

# Bayesian Linear Regression

data: observe pairs  $\{x_i, y_i\}_{i=1}^N$

goals: given new  $x_*$ , predict its  $y_*$  <sup>or sample plausible  $y_*$</sup>   
 across many possible  $x$ , what is shape of  $y(x)$ ?

model:

$$y_i = y(x_i) = f(x_i) + \epsilon_i$$

deterministic mean
"noise"
 $\epsilon_i \sim N(0, \sigma_n^2)$   

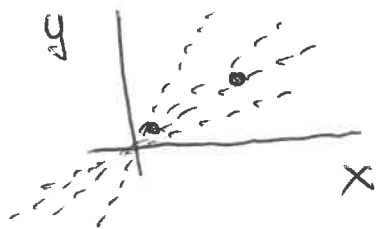

 $n$  for noise

Assume a linear model for mean function

$$f(x_i) = w * x_i \quad \text{in } 1D$$

$$= w^T x_i \quad \text{when } x_i \text{ is a } D\text{-dim vector}$$

Picture



several possible values for slope  $w$  are likely/plausible

Assume Gaussian prior on weights  $w$

$$p(w) = w \sim N(0, \sigma_p^2) \quad \text{in } 1D$$

$$\vec{w} \sim N(0, \Sigma_p) \quad \text{in } D\text{-dimensions}$$

Goal:

what is  $p(w|x, y)$ ?

Bayes thm says:

$$p(w|x, y) = \frac{p(y|x, w) p(w)}{p(y|x)} = \frac{\text{lik} * \text{prior}}{\text{marg. lik.}}$$

# 1D posterior computation

$$p(w|y, x) \propto \frac{1}{p(y|x)} \overset{\text{const}}{p(w)} \prod_{i=1}^N p(y_i|x_i, w)$$

$$\propto \mathcal{N}(w|\mu, \sigma_p^2) \prod_{i=1}^N \mathcal{N}(y_i|w \cdot x_i, \sigma_n^2)$$

$$\propto \exp\left\{-\frac{1}{2} \frac{1}{\sigma_p^2} w^2\right\} \prod_{i=1}^N \exp\left\{-\frac{1}{2} \frac{1}{\sigma_n^2} (y_i - wx_i)^2\right\}$$

$$\propto \exp\left\{-\frac{1}{2} \left( \frac{1}{\sigma_p^2} w^2 + \frac{1}{\sigma_n^2} \sum_{i=1}^N (y_i - wx_i)^2 \right)\right\}$$

$$\propto \left( \frac{1}{\sigma_p^2} w^2 + \frac{1}{\sigma_n^2} \left[ \sum_i y_i^2 - 2 \sum_i y_i wx_i + w^2 \sum_i x_i^2 \right] \right)$$

$$\propto \exp\left\{-\frac{1}{2} \left( \overset{\text{a}}{\frac{1}{\sigma_p^2} + \frac{1}{\sigma_n^2} \sum_i x_i^2} w^2 - \overset{\text{b}}{\frac{1}{\sigma_n^2} \sum_i y_i x_i} w \right)\right\}$$

$$\propto \exp\left\{-\frac{1}{2} a w^2 - b w\right\}$$

"complete the square"  
add  $-\frac{1}{2} a b^2$   
which is const wrt  $w$

$$\propto \exp\left\{-\frac{1}{2} a \left(w - \frac{b}{a}\right)^2\right\}$$

looks like Gaussian pdf!

$$\propto \exp\left\{-\frac{1}{2} \frac{1}{\sigma^2} (w - \mu)^2\right\}$$

$$\sigma^2 = \frac{1}{a} = \left( \frac{1}{\sigma_p^2} + \frac{1}{\sigma_n^2} \sum_i x_i^2 \right)^{-1} \quad \mu = \frac{b}{a} = \frac{\frac{1}{\sigma_n^2} (\sum_i y_i x_i)}{\frac{1}{\sigma_p^2} + \frac{1}{\sigma_n^2} \sum_i x_i^2}$$

Key Idea: Gaussian prior + Gaussian likelihood  $\rightarrow$  Gauss posterior  
 Assuming model is good, we now can do closed form predictions  
 Compute exact samples for this posterior

General multivariate posterior with  $D > 1$  dims

$$P(w) = N_D(w | 0, \Sigma_p)$$

$$P(y_i | x_i, w) = N_1(y_i | w^T x_i, \sigma_n^2) \quad (y_i: \text{still scalar})$$

We apply same process

$$P(w | y, X) \propto \frac{1}{P(y|X)} \overset{\text{const}}{P(w)} \prod_{i=1}^N P(y_i | x_i, w)$$

$$\propto N(w | 0, \Sigma_p) \prod_i N(y_i | w^T x_i, \sigma_n^2)$$

$$\propto \exp \left\{ -\frac{1}{2} w^T \Sigma_p^{-1} w \right\} \exp \left\{ -\frac{1}{2} \frac{1}{\sigma_n^2} (\vec{y} - Xw)^T (\vec{y} - Xw) \right\}$$

$X$ :  $D \times N$  matrix

$$\begin{bmatrix} x_1^1 & \dots & x_1^N \\ \vdots & & \vdots \\ x_D^1 & \dots & x_D^N \end{bmatrix}$$

$$\propto \exp \left\{ -\frac{1}{2} (w - \mu)^T \Sigma^{-1} (w - \mu) \right\}$$

~~Posterior~~ Posterior:  $P(w | y, X) = N(w | \mu, \Sigma)$

~~$$\Sigma = \Sigma_p + \frac{1}{\sigma_n^2} X X^T$$~~

$$A = \Sigma^{-1} = \Sigma_p^{-1} + \frac{1}{\sigma_n^2} X X^T$$

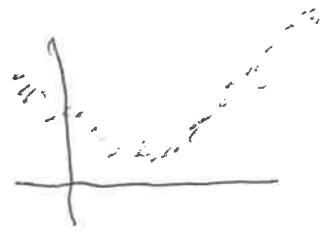
Concept check Does it match 1D case?  
 Do dims work?  $\mu = \frac{1}{\sigma_n^2} A^{-1} X \vec{y}$   
 What is posterior if no data observed?

# Beyond Linear Features Kernel Trick

In previous ML course, should have seen more flexible features



linear OK



linear will FAIL

Possible feature maps:

$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \end{bmatrix} \text{ "quadratic" } \quad \text{or} \quad \begin{bmatrix} 1 \\ x \\ \cos(x) \\ \sin(x) \end{bmatrix} \text{ "periodic"}$$

New regression model

$$y(x_i) = w^T \phi(x_i) + \epsilon_i$$

linear in FEATURE SPACE      noise

$$\phi(x_i) = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

F dim feature space

Now, can write posterior as

w also has length F

$$p(w | y, X) = N\left(w \mid \frac{1}{\sigma^2} A^{-1} \Phi^T \vec{y}, A^{-1}\right)$$

where  $A = \Sigma_p^{-1} + \frac{1}{\sigma^2} \Phi \Phi^T$   
check how big is A?

Can write posterior predictive as

$$p(f_* | x_*, X, y) = N\left(f_* \mid \frac{1}{\sigma^2} \phi(x_*)^T A^{-1} \Phi^T \vec{y}, \phi(x_*)^T A^{-1} \phi(x_*)\right)$$

What values will function take at new test points?

Consider a neat way to rewrite

A has size  $F \times F$

$$\begin{aligned} \sigma_{\text{posterior}}^2 &= \phi(x_*)^T A^{-1} \phi(x_*) \\ &= \phi(x_*)^T \left( \Sigma_p^{-1} + \frac{1}{\sigma_n^2} \Phi \Phi^T \right)^{-1} \phi(x_*) \\ &= \phi(x_*)^T \Sigma_p \phi(x_*) - \phi(x_*)^T \Sigma_p \Phi \left( K + \sigma_n^2 \mathbf{I} \right)^{-1} \Phi^T \Sigma_p \phi(x_*) \end{aligned}$$

$K = \Phi^T \Sigma_p \Phi$   
is  $N \times N$

Using MATRIX INVERSION LEMMA

When  $N < F$ , cheaper to solve using "K" formula  
 this is the kernel trick, never need to fully represent feature vectors

define function

$$k(x, x') = \phi(x)^T \Sigma_p \phi(x')$$

we'll call this a kernel function  
 aka covariance function  
 output: scalar,  $> 0$   
 larger values  $\Rightarrow$   $x$  and  $x'$  are closely correlated  
 zero  $\Rightarrow$   $x, x'$  not related

$$\sigma_{\text{posterior}}^2 = k(x_*, x_*) - k(x_*, X) \left[ k(X, X) + \sigma_n^2 \mathbf{I}_N \right]^{-1} k(X, x_*)$$

$$\mu_{\text{posterior}} = \underbrace{k(x_*, X)}_{1 \times N} \left[ \underbrace{k(X, X)}_{N \times N} + \underbrace{\sigma_n^2 \mathbf{I}_N}_{N \times N} \right]^{-1} \underbrace{y}_{N \times 1}$$

## For Next Time: Gaussian Process

Read R & W Ch. 2, esp. 2.2 <sup>which will be in-class focus</sup>  
 submit critical comments to Canvas!

GPs let us take advantage of kernel trick  
 to use very flexible feature embeddings  $\phi(x)$   
 that might be infinite in dimension

GPs use the Gaussian and its useful properties

- conjugacy
- marginalization

$$\text{if } p\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \mid \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \\ & \Sigma_{22} \end{bmatrix}\right)$$

$$\text{then } p(x_1) = \mathcal{N}(x_1 \mid \mu_1, \Sigma_{11})$$

Seen: Weights view

prior is on vector  $w$

$$w \sim p(w)$$

$$y|x \sim p(y \mid w^T \phi(x), \sigma_n^2)$$

New: Function view

prior is over possible  
function values

$$\text{specify } m(x) = \mathbb{E}[f(x)]$$

$$k(x, x') = \text{Cov}[f(x), f(x')]$$

$$\begin{array}{c|c} f_1 & x_1 \\ f_2 & \vdots \\ \vdots & \vdots \\ f_N & x_N \end{array} \sim \text{GP}()$$

$$\sim \mathcal{N}(\underline{\text{mean}}, \underline{\text{covar}})$$