

Part I: Gaussian Distribution.

1^o Gaussian Distribution in Vector form.

$$N(x|\mu, \Sigma) = \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu) \right\}$$

$x \in \mathbb{R}^D$; $\mu \in \mathbb{R}^D$; $\Sigma \in \mathbb{R}^{D \times D}$

2^o Conditional Gaussian

$$x \in \mathbb{R}^D \quad x \triangleq \begin{bmatrix} x_a \\ x_b \end{bmatrix} \quad x_a \in \mathbb{R}^m \quad x_b \in \mathbb{R}^{D-m}$$

Equiv.

$$\mu \in \mathbb{R}^D \quad \mu \triangleq \begin{bmatrix} \mu_a \\ \mu_b \end{bmatrix} \quad \mu_a \in \mathbb{R}^m \quad \mu_b \in \mathbb{R}^{D-m}$$

$$\Rightarrow \Sigma \in \mathbb{R}^{D \times D} \quad \Sigma = \begin{bmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{bmatrix}$$

$$\Sigma_{aa} \in \mathbb{R}^{m \times m} \quad \Sigma_{ab} \in \mathbb{R}^{m \times (D-m)} \quad \Sigma_{ba} \in \mathbb{R}^{(D-m) \times m} \quad \Sigma_{bb} \in \mathbb{R}^{(D-m) \times (D-m)}$$

$$\Sigma_{ab}^T = \Sigma_{ba}$$

$$\Lambda := \Sigma^{-1} \Leftarrow \text{Precision Matrix}$$

for the quadratic term.

$$\begin{aligned} -\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu) &= -\frac{1}{2} \left[\begin{pmatrix} x_a \\ x_b \end{pmatrix} - \begin{pmatrix} \mu_a \\ \mu_b \end{pmatrix} \right]^T \begin{bmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{bmatrix}^{-1} \left[\begin{pmatrix} x_a \\ x_b \end{pmatrix} - \begin{pmatrix} \mu_a \\ \mu_b \end{pmatrix} \right] \\ &= -\frac{1}{2} \begin{bmatrix} x_a - \mu_a \\ x_b - \mu_b \end{bmatrix}^T \begin{bmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{bmatrix}^{-1} \begin{bmatrix} x_a - \mu_a \\ x_b - \mu_b \end{bmatrix} \\ &= -\frac{1}{2} (x_a - \mu_a)^T \Sigma_{aa}^{-1} (x_a - \mu_a) - \frac{1}{2} (x_a - \mu_a)^T \Sigma_{ab}^{-1} (x_b - \mu_b) \\ &\quad - \frac{1}{2} (x_b - \mu_b)^T \Sigma_{ba}^{-1} (x_a - \mu_a) - \frac{1}{2} (x_b - \mu_b)^T \Sigma_{bb}^{-1} (x_b - \mu_b) \end{aligned}$$

Let the above be seen as a function of x_a :

$$\begin{aligned}
 P(x_a | x_b) &= -\frac{1}{2} (x_a - \mu_a)^T \tilde{\Sigma}_{aa}^{-1} (x_a - \mu_a) - \frac{1}{2} (x_a - \mu_a)^T \tilde{\Sigma}_{ab}^{-1} (x_b - \mu_b) \\
 &\quad - \frac{1}{2} (x_b - \mu_b)^T \tilde{\Sigma}_{ba}^{-1} (x_a - \mu_a) - \frac{1}{2} (x_b - \mu_b)^T \tilde{\Sigma}_{bb}^{-1} (x_b - \mu_b) \\
 &= -\frac{1}{2} x_a^T \Lambda_{aa} x_a + x_a^T \left\{ \Lambda_{aa} \mu_a - \Lambda_{ab} (x_b - \mu_b) \right\} + \text{const} \\
 &= -\frac{1}{2} x_a^T \Lambda_{aa} x_a + x_a^T \Lambda_{aa}^{-1} \Lambda_{aa} \left\{ \Lambda_{aa} \mu_a - \Lambda_{ab} (x_b - \mu_b) \right\} + \text{const}.
 \end{aligned}$$

For the quadratic term of the Gaussian:

$$-\frac{1}{2} (x - \mu)^T \tilde{\Sigma}^{-1} (x - \mu) = -\frac{1}{2} x^T \tilde{\Sigma}^{-1} x + x^T \tilde{\Sigma}^{-1} \mu + \text{const}$$

$$\therefore \tilde{\Sigma}_{ab} = \Lambda_{aa}^{-1} = \tilde{\Sigma}_{aa}$$

$$\begin{aligned}
 \mu_{ab} &= \Lambda_{aa}^{-1} \left\{ \Lambda_{aa} \mu_a - \Lambda_{ab} (x_b - \mu_b) \right\} \\
 &= \mu_a - \Lambda_{aa}^{-1} \Lambda_{ab} (x_b - \mu_b)
 \end{aligned}$$

$$\Rightarrow P(x_a | x_b) \sim N(\mu_{ab}, \tilde{\Sigma}_{ab}^{-1})$$

Lemma:

$$\begin{bmatrix} A & B \\ C & D \end{bmatrix}^{-1} = \begin{bmatrix} M & -M D^{-1} \\ -D^{-1} C M & D^{-1} + C M D^{-1} \end{bmatrix}$$

$$M = (A - BD^{-1}C)^{-1} \quad \text{"scher Komponente"}$$

$$\Rightarrow \begin{bmatrix} \tilde{\Sigma}_{aa} & \tilde{\Sigma}_{ab} \\ \tilde{\Sigma}_{ba} & \tilde{\Sigma}_{bb} \end{bmatrix} = \begin{bmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{bmatrix}$$

$$\Lambda_{aa} = (\bar{\Sigma}_{aa} - \bar{\Sigma}_{ab} \bar{\Sigma}_{bb}^{-1} \bar{\Sigma}_{ba})^{-1}$$

$$\Lambda_{ab} = -(\bar{\Sigma}_{aa} - \bar{\Sigma}_{ab} \bar{\Sigma}_{bb}^{-1} \bar{\Sigma}_{ba})^{-1} \bar{\Sigma}_{ab} \bar{\Sigma}_{bb}^{-1}$$

$$\Rightarrow \bar{\Sigma}_{a|b} = \Lambda_{aa}^{-1} = \bar{\Sigma}_{aa} - \bar{\Sigma}_{ab} \bar{\Sigma}_{bb}^{-1} \bar{\Sigma}_{ba} = \text{Var}(a) - \frac{\text{cov}(a, b)}{\text{Var}(b)}$$

$$m_{a|b} = m_a - \Lambda_{aa}^{-1} \Lambda_{ab} (x_b - m_b)$$

$$= m_a + \bar{\Sigma}_{ab} \bar{\Sigma}_{bb}^{-1} (x_b - m_b)$$

4

Lemma: If a joint distribution $p(x_a, x_b)$ is Gaussian, then the conditional distribution is Gaussian as well.

3^o Marginal Gaussian

$$p(x_a) = \int p(x_a, x_b) dx_b$$

For the quadratic term in $p(x_a, x_b)$:

$$-\frac{1}{2} (x_a - m_a)^T \bar{\Sigma}_{aa}^{-1} (x_a - m_a) - \frac{1}{2} (x_a - m_a)^T \bar{\Sigma}_{ab}^{-1} (x_b - m_b)$$

$$-\frac{1}{2} (x_b - m_b)^T \bar{\Sigma}_{bb}^{-1} (x_a - m_a) - \frac{1}{2} (x_b - m_b)^T \bar{\Sigma}_{bb}^{-1} (x_b - m_b)$$

For terms involving x_b :

$$-\frac{1}{2} x_b^T \bar{\Sigma}_{bb}^{-1} x_b + x_b^T [\Lambda_{bb} m_b - \Lambda_{ba} (x_a - m_a)]$$

$$\text{let } m := \Lambda_{bb} m_b - \Lambda_{ba} (x_a - m_a)$$

$$\Rightarrow -\frac{1}{2} x_b^T \bar{\Sigma}_{bb}^{-1} x_b + x_b^T m = -\frac{1}{2} (x_b - \Lambda_{bb}^{-1} m)^T \Lambda_{bb} (x_b - \Lambda_{bb}^{-1} m) + \frac{1}{2} m^T \Lambda_{bb}^{-1} m$$

$$\int \exp \left\{ -\frac{1}{2} (x_b - \Lambda_{bb}^{-1} m)^T \Lambda_{bb} (x_b - \Lambda_{bb}^{-1} m) \right\} dx_b$$

= the reciprocal of the normalization coefficient of
the Gaussian distribution.

For the left term $\frac{1}{2} m^T \Lambda_{bb} m$ and the remaining terms in the quadratic form:

$$\frac{1}{2} [\Lambda_{bb} x_b - \Lambda_{ba} (x_a - m_a)]^T \Lambda_{bb}^{-1} [\Lambda_{bb} x_b - \Lambda_{ba} (x_a - m_a)]$$

$$-\frac{1}{2} x_a^T \Lambda_{aa} x_a + x_a^T (\Lambda_{aa} m_a + \Lambda_{ab} m_b) + \text{const.}$$

$$= -\frac{1}{2} x_a^T (\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba}) x_a + x_a^T (\Lambda_{aa} m_a - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba}) m_a + \text{const}$$

\Rightarrow The mean and variance of $p(x_a)$:

$$\bar{x}_a = (\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba})^{-1}$$

$$m_a = \bar{x}_a (\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba}) m_a$$

$$\therefore \begin{bmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{bmatrix}^{-1} = \begin{bmatrix} \bar{x}_{aa} & \bar{x}_{ab} \\ \bar{x}_{ba} & \bar{x}_{bb} \end{bmatrix}$$

$$\therefore \bar{x}_{aa} = (\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba})^{-1}$$

$$\therefore E[x_a] = m_a \quad \text{cov}(x_a) = \bar{x}_{aa}.$$

Summary:

$$x \sim N(x|\mu, \Sigma) ; \Lambda := \Sigma^{-1}$$

$$x = \begin{pmatrix} x_a \\ x_b \end{pmatrix} \quad \mu = \begin{pmatrix} \mu_a \\ \mu_b \end{pmatrix} \quad \Sigma = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}$$

$$\Lambda = \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix}$$

$$\Rightarrow P(x_a|x_b) = N(x| \mu_{ab}, \Lambda_{aa})$$

$$\mu_{ab} = \mu_a - \Lambda_{aa}^{-1} \Lambda_{ab} (x_b - \mu_b)$$

$$P(x_a) = \int P(x_a, x_b) dx_b = N(x_a | \mu_a, \Sigma_{aa})$$

4^o Bayesian's theorem for Gaussian Random Variables.

Suppose, we are given a Gaussian marginal distribution $p(x)$ and a Gaussian conditional distribution $P(y|x)$; $P(y|x)$ has a mean that is linear function of x and a covariance that is independent of x "Linear Gaussian Model".

We want to find $P(y)$ and $P(x|y)$.

$$P(x) = N(x | \mu, \Lambda)$$

$$P(y|x) = N(y | Ax + b, L)$$

$$x \in \mathbb{R}^D, \quad y \in \mathbb{R}^M \quad A \in \mathbb{R}^{M \times D}$$

$$\text{let } z := \begin{bmatrix} x \\ y \end{bmatrix} \in \mathbb{R}^{D+M}$$

$$\Rightarrow P(z) = P(x) P(y|x)$$

$$\ln P(z) = \ln P(x) + \ln P(y|x)$$

$$= -\frac{1}{2}(x-a)^T \Lambda (x-a) - \frac{1}{2}(y - Ax - b)^T L (y - Ax - b) + \text{const}$$

For the second-order term:

$$\begin{aligned} & -\frac{1}{2}x^T (\Lambda + A^T \Lambda A) x - \frac{1}{2}y^T L y + \frac{1}{2}y^T L A x + \frac{1}{2}x^T A^T L y \\ &= -\frac{1}{2} \begin{bmatrix} x \\ y \end{bmatrix}^T \begin{bmatrix} \Lambda + A^T \Lambda A & -A^T L \\ -LA & L \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} \end{aligned}$$

$$\text{let } R = \begin{bmatrix} \Lambda + A^T \Lambda A & -A^T L \\ -LA & L \end{bmatrix}$$

$$\Rightarrow \text{The second-order term: } -\frac{1}{2}z^T R z$$

$$\text{By the Lemma: } \begin{pmatrix} A & B \\ C & D \end{pmatrix}^{-1} = \begin{pmatrix} M & -MBD^{-1} \\ -D^{-1}CA & D^{-1} + D^{-1}CA B D^{-1} \end{pmatrix}$$

$$M = (A - BD^{-1}C)^{-1}$$

$$\begin{aligned} \text{cov}(z) &= R^{-1} = \begin{bmatrix} \Lambda + A^T \Lambda A & -A^T L \\ -LA & L \end{bmatrix}^{-1} \\ &= \begin{bmatrix} \Lambda^{-1} & \Lambda^{-1} A^T \\ A \Lambda^{-1} & L^{-1} + A \Lambda^{-1} A^T \end{bmatrix} \end{aligned}$$

\Rightarrow For the linear terms:

$$x^T \Lambda u - x^T A^T L b + y^T L b = \begin{bmatrix} x \\ y \end{bmatrix}^T \begin{bmatrix} \Lambda u - A^T L b \\ L b \end{bmatrix}$$

$$\Rightarrow E[z] = R^{-1} \begin{bmatrix} \Lambda u - A^T L b \\ L b \end{bmatrix}$$

$$= \begin{bmatrix} I^{-1} & I^{-1}A^T \\ A I^{-1} & L^{-1} + A I^{-1} A^T \end{bmatrix} \begin{bmatrix} Iu - A^T L b \\ Lb \end{bmatrix}$$

$$= \begin{bmatrix} u \\ A u + b \end{bmatrix}$$

$$P(y) = \int P(z) dz$$

$$\mathbb{E}[y] = Au + b$$

$$\text{cov}(y) = \mathbb{E}_{yy} = L^{-1} + A I^{-1} A^T$$

For $x|y$:

$$\begin{aligned} \mathbb{E}[x|y] &= u_x - I_{xx}^{-1} I_{xy} (y - u_y) = \mathbb{E}[x] + \frac{\text{cov}(x,y)}{\text{var}(y)} (y - u_y) \\ &= (I + A^T L A)^{-1} \{ A^T L (y - b) + Iu \} \end{aligned}$$

$$\text{cov}(x|y) = (I + A^T L A)^{-1}$$

Part II. Bayesian Curve-fitting.

Setting: given N input values $(x_1, \dots, x_N)^T$, and the respective target values $(t_1, \dots, t_N)^T$

The uncertainty of the target values are expressed as probability distribution.

$$P\{t|x, w, \beta\} = N\{t|y(x, w), \beta^{-1}\}, w \text{ is the parameter}$$

We have a belief of the parameter w , we express this belief as the form of a prior distribution.

$$P(w|\alpha) = N(w|0, \alpha^{-1}) = \left(\frac{\alpha}{2\pi}\right)^{\frac{M}{2}} \exp\left\{-\frac{\alpha}{2} w^T w\right\}$$

Now, we are given a known data set X_{data} and their target value T_{data} along with a new point x , we want to predict the distribution of y

\Rightarrow which is

$$P\{t| x, X_{\text{data}}, T_{\text{data}}\} = \int P\{t|x, w\} P\{w| X_{\text{data}}, T_{\text{data}}\} dw$$

$$P\{t|x, w\} \sim N(t|y(x, w), \beta^{-1})$$

$P\{w| X_{\text{data}}, T_{\text{data}}\}$ is the posterior distribution over w

Part II: Bayesian Regression.

\therefore We assume the distribution of our target value is:

$$P\{t|x, w, \beta\} = \prod_{n=1}^N N(t_n|w^T \phi(x_n), \beta^{-1})$$

which is the exponential of a quadratic function of w .

\therefore The corresponding of a conjugate prior is:

$$P(w) = N(w|w_0, S_0)$$

By the conjugate property of exponentials, the posterior is also a Gaussian.

$$P(w|t) = N(w|m_N, S_N)$$

∴ By the conditional Gaussian:

$$\mathbb{E}[x|y] = \Sigma \{ A^T L(y - b) + \Lambda u \}$$

$$\text{cov}[x|y] = \Sigma = (\Lambda + A^T L A)^{-1}$$

$$S_N^{-1} = S_0^{-1} + \beta \Phi^T \Phi$$

$$m_N = S_N (S_0^{-1} m_0 + \beta \Phi^T t)$$

∴ The mode of a Gaussian is its mean.

∴ The maximum posterior estimation of the weight vector is,

$$w_{MAP} = m_N = S_N (S_0^{-1} m_0 + \beta \Phi^T t)$$

If the data points arrive sequentially, then the posterior distribution at any stage serve as the prior distribution of the next point.

We can view this mode as a point estimate of w .

Formally, a point estimate in the Bayesian model is defined as the estimator that minimizes the Bayesian Risk. (Bayesian Estimator)

$$\text{Bayesian Risk: } R(\theta, \delta) = \mathbb{E}_{\theta} [L(\theta, \delta(x))] = \int_X L(\theta, \delta(x)) f(x|\theta) dx$$

$$B_{\pi}(\theta) = \int R(\theta, \hat{\theta}) \pi(\theta) d\theta = \mathbb{E} [\mathbb{E} [L(\theta, \hat{\theta}) | \theta]]$$

$$\hat{\theta}^* = \arg \min_{\hat{\theta}} B_{\pi}(\hat{\theta})$$

By the tower property:

$$B_\theta(\theta) = \mathbb{E}[L(\theta, \delta(x))] = \mathbb{E}[\mathbb{E}[L(\theta, \delta(x))|\theta]]$$

Theorem:

If: 1° there exists δ_0 an estimator of θ with finite MSE
for all θ

2° there exists a $\delta'(x)$ that minimizes

$$\mathbb{E}[L(\theta, \delta(x))|x=x] \text{ for almost every } x$$

then $\delta'(x)$ is a Bayes estimator for θ

Proof:

$$\mathbb{E}[L(\theta, \delta'(x))|x=x] \leq \mathbb{E}[L(\theta, \delta(x))|x=x] \text{ for almost every } x$$

By taking expectation over x :

$$\mathbb{E}[L(\theta, \delta'(x))] \leq \mathbb{E}[L(\theta, \delta(x))]$$

$$\therefore \hat{\theta}^* = \underset{\hat{\theta}}{\operatorname{arg\,min}} B_{\pi}(\hat{\theta}) = \underset{\hat{\theta}}{\operatorname{arg\,min}} \mathbb{E}[L(\theta, \hat{\theta})|x=x]$$

If we take the loss function as MSE:

we want to minimize the posterior MSE:

$$\min \mathbb{E}[\pi(\theta)(\hat{\theta} - \theta)^2|x=x]$$

The expectation has nothing to do with $\hat{\theta}$

Take $\hat{\theta}$ out of the expectation:

$$\hat{\theta}^2 \mathbb{E}[\pi(\theta) | x=x] - 2\hat{\theta} \mathbb{E}[\pi(\theta)\theta | x=x] + \mathbb{E}[\pi(\theta)\theta^2 | x=x]$$

A convex quadratic function of $\hat{\theta}$

$$\hat{\theta}^* = \frac{\mathbb{E}[\pi(\theta)\theta | x=x]}{\mathbb{E}[\pi(\theta) | x=x]}$$

$$\Rightarrow \hat{\theta}^* = \hat{\mathbb{E}}_{\theta}(\theta | \text{data}) = \frac{\int_{\theta} \theta P(x=x|\theta) \pi(\theta) d\theta}{\int_{\theta} P(x=x|\theta) \pi(\theta) d\theta}$$

which is just the posterior mean.

The predictive distribution:

\tilde{x}, \tilde{t} are from the training set, we want

to predict the label of the testing set x

$$\Rightarrow P(+ | x, \tilde{x}, \tilde{t}) = \int P(+ | w, \beta) P(w | x, \beta) dw$$

$$P(+ | w, \beta) = N(t | y_i x, \beta^{-1})$$

$$P(w | \alpha, \beta) = N(w | \mu_N, \sigma_N)$$

$P(+ | x, \tilde{x}, \tilde{t})$ equals a convolution of two Gaussian distributions.

By marking the result of the conditional Gaussian.

$$p(y) = \int p(x, y) dx = \int p(y|x) p(x) dx$$

$$= N(y | Ax + b, I^{-1} + A\Lambda^{-1}A^T)$$

$$\Rightarrow p(t|x, \hat{x}, \hat{f}) = N\left(t | \hat{w}^\top \phi(x), \frac{1}{p} + \phi(x)^\top S_N \phi(x)\right)$$

Part IV: Ridge Regression and Bayesian Regression.

The Bayesian regression assumes the parameter w and σ^2 to be the random variables, while take x and y to be fixed.

for $w \in \mathbb{R}^D$:

$$w | \sigma^2 \sim N(\hat{w}_D, \sigma^2 \Lambda^{-1} I_{DD})$$

$$\sigma^2 \sim IG(\alpha_0, \beta_0)$$

IG is the inverse Gamma distribution. why choose IG?

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

Due to the computation advantage of the conjugate priors.

$$\Gamma(\alpha) = \int_0^\infty x^{\alpha-1} e^{-x} dx$$

\Rightarrow The joint probability of w and σ^2 are:

$$P\{w, \sigma^2 | \tilde{t}, \tilde{x}\} \propto P\{\tilde{t} | \tilde{x}, w, \sigma^2\} \cdot P\{w | \sigma^2\} \cdot P\{\sigma^2\}$$

Due to the Gaussian assumption of the delta distribution.

$$\begin{aligned} & \propto \sigma^{-n} \exp \left[-\frac{1}{2} \sigma^2 (\tilde{t} - \tilde{x}w)^T (\tilde{t} - \tilde{x}w) \right] \\ & \cdot \sigma^{-D} \exp \left[-\frac{1}{2} \sigma^2 \lambda w^T w \right] \cdot (\sigma^2)^{-\alpha_0-1} \exp \left[-\frac{1}{2} \sigma^2 \beta_0 \right] \\ & = \sigma^{-n+D-\frac{1}{2}\alpha_0-2} \exp \left\{ -\frac{1}{2} \sigma^2 \left[(\tilde{t} - \tilde{x}w)^T (\tilde{t} - \tilde{x}w) + \lambda w^T w + \beta_0 \right] \right\}. \end{aligned}$$

¶

We can consider the above as a multivariate Gaussian distribution with respective to w , with terms involving

$$\begin{aligned} & \Rightarrow (\tilde{t} - \tilde{x}w)^T (\tilde{t} - \tilde{x}w) + \lambda w^T w \\ & = \tilde{t}^T \tilde{t} - w^T \tilde{x}^T \tilde{t} - \tilde{t}^T \tilde{x} w + w^T \tilde{x}^T \tilde{x} w + \lambda w^T w \\ & = \tilde{t}^T \tilde{t} - w^T (\tilde{x}^T \tilde{x} + \lambda I_{DD}) (\tilde{x}^T \tilde{x} + \lambda I_{DD})^{-1} \cdot \tilde{x}^T \tilde{t} \\ & \quad - \tilde{t}^T \tilde{x} (\tilde{x}^T \tilde{x} + \lambda I_{DD})^{-1} (\tilde{x}^T \tilde{x} + \lambda I_{DD}) \cdot w + w^T (\tilde{x}^T \tilde{x} + \lambda I_{DD}) w. \end{aligned}$$

Recall that in Ridge regression, the analytical solution for w is:

$$w(\lambda) = (\tilde{x}^T \tilde{x} + \lambda I_{DD})^{-1} \cdot \tilde{x} \cdot \tilde{t}$$

\therefore The above is equivalent to:

$$\begin{aligned} & \tilde{t}^T \tilde{t} - w^T (\tilde{x}^T \tilde{x} + \lambda I_{DD}) \cdot w(\lambda) - w(\lambda)^T (\tilde{x}^T \tilde{x} + \lambda I_{DD}) w + w^T (\tilde{x}^T \tilde{x} + \lambda I_{DD}) w \\ & = \tilde{t}^T \tilde{t} - \tilde{t}^T \tilde{x} (\tilde{x}^T \tilde{x} + \lambda I_{DD})^{-1} \tilde{x}^T \tilde{t} + [w - w(\lambda)]^T (\tilde{x}^T \tilde{x} + \lambda I_{DD}) [w - w(\lambda)] \end{aligned}$$

$$\Rightarrow P\{w, \epsilon^2 | \bar{t}, \bar{x}\} \propto g_w(w | \bar{t}, \bar{x}) \cdot g(\epsilon^2 | \bar{t}, \bar{x})$$

with $g_w(w | \bar{t}, \bar{x}) \propto \exp\left\{-\frac{1}{2}\epsilon^2[\bar{w} - w_{(0)}]^T [\bar{x}^T \bar{x} + \lambda I_{DD}]^{-1} [\bar{w} - w_{(0)}]\right\}$

$$g(\epsilon^2 | \bar{t}, \bar{x}) \propto (\epsilon^2)^{-k_0-1} \exp\left\{-\frac{1}{2}\epsilon^2 \beta_0\right\}$$

\Rightarrow The conditional prior mean of w is $E[w | \epsilon^2, \bar{x}, \bar{t}] = w_{(0)}$

$$\Rightarrow w_{MAP} = w_{(0)} = (\bar{x}^T \bar{x} + \lambda I_{DD})^{-1} \bar{x}^T \bar{t}$$

The penalty parameter λ can be seen as the belief of the precision of the prior distribution. (the scaled parameter of the distribution)

If λ is large, we are confident about the prior. Actually, any penalized estimator has a Bayesian equivalent.

The Bayesian estimator is:

$$\hat{w}_{MAP} = \underset{w}{\operatorname{argmax}} \pi\{w | X=x\} = \underset{w}{\operatorname{argmax}} \frac{P\{X=x | w\} \pi(w)}{\int P\{X=x | w\} \pi(w) dw}$$

$$\Rightarrow \hat{w}_{MAP} = \underset{w}{\operatorname{argmax}} \log P\{X=x | w\} + \log [\pi(w)]$$

\uparrow log likelihood

- $\approx \log(\pi(w)) \propto \sigma^{-2} \|w\|^2$ If we use normal prior
- If we use other prior distribution, we can model other penalty terms.

Part V: Conjugate priors.

Recall the general process of the Bayesian curve-fitting:

prior distribution: $P(\theta)$

$$P(\theta|x) = \frac{P(x|\theta) \cdot P(\theta)}{P(x)}$$

$$P(x) = \int P(x|\theta) P(\theta) d\theta$$

$$P(x_{new}|x) = \int P(x_{new}|\theta) P(\theta|x) d\theta$$

\therefore Compute the integral every time is computational complex

\therefore We want to

1° $P(x) = \int P(x|\theta) P(\theta) d\theta$ can be obtained tractably.

2° $P(x_{new}|x) = \int P(x_{new}|\theta) P(\theta|x) d\theta$ can be obtained tractably.

Take univariate Gaussian as an example:

$$P(x|u, \sigma^2) \propto (\sigma^2)^{-\frac{1}{2}} \exp\left\{-\frac{1}{2}\sigma^{-2}(x-u)^2\right\}.$$

10 Conjugacy for the mean.

Assume the conjugate prior of the mean is Gaussian:

$$p\{u|u_0, \sigma_0^2\} = (\sigma_0^2)^{-\frac{1}{2}} \exp\left\{-\frac{1}{2}\sigma_0^{-2}(u-u_0)^2\right\}$$

To obtain the posterior $p\{u|x\}$

If only one data point is observed:

$$\mathbb{E}[u|x] = \mathbb{E}[u] + \frac{\text{Cov}(x,u)}{\text{Var}(x)} (x - \mathbb{E}[x])$$

$$\text{Var}(u|x) = \text{Var}(u) - \frac{\text{Cov}^2(x,u)}{\text{Var}(x)}$$

Write: $x = u + \epsilon \varepsilon \quad \varepsilon \sim N(0,1)$

$$u = u_0 + \sigma_0 \cdot \varepsilon \quad \varepsilon \sim N(0,1)$$

$$\mathbb{E}[x] = \mathbb{E}[u] + \mathbb{E}[\varepsilon] = u_0$$

$$\text{Var}(x) = \mathbb{E}[(x-u_0)^2] = \mathbb{E}[u-u_0+\epsilon\varepsilon]^2 = \sigma^2 + \sigma_0^2$$

$$\text{Cov}(x,u) = \mathbb{E}[(x-u_0)(u-u_0)] = \mathbb{E}[(u-u_0+\epsilon\varepsilon)(u-u_0)] = \sigma_0^2$$

$$\therefore u_{\text{post}} := \mathbb{E}[u|x=x] = u_0 + \frac{\sigma_0^2}{\sigma^2 + \sigma_0^2} (x - u_0) = \frac{\sigma_0^2}{\sigma^2 + \sigma_0^2} x + \frac{\sigma^2}{\sigma^2 + \sigma_0^2} u_0$$

$$\sigma_{\text{post}}^2 := \text{Var}(u|x=x) = \sigma_0^2 - \frac{\sigma_0^4}{\sigma^2 + \sigma_0^2} = \frac{\sigma^2 \sigma_0^2}{\sigma^2 + \sigma_0^2}$$

Expressed in precision:

$$\tau = \frac{1}{\sigma^2} \quad ; \quad \tau_0 = \frac{1}{\sigma_0^2}$$

$$\mathbb{E}[u|x=x] = \frac{\tau}{\tau + \tau_0} x + \frac{\tau_0}{\tau + \tau_0} u_0$$

$$T_{\text{post}} = T + T_0.$$

For the multiple data input:

$$\mathbf{x} = (x_1, \dots, x_n), x_i \sim N(\mu, \sigma^2)$$

Expressed in precision:

$$P(x|u, T) \propto T^{\frac{1}{2}} e^{-\frac{T}{2} \sum_{i=1}^n (x_i - u)^2}$$

$$\sum_{i=1}^n (x_i - u)^2 = \sum_{i=1}^n (x_i - \bar{x} + \bar{x} - u)^2$$

$$= \sum_{i=1}^n (x_i - \bar{x})^2 + n(\bar{x} - u)^2$$

$$\sum_{i=1}^n (x_i - \bar{x})^2 \text{ is constant}$$

\bar{x} is a sufficient statistics with respect to u .

∴ The problem is reduced to univariate case.

$$\bar{x} \sim N(\mu, \frac{\sigma^2}{n})$$

$$u_{\text{post}} = \frac{nT}{nT+T_0} \bar{x} + \frac{T_0}{nT+T_0} u_0$$

$$T_{\text{post}} = nT + T_0.$$

$$\begin{aligned} P\{x_{\text{new}} | x, T\} &= \int P\{x_{\text{new}} | x, u, T\} P\{u | x, T\} du \\ &= \int P\{x_{\text{new}} | u, T\} P\{u | x, T\} du \end{aligned}$$

This is again the convolution of two Gaussian

∴ $x_{\text{new}} | x, T$ is Gaussian distributed with new post ,

$$\text{Variance } \sigma_{\text{post}}^2 + \sigma^2.$$

20 Conjugacy for variance

$$P\{x|u, \sigma^2\} \sim (\sigma^2)^{-\frac{1}{2}} \exp\left\{-\frac{1}{2}\sigma^{-2}(x-u)^2\right\}$$

$$= (\sigma^2)^\alpha e^{-b/\sigma^2}$$

Assume the conjugate prior distribution of σ^2 is
inverse Gamma.

$$P\{\sigma^2|\alpha, \beta\} = \frac{\beta^\alpha}{\Gamma(\alpha)} (\sigma^2)^{-\alpha-1} e^{-\frac{\beta}{\sigma^2}}$$

The posterior:

$$P\{\sigma^2|x, u, \alpha, \beta\} \propto P\{x|u, \sigma^2\} P\{\sigma^2|\alpha, \beta\}$$

$$\propto (\sigma^2)^{-\frac{n}{2}} \exp\left\{-\frac{1}{2}\sum_{i=1}^n \frac{(x_i-u)^2}{\sigma^2}\right\} (\sigma^2)^{-\alpha-1} e^{-\beta/\sigma^2}$$

$$= (\sigma^2)^{-\left(\alpha + \frac{n}{2}\right)-1} \exp\left\{-\left(\beta + \frac{1}{2}\sum_{i=1}^n (x_i-u)^2/\sigma^2\right)\right\}$$

$$\sim I_G(\alpha + \frac{n}{2}, \beta + \frac{1}{2}\sum_{i=1}^n (x_i-u)^2)$$

Expressed in terms of the precision:

$$\therefore \sigma^2 \sim I_G(\alpha, \beta)$$

$$\therefore \tau = \frac{1}{\sigma^2} \sim G_\alpha(\alpha, \beta)$$

$$P\{\tau|x, u, \alpha, \beta\} \propto P\{x|u, \tau\} P\{\tau|\alpha, \beta\}$$

$$\propto \tau^{\frac{n}{2}} \exp\left\{-\frac{1}{2}\sum_{i=1}^n (x_i-u)^2\right\} \tau^{\alpha-1} \cdot \exp\{-\beta\tau\}$$

$$= T^{\alpha + \frac{1}{2} - 1} \exp \left\{ -(\beta + \frac{T}{2} \sum_{i=1}^n (x_i - \mu)^2) \right\}$$

$$\Rightarrow T|x \sim \text{Ga}(x + \frac{n}{2}, \beta + \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^2)$$

$$P\{X_{new} | x, \mu, \alpha, \beta\} = \int P\{X_{new} | x, \mu, T\} P\{T | x, \mu, \alpha, \beta\} dT$$

$$= \int P\{X_{new} | \mu, T\} P\{T | x, \mu, \alpha, \beta\} dT$$

$$= \int \frac{\beta^\alpha}{\Gamma(\alpha)} T^{\alpha-1} e^{-\beta T} \left(\frac{T}{2\pi}\right)^{\frac{1}{2}} e^{-\frac{T}{2}(x-\mu)^2} dT$$

$$= \frac{\beta^\alpha}{\Gamma(\alpha)} \frac{1}{(2\pi)^{\frac{1}{2}}} \int T^{\alpha-\frac{1}{2}} e^{-\beta T} e^{-\frac{1}{2}(x-\mu)^2} dT$$

↑ An unnormalized Gamma integral

$$= \frac{\beta^\alpha}{\Gamma(\alpha)} \frac{1}{(2\pi)^{\frac{1}{2}}} \int T^{\alpha-\frac{1}{2}} e^{-(\beta + \frac{1}{2}(x-\mu)^2)T} dT$$

$$= \frac{\beta^\alpha}{\Gamma(\alpha)} \frac{1}{(2\pi)^{\frac{1}{2}}} \frac{\Gamma(\alpha + \frac{1}{2})}{(\beta + \frac{1}{2}(x-\mu)^2)^{\alpha + \frac{1}{2}}}$$

$$= \frac{\Gamma(\alpha + \frac{1}{2})}{\Gamma(\alpha)} \frac{1}{(2\pi\beta)^{\frac{1}{2}}} \frac{1}{(1 + \frac{1}{2\beta}(x-\mu)^2)^{\alpha + \frac{1}{2}}}$$

↑

Generalized student's t.

$$\text{let } p := \frac{\alpha}{2}, \quad \lambda := \frac{\alpha}{\beta}.$$

$$P\{X | \mu, \lambda, p\} = \frac{\Gamma(\frac{\lambda+1}{2})}{\Gamma(\frac{\lambda}{2})} \left(\frac{\lambda}{\pi p}\right)^{\frac{1}{2}} \frac{1}{(1 + \frac{\lambda}{p}(x-\mu)^2)^{\frac{\lambda+1}{2}}}$$

Part VI: Markov Chain Monte Carlo.

When we don't use conjugate prior, we can't yield a tractable closed form formula of the posterior.

△ How to use Markov Chain:

1° Construct a Markov chain whose stationary probability equals to the desired probability

2° Iterate until the chain is converged and reaches stationary.

3° Draw sample from this stationary chain.

△ Monte-Carlo Integration.

Assume the analytical solution of $\hat{E}[h(y)] = \int h(y)\pi(y)dy$ is impossible

Draw enough samples from $\pi(\cdot)$: $y^{(1)}, \dots, y^{(n)}$

⇒ The estimated expectation: $\hat{E}[h(y)] = n^{-1} \sum_{i=1}^n h(y_i)$

By the Law of large numbers, it's a consistent estimator.

For the Bayesian framework:

$$P(w|\bar{x}) = \int_0^\infty P(w, \sigma^2 | \bar{x}) d\sigma^2 = \int_0^\infty P(w|\sigma^2, \bar{x}) P(\sigma^2 | \bar{x}) d\sigma^2$$

∴ The prior distribution of σ^2 is known.

- 1° Draw $\epsilon^{(1)}, \epsilon^{(2)}, \dots, \epsilon^{(T)}$ from $P\{\epsilon^t | \bar{t}, \bar{x}\}$
- 2° Draw $\{w^{(t)}\}_{t=1}^T$ from $P\{w | \epsilon^{(t)}, \bar{x}, \bar{t}\}$, $t = 1, \dots, T$.
- 3°. From the above procedure, ϵ^t has been integrated out.
 i.e. the resulting $p^{(1)}, \dots, p^{(T)}$ is from $P\{p | \bar{t}, \bar{x}\}$

How to sample from $\pi(y)$ is the key in this process.

Metropolis-Hastings sampler as an example.

Task 1: Construct a Markov chain to model $\pi(\cdot)$.

For a random chain $\{Y_t\}$, $Y_t \in \mathbb{R}^D$

The stationary distribution satisfies:

$$\varphi(y_a) = \int_{\mathbb{R}^D} f_{Y_{t+1}|Y_t}(y_{t+1} | y_a) \varphi(y_b) dy_b$$

The stationary distribution must satisfy the balance equation.

$$f_{Y_{t+1}|Y_t}(y_b | y_a) \varphi(y_a) = f_{Y_t|Y_{t+1}}(y_a | y_b) \varphi(y_b)$$

For MCMC, the stationary distribution is the target density where the sample is generated.

However, the transition kernel is unknown.

To have the chain to work, we need to compute its transition probability.

For an arbitrary kernel, it's unlikely to be the valid one.

$$\int_{Y_{t+1}|Y_t} (y_b | y_a) \varphi(y_a) = \int_{Y_t|Y_{t+1}} (y_a | y_b) \varphi(y_b)$$

We need to correct this arbitrary kernel.

Given the current state $y_t = y_t$, this kernel produce a suggestion y_{t+1} for the next state of the Markov chain. Check this by the balance equation. If the value is too far from the stationary distribution, reject it and repeat until we find an acceptable one.

The probability of acceptance given the current state $y_t = y_t$: $A(y_{t+1}|y_t)$

⇒ The constructed transition probability is equivalent to:

$$h_{t+1|t} (y_a | y_b) = g_{t+1|t} (y_a | y_b) A(y_a | y_b) + r(y_b) \delta_{y_b}(y_a)$$

$A(y_a | y_b)$ is the acceptance probability of the suggestion y_a :

$$A(y_a | y_b) = \min \left\{ 1, \frac{\ell(y_a)}{\ell(y_b)} \frac{g_{t+1|t} (y_b | y_a)}{g_{t+1|t} (y_a | y_b)} \right\}$$

$\delta_{y_b}(y_a)$ is a Dirac function.

Dirac function is a generalized function that can be used in describing random variables with mixture of a continuous part and a discrete part. It is characterized by returning the value 0 anywhere besides its mass.

In our case, the discrete part is y_b since we want to model

the case that the chain remains the same state at $t+1$

$$\pi_{ya} = 1 - \int_{\mathbb{R}^D} p_{t+1|t} (y_a | y_b) A(y_a | y_b) dy_a$$

The probability to reject y_a .

We want to verify if the above defined transition kernel is valid for the Markov Chain:

$$h_{t+1|t} (y_a | y_b) Q(y_b) = h_{t+1|t} (y_b | y_a) Q(y_a)$$

$$\Rightarrow \int_{\mathbb{R}^D} p_{t+1|t} (y_a | y_b) A(y_a | y_b) Q(y_b) + \boxed{\pi_{ya} \delta_{y_b}(y_a) Q(y_b)}$$

"exists only when $y_a = y_b$ "

$$= \int_{\mathbb{R}^D} p_{t+1|t} (y_b | y_a) A(y_b | y_a) Q(y_a) + \boxed{\pi_{yb} \delta_{ya}(y_b) Q(y_a)}$$

$$\Rightarrow p_{t+1|t} (y_a | y_b) A(y_a | y_b) Q(y_b) = \int_{\mathbb{R}^D} p_{t+1|t} (y_b | y_a) A(y_b | y_a) Q(y_a)$$

To verify this equation.

If $A(y_a | y_b) < 1$:

$$\Rightarrow \int_{\mathbb{R}^D} p_{t+1|t} (y_a | y_b) \frac{Q(y_a)}{\int_{\mathbb{R}^D} p_{t+1|t} (y_b | y_a) Q(y_b)} \cdot Q(y_b)$$

$$= Q(y_a) \int_{\mathbb{R}^D} p_{t+1|t} (y_b | y_a) Q(y_b)$$

$$\because A(y_a | y_b) < 1 \quad \therefore A(y_b | y_a) = 1$$

The equality holds.

π_{t+1} is the valid stationary distribution can be verified by:

$$\begin{aligned} \int_{\mathbb{R}^D} h_{t+1|t} f_t(y_b|y_a) \varrho(y_b) dy_b &= \int_{\mathbb{R}^D} h_{t+1|t} f_t(y_b|y_a) \varrho(y_a) dy_b \\ &= \int_{\mathbb{R}^D} g_{t+1|t} f_t(y_b|y_a) A(y_b|y_a) \varrho(y_a) dy_b + \int_{\mathbb{R}^D} r(y_b) \delta_{y_a(y_b)} \varrho(y_a) dy_b \\ &= \varrho(y_a) \int_{\mathbb{R}^D} g_{t+1|t} f_t(y_b|y_a) A(y_b|y_a) dy_b + \varrho(y_a) \int_{\mathbb{R}^D} r(y_b) \delta_{y_a(y_b)} dy_b \\ &= \varrho(y_a) [1 - r(y_a)] + r(y_a) \varrho(y_a) = \varrho(y_a) \end{aligned}$$

Part VII Empirical Bayes

Does not fully specify the priors, but only the form of the priors, leaving the hyperparameters to be estimated.

Assume: w and σ^2 are endowed with the conjugate priors.

$$w|\sigma^2 \sim N(\theta_0, \sigma^2 \lambda^{-1} I_D)$$

$$\sigma^2 \sim IG(\alpha_0, \beta_0)$$

Thus three parameters to be estimated.

$$\lambda, \alpha_0, \beta_0$$

The basic idea is to integrate out the model parameter and results only a marginal posterior with only data and hyperparameters.

$$\begin{aligned}
 \hat{\lambda}_{\text{eb}} &= \underset{\lambda}{\arg \max} \int_0^{\infty} \int_{\mathbb{R}^D} p(w, \sigma^2 | \tilde{t}, \tilde{x}) dw d\sigma^2 \\
 &= \underset{\lambda}{\arg \max} \int_0^{\infty} \int_{\mathbb{R}^D} \sigma^{-n} \exp \left\{ -\frac{1}{2} \sigma^{-2} (\tilde{t} - \tilde{x}w)^T (\tilde{t} - \tilde{x}w) \right\} \\
 &\quad \cdot \sigma^{-D} \exp \left\{ -\frac{1}{2} \sigma^2 \lambda w^T w \right\} \cdot (\sigma^2)^{-\alpha_0 - 1} \exp \left\{ -\beta_0 \sigma^{-2} \right\} dw d\sigma^2 \\
 &= \underset{\lambda}{\arg \max} \int_0^{\infty} \int_{\mathbb{R}^D} \underbrace{\sigma^{-n} \exp \left\{ -\frac{1}{2} \sigma^{-2} \left[\tilde{t}^T \tilde{t} - \tilde{t}^T \tilde{x} (\tilde{x}^T \tilde{x} + \lambda I_{DD})^{-1} \tilde{x}^T \tilde{t} \right] \right\}}_{\text{"unnormalized Gaussian"}}
 \end{aligned}$$

invacant to β

$$\begin{aligned}
 &\cdot \underbrace{\sigma^{-D} \exp \left\{ -\frac{1}{2} \sigma^{-2} [w - w_{\text{m}}]^T (\tilde{x}^T \tilde{x} + \lambda I_{DD})^{-1} [w - w_{\text{m}}] \right\}}_{\text{"unnormalized } I_{\text{G}}\text{"}} \\
 &\cdot \underbrace{(\sigma^2)^{-\alpha_0 - 1} \exp \left\{ -\beta_0 \sigma^{-2} \right\}}_{\text{"unnormalized } I_{\beta}\text{"}} d\beta d\sigma^2 \\
 &= \underset{\lambda}{\arg \max} \int_0^{\infty} \sigma^{-n} \exp \left\{ -\frac{1}{2} \sigma^{-2} \left[\tilde{t}^T \tilde{t} - \tilde{t}^T \tilde{x} (\tilde{x}^T \tilde{x} + \lambda I_{DD})^{-1} \tilde{x}^T \tilde{t} \right] \right\} \\
 &\quad \cdot |\tilde{x}^T \tilde{x} + \lambda I_{DD}|^{-\frac{1}{2}} (\sigma^2)^{-\alpha_0 - 1} \exp \left\{ -\beta_0 \sigma^{-2} \right\} d\sigma^2 \\
 &= \underset{\lambda}{\arg \max} \left| \tilde{x}^T \tilde{x} + \lambda I_{DD} \right|^{-\frac{1}{2}} \int_0^{\infty} \exp \left\{ -\sigma^{-2} \left[\beta_0 + \frac{1}{2} \left[\tilde{t}^T \tilde{t} - \tilde{t}^T \tilde{x} (\tilde{x}^T \tilde{x} + \lambda I_{DD})^{-1} \tilde{x}^T \tilde{t} \right] \right] \right\}^2 \\
 &\quad \cdot (\sigma^2)^{-\alpha_0 - \frac{1}{2} - 1} d\sigma^2 \quad \text{"unnormalized } I_{\beta}\text{"}
 \end{aligned}$$

$$= \underset{\lambda}{\operatorname{arg\,max}} \quad (\bar{x}^T \bar{x} + \lambda I_D)^{-\frac{1}{2}} b_1^{-\alpha_0 - \frac{n}{2}}$$

$$\text{where } b_1 = b_0 + \frac{1}{2} \left[\bar{f}^T \bar{f} - \bar{f}^T \bar{x} (\bar{x}^T \bar{x} + \lambda I_D)^{-1} \bar{x}^T \bar{f} \right]$$