

Probabilistic derivation of ω (point estimate

$$p(x) = \frac{1}{(2\pi)^{m/2} |\Sigma|^{1/2}} e^{-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)}$$

of ω using probability)

x has m dimensions

$$\mu = [\mu_1, \mu_2]$$

$$\Sigma = \begin{bmatrix} \text{Cov}(x_1, x_1) & \text{Cov}(x_1, x_2) \\ \text{Cov}(x_2, x_1) & \text{Cov}(x_2, x_2) \end{bmatrix}$$

$$\text{Cov}(z_i, z_j) = E[(z_i - E(z_i))(z_j - E(z_j))]$$

$$N(z|\mu, \Sigma) = \frac{1}{(2\pi)^{m/2} |\Sigma|^{1/2}} e^{-\frac{1}{2} (z-\mu)^T \Sigma^{-1} (z-\mu)}$$

no. of dimensions



$$N(z|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(z-\mu)^2}$$

$$\text{Cov}(z_i, z_i) = E[(z_i - E(z_i))(z_i - E(z_i))] = E[(z_i - E(z_i))^2] = \text{Var}(z_i)$$

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

$$D = (x_1, t_1) (x_2, t_2) \dots (x_n, t_n)$$

$$y(x, \omega) = \omega_0 + \omega_1 x + \omega_2 x^2 + \dots + \omega_D x^D$$

$$p(\omega|D) = \frac{p(D|\omega) p(\omega)}{\int p(D|\omega) p(\omega) d\omega}$$

$$\omega = (\omega_0, \omega_1, \omega_2, \dots, \omega_m)$$

$$p(\omega|D) \propto \left(\frac{\beta^{N/2}}{(2\pi)^{N/2}} e^{-\frac{\beta}{2} (t_n - y(x_n, \omega))^2} \right) p(\omega)$$

$$\omega = (\omega_0, \omega_1, \omega_2, \dots, \omega_m)$$

$$\begin{bmatrix} \omega_0 \\ \omega_1 \\ \omega_2 \\ \vdots \\ \omega_m \end{bmatrix} \sim \Sigma \begin{bmatrix} \circ & \circ & \circ & \circ & \circ \\ 0 & \circ & \circ & \circ & \circ \\ 0 & 0 & \circ & \circ & \circ \\ 0 & 0 & 0 & \circ & \circ \\ 0 & 0 & 0 & 0 & \circ \end{bmatrix}$$

variance = α^{-1} (assumption)
Covariance = 0

$$\therefore \text{we assume } \omega \sim N(0, \alpha^{-1} I)$$

we get this α for 1000 regression models

$$p(\omega) = \frac{1}{(2\pi)^{\frac{(M+1)}{2}} |\Sigma|^{\frac{1}{2}}} e^{-\frac{1}{2}(\omega-0)^T (\alpha^{-1} I)^{-1} (\omega-0)}$$

$$= \frac{1}{(2\pi)^{\frac{(M+1)}{2}} (\alpha^{-1})^{\frac{(M+1)}{2}}} e^{-\frac{1}{2} \omega^T (\alpha I) \omega}$$

$$p(\omega) = \frac{\alpha^{\frac{m+1}{2}}}{(2\pi)^{\frac{m+1}{2}}} e^{-\frac{\alpha}{2} \omega^T \omega}$$

we substitute:

$$p(\omega/D) \propto \frac{\beta^{N/2}}{(2\pi)^{N/2}} \frac{\alpha^{\frac{M+1}{2}}}{(2\pi)^{\frac{M+1}{2}}} e^{-\frac{\beta}{2} \left(\sum_{n=1}^N (t_n - y(x_n, \omega))^2 \right) - \frac{\alpha}{2} \omega^T \omega}$$

$$\max p(\omega/D) \Rightarrow \max \left(-\frac{\beta}{2} \sum_{n=1}^N (t_n - y(x_n, \omega))^2 - \frac{\alpha}{2} \omega^T \omega \right)$$

$$\min \left(\frac{\beta}{2} \sum_{n=1}^N (t_n - y(x_n, \omega))^2 + \frac{\alpha}{2} \omega^T \omega \right)$$

$$\min \left(\sum_{n=1}^N (t_n - y(x_n, \omega))^2 + \lambda \omega^T \omega \right) \quad \boxed{\lambda = \frac{\alpha}{2\beta}}$$

This is ridge regression

We made a weird assumption for true ω , which is $\omega \sim N(0, \sigma^2 I)$

by using a probabilistic model, we might get a better fit

we ultimately get a point estimate of ω 's (hence t 's) anyways

but we can also get a distribution of t 's
 \Rightarrow distribution of ω

Handwritten derivation of the posterior distribution $p(\omega)$ for ridge regression:

$$\begin{aligned}
 & p(D|\omega) p(\omega) \\
 &= \int p(D|\omega) p(\omega) d\omega \\
 &= p(t_1, t_2, \dots, t_N | x_1, x_2, \dots, x_N, \omega) p(\omega) \\
 &= \left(\prod_{n=1}^N p(t_n | x_n, \omega) \right) p(\omega) \\
 &= \left(\prod_{n=1}^N \frac{\beta}{\sqrt{2\pi}} e^{-\frac{\beta}{2} (t_n - y(x_n, \omega))^2} \right) p(\omega) \\
 &= \left(\frac{\beta^{N/2}}{(2\pi)^{N/2}} e^{-\frac{\beta}{2} \sum_{n=1}^N (t_n - y(x_n, \omega))^2} \right) p(\omega) \\
 &= \left(\frac{\beta^{N/2}}{(2\pi)^{N/2}} e^{-\frac{\beta}{2} \sum_{n=1}^N (t_n - y(x_n, \omega))^2} \right) \frac{1}{(2\pi)^{D/2}} e^{-\frac{\alpha}{2} \omega^T \omega}
 \end{aligned}$$

Assumptions and definitions shown in the notes:

- $\omega = (\omega_0, \omega_1, \omega_2, \dots, \omega_M)$
- $\omega \sim N(0, \sigma^2 I)$
- $p(\omega) = \frac{1}{(2\pi)^{D/2}} e^{-\frac{1}{2\sigma^2} \omega^T \omega}$

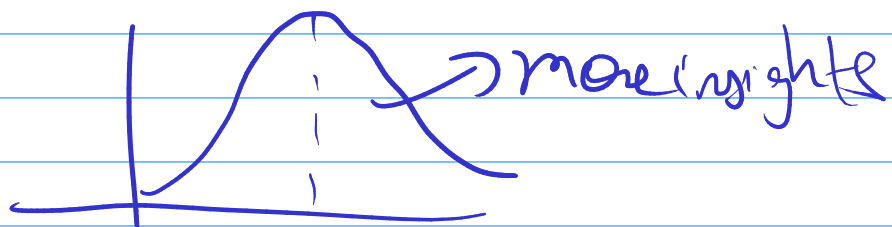
$$p(\omega/D) \propto \beta^{N/2} \left(\frac{\alpha}{2\pi} \right)^{\frac{M+1}{2}} e^{-\frac{\beta}{2} \left(\sum t_n - y(x, \omega) \right)^2 - \frac{\alpha}{2} \omega^T \omega}$$

Normalize to make α to =

(A)
$$p(\omega/D) = \frac{\left(\frac{\beta}{2\pi} \right)^{N/2} \left(\frac{\alpha}{2\pi} \right)^{\frac{M+1}{2}} e^{-\frac{\beta}{2} \left(\sum t_n - y(x, \omega) \right)^2 - \frac{\alpha}{2} \omega^T \omega}}{\int d\omega}$$

now if we take a point from here, we get a regression model!

we can even make target distribution



$$p(t/x, D) = p(t/x, (\alpha, t))$$

$$= \int \underbrace{p(t/x, \omega)}_{\text{likelihood}} p(\omega/x, t) d\omega$$

Assumption \Rightarrow likelihood
 $t \sim N$
 given x & ω

$$= \int \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{1}{2\sigma^2} (t - y(x, \omega))^2} \left[\boxed{} \right] d\omega$$

from (A) \uparrow

Closest approx:

$$p(t/x, 0) \propto N(t/m(x), s^2(x))$$

where $m(x) = \beta \phi(x)^T S \sum_{n=1}^N \phi(x_n) t_n$

$$s^2(x) = \beta^{-1} + \phi(x)^T S \phi(x)$$

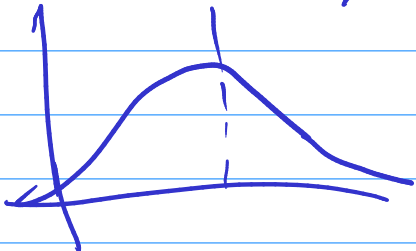
$$\Rightarrow S^{-1} = \alpha I + \beta \sum_{n=1}^N \phi(x_n) \phi(x_n)^T$$

$$\phi = [1, x, x^2, \dots, x^k]$$

Order of the polynomial regression

$$x = 1450$$

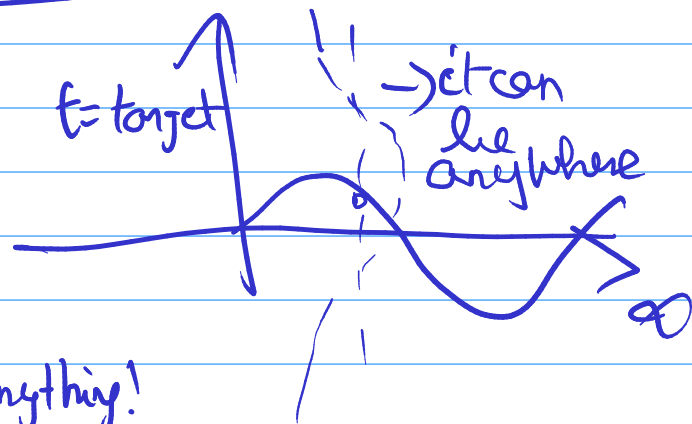
$$m(1450), s^2(1450)$$



x : univariate for now

Polynomial regression model

if we have $p(x, t)$ joint distribution



(we can get $p(t/x)$ anything!)

very hard so we find $p(t/x)$

or even just a line separating 2 classes/
best fit etc.

By knowing the joint probability distribution
how can we find $y(x)$

We are quite bad at finding $y(x)$ so we assumed

$$y(x) = w_0 + w_1 x \\ = w_0 + w_1 x + w_2 x^2$$

\vdots

$$= w_0 + w_1 x + \dots + w_9 x^9$$

And found the best of these families!

But can we find a function NOT belonging to any family??

$$p(x, t) \rightarrow y(x)$$

Now what's the loss? $L(y(x), t)$

$$\text{Wrt } y_2(x) : L(y_2(x), t)$$

(Average = Expectation)

if $E(L(y_1(x), t)) < E(L(y_2(x), t))$ then y_1 is better

what is this AVERAGE?

$$\min \iint L(y(x), t) p(x, t) dx dt = E(L)$$

$$\downarrow \text{if } L = (y(x) - t)^2$$

$$\min_{y(x)} \iint (y(x) - t)^2 p(x, t) dx dt$$

$$\frac{\partial E(L)}{\partial (y(x))} = 2 \int (y(x) - t) p(x, t) dt$$

$$\frac{\partial E(L)}{\partial (y(x))} = 0 \Rightarrow \int (y(x) - t) p(x, t) dt = 0$$

$$\Rightarrow \int y(x) p(x, t) dt = \int t p(x, t) dt$$

$$y(x) \int p(x, t) dt = \int t p(x, t) dt$$

$$y(x) p(x) = \int t (p(x, t)) dt$$

$$y(x) = \frac{1}{p(x)} \int t p(t/x) p(x) dt$$

$$y(x) = E(t/x)$$

So if we know joint distribution, we can find a function $E(t/x)$ that minimizes Loss

But finding the joint probability distribution then there's no issue!