{L}_n \mid \beta, \alpha, \gamma)\, f(\beta \mid \alpha, \gamma)\, d\beta$

↳ our latent variable $z$

↳ The complete (log) likelihood is:

$\log f(\mathcal{L}_n, \beta \mid \alpha, \gamma) = \boxed{\log f(\mathcal{L}_n \mid \beta, \alpha, \gamma)}$  "$\mathcal{N}(y \mid X\beta, \gamma^{-1} I_n)$"
$+$
$\boxed{\log f(\beta \mid \alpha, \gamma)}$  "$\mathcal{N}(\beta \mid 0, \alpha^{-1} I_d)$"

$= \boxed{\dfrac{d}{2}\log\left(\dfrac{\alpha}{2\pi}\right) - \dfrac{\alpha}{2}\beta^t\beta} + \boxed{\dfrac{n}{2}\log\dfrac{\gamma}{2\pi} - \dfrac{\gamma}{2}\sum_{i=1}^{n}(y_i - \beta^t x_i)^2}$

↪ E-step =

$$\mathbb{E}\left\{ \log f(\mathcal{L}_n, \beta \mid \alpha, \gamma) \mid \mathcal{L}_n, \alpha^{(m)}, \gamma^{(m)} \right\}$$

<span style="color:green">current parameter values.</span>

$$= \frac{d}{2} \log\left(\frac{\alpha}{2\pi}\right) + \frac{n}{2}\frac{\gamma}{2\pi} - \frac{\alpha}{2}\mathbb{E}\left\{\beta^t\beta \mid \mathcal{L}_n, \alpha^{(m)}, \gamma^{(m)}\right\}$$

$$- \frac{\gamma}{2}\sum_{i=1}^{n}\mathbb{E}\left\{(y_i - \beta^t x_i)^2 \mid \mathcal{L}_n, \alpha^{(m)}, \gamma^{(m)}\right\}$$

where

- $\mathbb{E}\left\{\beta^t\beta \mid \mathcal{L}_n, \alpha^{(m)}, \gamma^{(m)}\right\} = m_n^t m_n + S_n$,

with $\begin{pmatrix} m_n = \gamma^{(m)} S_n X^t y \; (= \hat\beta) \\ S_n^{-1} = \alpha^{(m)} I_d + \gamma^{(m)} X^t X \end{pmatrix}$, see page 3.

- $\mathbb{E}\left\{(y_i - \beta^t x_i)^2 \mid \mathcal{L}_n, \alpha^{(m)}, \gamma^{(m)}\right\}$

$$= y_i^2 - 2y_i x_i^t m_n + \underbrace{\mathbb{E}\left\{\mathrm{Tr}(x_i x_i^t \beta\beta^t) \mid -''-\right\}}$$

$$\mathrm{Tr}\left\{ x_i x_i^t \, \mathbb{E}(\beta\beta^t \mid -''-)\right\}$$

$$\mathrm{Tr}(x_i x_i^t S_n) + m_n^t x_i x_i^t m_n$$

$$= (y_i - m_n^t x_i)^2 + x_i^t S_n x_i$$

⇒ $$Q(\theta, \theta^{(m)}) = \frac{d}{2}\log\left(\frac{\alpha}{2\pi}\right) + \frac{n}{2}\log\frac{\gamma}{2\pi} - \frac{\alpha}{2}\left(m_n^t m_n + S_n\right)$$

$$- \frac{\gamma}{2}\sum_{i=1}^{n}\left((y_i - m_n^t x_i)^2 + x_i^t S_n x_i\right)$$

---

↪ M-step = • $\dfrac{\partial}{\partial \alpha} Q(\theta, \theta^{(m)}) = \dfrac{d}{2\alpha^{(m+1)}} - \dfrac{1}{2}(m_n^t m_n + S_n)$ <span style="color:red">= 0</span>

$$\Rightarrow \quad \alpha^{(m+1)} = \frac{d}{m_n^t m_n + S_n}$$

• $\dfrac{\partial}{\partial \gamma} Q(\theta, \theta^{(m)}) = \dfrac{n}{2\gamma^{(m+1)}} - \dfrac{1}{2}\sum_{i=1}^{n}\left((y_i - m_n^t x_i)^2 + x_i^t S_n x_i\right)$ <span style="color:red">= 0</span>

$$\Rightarrow \quad \gamma^{(m+1)} = \frac{n}{\sum_{i=1}^{n}(y_i - m_n^t x_i)^2 + x_i^t S_n x_i}$$

Now compare the M-step of the EM algorithm for $\alpha$ with recursion <span style="color:red">(*)</span> page 17: $\hat\alpha \, m_n^t m_n = d - \hat\alpha \underbrace{\sum_{i=1}^{d}\frac{1}{l_i + \hat\alpha}}_{= \mathrm{Tr}(S_n)}$

Re-arranging terms, this is precisely the recursion above for $\alpha$.

since $l_i = $ eigenvalue of $\gamma X^t X$ and $S_n^{-1} = \alpha I_d + \gamma X^t X$.

## II. BAYESIAN LOGISTIC REGRESSION

We consider the binary classification problem; $X \in \mathbb{R}^d$, $Y \in \{0, 1\}$, and

$$\log\left\{\frac{\mathbb{P}(Y=1 \mid X=x)}{\mathbb{P}(Y=0 \mid X=x)}\right\} = \beta^t x \; ; \; \beta \in \mathbb{R}^d$$

In a Bayesian framework, we put a Gaussian prior on $\beta$,

and assume that $\beta \sim f(\beta) = \mathcal{N}(\beta \mid m_0, S_0)$

## II.1. Posterior distribution.

- The posterior distribution of $\beta$ given $\mathcal{L}_n$ is proportional to the product $f(\mathcal{L}_n \mid \beta) f(\beta)$, so that

$$\log f(\beta \mid \mathcal{L}_n) = \sum_{i=1}^{n} \left( y_i \log \sigma_i + (1-y_i) \log(1-\sigma_i) \right)$$
$$- \frac{1}{2}(\beta - m_0)^t S_0^{-1}(\beta - m_0) + \text{constant in } \beta,$$

where $\sigma_i = \sigma(\beta^t x_i)$, $\sigma$ = sigmoid function.

- We are looking for a Gaussian approximation of the posterior. Denote it $q(\beta) = \mathcal{N}(\beta \mid m_n, S_n) \approx f(\beta \mid \mathcal{L}_n)$.

↗ Laplace approximation

⇘ $m_n$ = value of $\beta$ maximizing $\log f(\beta \mid \mathcal{L}_n)$
  = MAP estimate

⇘ $S_n$ is obtained by considering a Taylor expansion of log of $f(\beta \mid \mathcal{L}_n)$ around its mode $m_n$:

$$\log f(\beta \mid \mathcal{L}_n) \simeq \log f(m_n \mid \mathcal{L}_n) - \frac{1}{2}(\beta - m_n)^t S_n^{-1}(\beta - m_n),$$

since the derivative / gradient of $\log f(\beta \mid \mathcal{L}_n)$ vanishes at its mode $m_n$.

where $S_n^{-1} = -\nabla_\beta^2 \{ -\log f(\beta \mid \mathcal{L}_n) \}$

$$= S_0^{-1} + X^t W X; \qquad W = \begin{pmatrix} \sigma_1(1-\sigma_1) & & 0 \\ & \ddots & \\ 0 & & \sigma_n(1-\sigma_n) \end{pmatrix}$$

see p.14 in SL: LINEAR CLASSIFIERS

evaluated at $m_n$

---

## II.2. Predictive distribution.

The predictive distribution is given by

$$\mathbb{P}(Y=1 \mid \mathcal{L}_n, x) = \int \underbrace{\mathbb{P}(Y=1 \mid x, \beta)}_{= \ \sigma(\beta^t x)} \underbrace{f(\beta \mid \mathcal{L}_n)}_{SS} \, d\beta$$

new point

$q(\beta) = \mathcal{N}(\beta \mid m_n, S_n)$

① **Plug-in approximation** =

$$\mathbb{P}(Y=1 \mid \mathcal{L}_n, x) \approx \mathbb{P}(Y=1 \mid x, \beta=m_n) = \sigma(m_n^t x)$$

② **MC simulations** =

Consider independent samples $\beta_i \sim q(\beta)$. Then

$$\mathbb{P}(Y=1 \mid \mathcal{L}_n, x) \approx \frac{1}{M} \sum_{i=1}^{M} \sigma(\beta_i^t x).$$
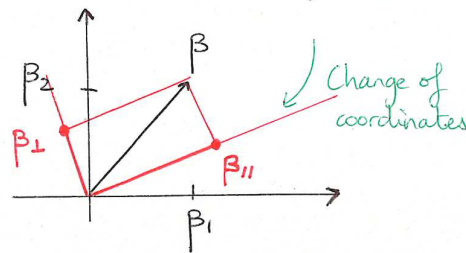
③ **Probit approximation** =

First, we compute a more tractable expression for the integral

$$\int \sigma(\beta^t x) \, q(\beta) \, d\beta,$$

by expressing $\beta = \begin{pmatrix} \beta_{\parallel} \\ \beta_{\perp} \end{pmatrix} \begin{smallmatrix} \Uparrow 1 \\ \Uparrow d-1 \end{smallmatrix}$ with $\beta_{\parallel} := \dfrac{\langle \beta, x \rangle}{\sqrt{\langle x, x \rangle}}$

= projection of $\beta$ onto $x$

Change of coordinates

$$\int \sigma(\beta^t x) \, q(\beta) \, d\beta = \iint \sigma(\beta_{\parallel} \|x\|) \underbrace{q(\beta_\perp \mid \beta_{\parallel})} \, q(\beta_{\parallel}) \, d\beta_{\parallel} \, \underbrace{d\beta_\perp}$$

integrates to 1

$\Rightarrow \int \sigma(\beta^t x)\, q(\beta)\, d\beta = \int \sigma(\beta_{\parallel}\, \|x\|)\, q(\beta_{\parallel})\, d\beta_{\parallel}$

Since $\beta \sim$ multivariate gaussian, $q(\beta_{\parallel})$ is $\mathcal{N}(\beta_{\parallel} \mid \mu_n, \sigma_n^2)$, with

$\mu_n := \mathbb{E}(\beta_{\parallel} \mid \mathcal{L}_n) = \mathbb{E}\left(\dfrac{x^t \beta}{\|x\|} \mid \mathcal{L}_n\right) = \dfrac{x^t m_n}{\|x\|}$, where

$\qquad\qquad\qquad\qquad m_n = \text{posterior mean} = \text{MAP}$

$\sigma_n^2 := \mathbb{E}\left\{\left(\dfrac{x^t(\beta - m_n)}{\|x\|}\right)^2 \mid \mathcal{L}_n\right\} \quad \left(= \mathbb{E}\{(\beta_{\parallel} - \mathbb{E}\beta_{\parallel})^2 \mid \mathcal{L}_n\}\right)$

$\qquad = \dfrac{x^t}{\|x\|}\, \mathbb{E}\{(\beta - m_n)(\beta - m_n)^t \mid \mathcal{L}_n\}\, x = \dfrac{x^t S_n x}{\|x\|^2}$.

Thus, $\quad q(\beta_{\parallel}) = \mathcal{N}\left(\beta_{\parallel} \mid \dfrac{x^t m_n}{\|x\|},\ \dfrac{x^t S_n x}{\|x\|^2}\right)$.

$\Rightarrow$ With $s := \beta_{\parallel}\, \|x\|$,

$\int \sigma(\beta_{\parallel}\, \|x\|)\, q(\beta_{\parallel})\, d\beta_{\parallel} = \int \sigma(s)\, \underbrace{q\left(\dfrac{s}{\|x\|}\right) \dfrac{ds}{\|x\|}}$

$\qquad\qquad\qquad$ density of

$\qquad\qquad\qquad \beta_{\parallel}\, \|x\| \sim \mathcal{N}(x^t m_n,\ x^t S_n x)$

We finally get

$$\boxed{\mathbb{P}(Y=1 \mid \mathcal{L}_n, x) = \int \sigma(s)\, \mathcal{N}(s \mid x^t m_n,\ x^t S_n x)\, ds}$$

The probit approximation of this integral uses $\Phi$ in place of $\sigma$. Indeed, $\sigma(s) \approx \Phi\left(\sqrt{\dfrac{\pi}{8}}\, s\right)$, which is easily obtained by equating the slope of the two functions at the origin.

$\Rightarrow \mathbb{P}(Y=1 \mid \mathcal{L}_n, x) \approx \int \Phi\left(\sqrt{\dfrac{\pi}{8}}\, s\right) \mathcal{N}(s \mid x^t m_n,\ x^t S_n x)\, ds$,

which can be analytically computed.

Indeed, consider $X \sim \mathcal{N}(0, \lambda^{-2})$, $Y \sim \mathcal{N}(m, \sigma^2)$, independent.

$\mathbb{P}(X \leq Y) = \mathbb{E}_Y\, \mathbb{P}(X \leq y) = \int \Phi(\lambda y)\, \mathcal{N}(y \mid m, \sigma^2)\, dy$.

On the other hand, $X - Y \sim \mathcal{N}(-m, \lambda^{-2} + \sigma^2)$, so that

$\mathbb{P}(X \leq Y) = \mathbb{P}(X - Y \leq 0) = \Phi\left(\dfrac{m}{(\lambda^{-2} + \sigma^2)^{1/2}}\right)$.

$$\Rightarrow \boxed{\mathbb{P}(Y=1 \mid \mathcal{L}_n, x) \approx \Phi\left(\dfrac{x^t m_n}{(\lambda^{-2} + x^t S_n x)^{1/2}}\right) ;\quad \lambda = \sqrt{\dfrac{\pi}{8}}}$$

use the sigmoid approximation again $\qquad \approx \sigma\left(x^t m_n\left(1 + \dfrac{\pi}{8} x^t S_n x\right)^{-1/2}\right)$

Remark: Classify a new observation $x$ as $1$ if

$\mathbb{P}(Y=1 \mid x, \mathcal{L}_n) \geq \mathbb{P}(Y=0 \mid x, \mathcal{L}_n)$

$\Updownarrow$

$\sigma\left(x^t m_n\, (\text{---}''\text{---})^{-1/2}\right) \geq \frac{1}{2}$

$\Updownarrow$

$x^t m_n \geq 0$

If the objective is the minimization of the misclassification rate, with equal prior probabilities, then the marginalization over $\beta$ in the computation of the predictive distribution has no effect.

## Metropolis - Hastings algorithm.

The previous techniques required the approximation of the posterior distribution. There exists techniques to sample directly from the posterior without requiring to approximate it first. Metropolis - Hastings algorithm (see MS : MCMC) is one of them.

×  **Goal**: to approximate the integral $\int h(\theta) f(\theta) d\theta$

×  **Idea**: generate samples $\sim$ density $f$ : $\theta_1, \theta_2, \dots$
the integral is then $\approx \frac{1}{M} \sum_{i=1}^{M} h(\theta_i)$.

×  **How**: start from $\theta_0$, and generate $\theta_i$ using a transition kernel, with target density $f$.

  ≫ $f$ is known up to a multiplicative constant
  ↘ choose a proposal density $q(y \mid \theta)$
and proceed as **follows**

---

Given $\theta_i$

(i) generate $y_i \sim q(y \mid \theta_i)$

(ii) accept / reject

$$\theta_{i+1} = \begin{cases} y_i & \text{w.p. } \rho(\theta_i, y_i) \\ \theta_i & \text{w.p. } 1 - \rho(\theta_i, y_i) \end{cases}$$

where

$$\rho(\theta, y) = \min \left\{ \frac{f(y)}{f(\theta)} \frac{q(\theta \mid y)}{q(y \mid \theta)}, 1 \right\}$$

**METROPOLIS – HASTINGS  ALGORITHM**

---

*the sequence can take several times the same value (non iid sample)*

*If $q(\theta \mid y) = q(y \mid \theta)$ (symmetric case), always accept points $y_i$ increasing the "likelihood".*

---

Under some general conditions [ such as the event $\{\theta_i = \theta_{i+1}\}$ is possible, $q(y \mid \theta) > 0 \quad \forall (\theta, y)$ ], then we have

- ergodicity $\quad \frac{1}{M} \sum_{i=1}^{M} h(\theta_i) \longrightarrow \int h(\theta) f(\theta) d\theta$

- convergence in total variation. In particular,
$$\mathbb{P}(\theta_i \in B) \longrightarrow \int_B f(\theta) d\theta$$

→ In Bayesian logistic Regression, the target density is the posterior
$$f(\beta \mid \mathcal{L}_n) = \frac{f(\mathcal{L}_n \mid \beta) f(\beta)}{f(\mathcal{L}_n)}.$$

In MH, we need to compute the ratios,

$$\frac{f(\beta_1 \mid \mathcal{L}_n)}{f(\beta_2 \mid \mathcal{L}_n)} = \frac{f(\mathcal{L}_n \mid \beta_1) f(\beta_1)}{f(\mathcal{L}_n \mid \beta_2) f(\beta_2)}, \text{ which is known to us.}$$

In addition, using the symmetric proposal distribution
$$q(\beta_1 \mid \beta_2) = \mathcal{N}(\beta_1 \mid \beta_2, \varsigma^2 I_d)$$
$$= \mathcal{N}(\beta_2 \mid \beta_1, \varsigma^2 I_d) = q(\beta_2 \mid \beta_1),$$

we obtain

---

**MH FOR BAYESIAN LOGISTIC REG.**

Given $\beta_i$

(i) Generate $\gamma_i \sim q(\gamma \mid \beta_i)$

(ii) Take
$$\beta_{i+1} = \begin{cases} \gamma_i & \text{w.p. } \rho(\beta_i, \gamma_i) \\ \beta_i & \text{w.p. } 1 - \rho(\beta_i, \gamma_i) \end{cases}$$
where
$$\rho(\beta, \gamma) = \min \left\{ \frac{f(\mathcal{L}_n \mid \gamma) f(\gamma)}{f(\mathcal{L}_n \mid \beta) f(\beta)}, 1 \right\}$$

---

And compute
$$\mathbb{P}(Y = 1 \mid \mathcal{L}_n, x) \leftarrow$$

$$\frac{1}{M} \sum_{i=1}^{M} \overset{\sim}{\sigma}(\beta_i^t x).$$