

# SPR Day 8

## Bayesian Linear Regression

Concepts: Review Gaussian special properties

from A conditional to B joint  
from joint to conditional  
from joint to marginal  
from prior + lik to posterior

if A is Gaussian, then B is too

Posterior for Gaussian linear regression

MAP for weights

Posterior Predictive for next output

What is so special about the Gaussian?

Given

Two independent  
Gaussian r.v.s

$$x \sim N(\mu_x, \sigma_x^2)$$

$$y \sim N(\mu_y, \sigma_y^2)$$

Can Show That

Their joint is Gaussian

$$\begin{bmatrix} x \\ y \end{bmatrix} \sim N \left( \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \begin{bmatrix} \sigma_x^2 & 0 \\ 0 & \sigma_y^2 \end{bmatrix} \right)$$

Two linearly-dependent  
Gaussian r.v.s

$$x \sim N(\mu_x, \beta_x^{-1})$$

$$y \sim N(m x + b, \beta_y^{-1})$$

key:  $y$  only depends on  $x$  via mean  
which is linear func. of  $x$

Their joint is Gaussian

$$\begin{bmatrix} x \\ y \end{bmatrix} \sim N \left( \begin{bmatrix} \mu_x \\ m\mu_x + b \end{bmatrix}, \begin{bmatrix} \beta_x + \beta_y m^2 & -\beta_y m \\ -\beta_y m & \beta_y \end{bmatrix}^{-1} \right)$$

See PRML 2.103  
and 2.113

A joint distribution  
over a partitioned vector

$$\begin{bmatrix} x \\ y \end{bmatrix}^T = \begin{bmatrix} x_A \\ x_B \end{bmatrix}^T$$

$$\begin{bmatrix} x_A \\ x_B \end{bmatrix} \sim N \left( \begin{bmatrix} \mu_A \\ \mu_B \end{bmatrix}, \begin{bmatrix} \Sigma_{AA} & \Sigma_{AB} \\ \Sigma_{AB}^T & \Sigma_{BB} \end{bmatrix} \right)$$

Marginal is Gaussian  
PRML Eq. 2.98

$$P(x_A) = N(\mu_A, \Sigma_{AA})$$

Conditional is Gaussian  
PRML Eq. 2.96

$$P(x_A | x_B) = N(\mu_{A|B}, \Sigma_{AA}^{-1})$$

$$\mu_{A|B} = \mu_A - \Sigma_{AA}^{-1} \Sigma_{AB} (x_B - \mu_B)$$

Two linearly dependent  
Gaussians

$$x \sim N(\mu, \Delta^{-1})$$

$$y | x \sim N(Ax + b, L^{-1})$$

Posterior is Gaussian

$$P(x | y) = N(\mu_{x|y}, \Sigma)$$

see formula in  
PRML Eq. 2.116

# From Joint Gaussian Distribution to Marginal + Conditional

We have vector  $x = [x_1, x_2, x_3, \dots, x_{D-1}, x_D]$  in  $D$  dimensions

Consider any exclusive partition of indices into  $A$  and  $B$

e.g. if  $D=4$ , then  $A = \{1, 3\}$ ,  $B = \{2, 4\}$

or  
 $A = \{2, 3\}$ ,  $B = \{1, 4\}$

Define  $x \sim N(\mu, \Sigma)$

assuming we've reordered dims so we count first thru  $A$  then  $B$

this means:

$$\begin{pmatrix} x_A \\ x_B \end{pmatrix} \sim N \left( \begin{pmatrix} \mu_A \\ \mu_B \end{pmatrix}, \begin{pmatrix} \Sigma_{AA} & \Sigma_{AB} \\ \Sigma_{BA} & \Sigma_{BB} \end{pmatrix} \right)$$

$D \times 1$  vector

$D \times D$  cov matrix

Section 2.3.1 + 2.3.2

Bishop has lots of math deriving the following:

What is the conditional  $P(x_A | x_B)$ ?

What is the marginal  $P(x_A)$ ?

You should understand the math,  $P(x_B)$ ?  
 but we'll skip math and show results

# From Joint Gaussian to Marginal + Conditional

$$\begin{matrix} \top \\ X_A \\ | \\ X_B \\ \perp \end{matrix} \sim N \left( \begin{matrix} \top \\ \mu_A \\ | \\ \mu_B \\ \perp \end{matrix}, \begin{matrix} \Sigma_{AA} & \Sigma_{AB} \\ \Sigma_{BA} & \Sigma_{BB} \end{matrix} \right)$$

Whole covariance matrix  $\Sigma$  is  $D \times D$

has an inverse

$$\Delta = \Sigma^{-1}$$

$$= \begin{matrix} \Delta_{AA} & \Delta_{AB} \\ \Delta_{BA} & \Delta_{BB} \end{matrix}$$

Marginal

$$\begin{aligned} P(X_A) &= \int P(X_A, X_B) dx_B \\ &= N(X_A | \mu_A, \Sigma_{AA}) \end{aligned}$$

Conditional

$$\begin{aligned} P(X_A | X_B) &= \frac{P(X_A, X_B)}{P(X_B)} \\ &= N(X_A | \mu_A - \Delta_{AA}^{-1} \Delta_{AB} (X_B - \mu_B), \Delta_{AA}^{-1}) \end{aligned}$$

dim check:  $A \times 1 - (A \times A)(A \times B)(B \times 1 - B \times 1), A \times A$   
 $= A \times 1$  ✓

The Gaussian-Gaussian model for regression  
Assume  $\beta^{-1} > 0$  known precision

Prior

$$p(w) = \mathcal{N}(m_0, S_0)$$

$$m_0 \in \mathbb{R}^M$$

$S_0$  is  $M \times M$   
covariance  
matrix  
(symmetric,  
pos. def.)

Likelihood

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

iid each example  $\leftarrow$  linear function of  $w$ !

Turns out, posterior is given by (see Bishop 2.116)

$$p(w|t, x, \beta) = \mathcal{N}(m_N, S_N)$$

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

$$S_N^{-1} = S_0^{-1} + \beta \Phi^T \Phi \quad \boxed{\text{Eq. 3.50}}$$

Check: dimensions,  
what if no data?

What is the MAP estimator  
for linear regression?

$$W_{\text{MAP}} = \operatorname{argmax}_{W \in \mathbb{R}^M} \log p(W | \{x_n, t_n\}_{n=1}^N)$$

we've shown this is  
Gaussian:  $N(m_N, S_N)$

thus,

$$W_{\text{MAP}} = m_N$$

Why? Use property that  
mode of any Gaussian  
is its mean parameter

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

Consider a prior that favors zero mean + diag. covariance,

$$m_0 = \vec{0} \quad \text{and} \quad S_0 = \alpha^{-1} I_M = \begin{bmatrix} \alpha^{-1} & & \\ & \alpha^{-1} & \\ & & \ddots \\ & & & \alpha^{-1} \end{bmatrix}$$

$$S_0^{-1} = \alpha I_M$$

then our MAP simplifies to

$$W_{\text{MAP}} = (\alpha I_M + \beta \Phi^T \Phi)^{-1} \beta \Phi^T t$$

$$= \frac{\beta}{\beta} \left( \frac{\alpha}{\beta} I_M + \Phi^T \Phi \right)^{-1} \Phi^T t$$

factoring  $\frac{1}{\beta}$   
out of inverse

$$= \left( \frac{\alpha}{\beta} I_M + \Phi^T \Phi \right)^{-1} \Phi^T t$$

Looks familiar! Equivalent to our sum-of-squares penalized  
linear regression w/  $\lambda = \frac{\alpha}{\beta}$

# Predictive distribution for Bayes linear regression

Assume: Linear Regression model with prior  $m_0 = \vec{0}$ ,  $S_0 = \kappa^{-1} I$ , precision  $\beta$  known See PRML Sec. 3.3.2

Goal: Probability of next output  $t_*$  given its input features  $x_*$  and all training data  $\{x_n, t_n\}_{n=1}^N$

$$P(t_* | x_*, \{x_n, t_n\}_{n=1}^N) = \int \underbrace{p(t_* | x_*, w, \beta)}_{\substack{\uparrow \\ \text{likelihood}}} \underbrace{p(w | \{x_n, t_n\}_{n=1}^N, \alpha, \beta)}_{\substack{\uparrow \\ \text{posterior} \\ \text{given} \\ \text{train data}}} dw$$

also condition on  $\alpha > 0$   
 $\beta > 0$

this is Gaussian, with mean linear in  $w$       this is Gaussian for  $w$

Applying our Gaussian rules, looking for marginal  $p(t_*)$  given joint Gaussian  $p(w, t_*)$ , we have

↙ scalar!

$$P(t_* | x_*, \{x_n, t_n\}_{n=1}^N) = \mathcal{N}\left(t_* \mid \underbrace{m_N^T \phi(x_*)}_{\substack{\text{scalar} \\ \text{mean}}}, \underbrace{\Sigma_N^2(x_*)}_{\substack{\text{scalar} \\ \text{variance}}}\right)$$

where

$$\Sigma_N^2(x_*) = \frac{1}{\beta} + \phi(x_*)^T S_N \phi(x_*)$$

Key message: Has simple form given posterior  $p(w|x,t) = \mathcal{N}(m_N, S_N)$

PRML Eq 3.59