

# SPPR Day 07

## Probabilistic Linear Regression

Reading: Bishop PRML Sec 3.1

and Sec. 3.3

Outline: (1) Linear Regression: A Probabilistic View

(2) ML estimator + Connections to "Least Squares"

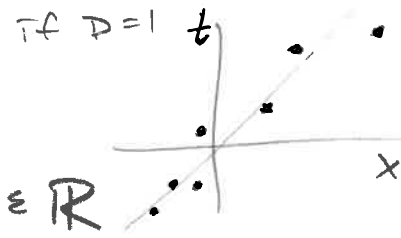
(3) Gaussian-Gaussian regression model  
to give uncertainty to weights

(4) MAP estimator + Connections to "Penalized Least Squares"  
aka Ridge Regression

# Probabilistic View of Linear Regression w/ Known Feature Transforms

## Linear Regression

## Standard Features



Goal: Given dataset  $\{x_n, t_n\}_{n=1}^N$   
with  $x_n \in \mathbb{R}^D$  and  $t_n \in \mathbb{R}$   
want to predict  $t_*$  given  $x_*$

Assume: "Linear" model, which means prediction function is linear function of input  $x_n$

$$\begin{aligned} y(x_n, w) &= w_0 + w_1 x_{n1} + w_2 x_{n2} + \dots + w_D x_{nD} \\ &= \sum_{d=0}^D w_d x_{nd} \quad (\text{defining } x_{n0} = 1 \forall n) \\ &= w^T x_n \quad \text{inner product of two } D+1 \text{ vectors} \end{aligned}$$

Often can define "smarter" features by transforming input  $x_n$  into another feature space via  $\phi(x_n)$

Define 
$$\phi(x_n) = [1 \quad \phi_1(x_n) \quad \phi_2(x_n) \quad \dots \quad \phi_{M-1}(x_n)]$$

$\phi_M(x_n)$  can be non-linear!  
sometimes called "basis" function  
M total entries, include "always 1" feature

$$x_{n1}^2 \text{ or } x_{n1} x_{n3} \text{ or } \cos(x_{n4}) \text{ or } \dots$$

Key idea is that we define a feature transform function  $\phi(x_n)$  in advance, with known size M.

"Featurized" model for prediction:

$$y(x_n, w) = \sum_{m=1}^M w_m \phi_m(x_n) = w^T \phi(x_n)$$

Note: our predictions will not be perfect.

Need to tolerate some noise.

Let's define a probabilistic approach.

Likelihood of observing output  $t_n$  given input  $x_n$

$$p(t_n | x_n, w, \beta) = N\left(t_n \mid \overbrace{w^T \phi(x_n)}^{\text{mean}}, \overbrace{\beta^{-1}}^{\text{variance}}\right)$$

notation for  
Normal PDF

1D  
variable  
in  $\mathbb{R}$

$\beta^{-1}$  is  $\text{Var}[t_n]$

$\beta$  is precision  
of  $t_n$

If we assume all  $N$  observations are i.i.d. from this distribution

likelihood

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

Taking log of both sides and simplifying

$$\log p(t | X, w, \beta) = \frac{N}{2} \log \beta - \frac{N}{2} \log(2\pi) - \beta \frac{1}{2} \underbrace{\sum_{n=1}^N (t_n - w^T \phi(x_n))^2}_{\text{sum of squared errors}}$$

log likelihood

Key idea: Can treat this as a log likelihood,  
apply ML ideas to estimate  $w_{ML}, \beta_{ML}$

# ML Estimation for weight vector + precision

View log likelihood as function of  $w, \beta$ , then try to maximize by taking gradients & setting equal to 0 & solving

step 1

step 2

step 3

Step 1

$$\alpha(w, \beta) = \frac{N}{2} \log \beta - \beta \frac{1}{2} \sum_n (t_n - w^T \phi(x_n))^2 + \text{const wrt } w, \beta$$

Step 2 Gradient wrt  $w, ML$

$$\begin{aligned} \nabla_w \alpha(w, \beta) &= \text{zero} + -\frac{1}{2} \beta \sum_n \nabla_w (t_n^2 - 2t_n w^T \phi(x_n) + w^T \phi(x_n)^2) \\ &= \text{zero} + + \beta \sum_n [t_n \phi(x_n) - \beta \frac{1}{2} \nabla_w [w^T \phi(x_n)^2]] \\ &= + \beta \sum_n [t_n \phi(x_n) - \beta \frac{1}{2} \nabla_w [w^T \phi(x_n)] \phi(x_n)] \\ &= \beta \sum_{n=1}^N (t_n - w^T \phi(x_n)) \phi(x_n)^T \end{aligned}$$

by chain rule

$\uparrow$  scalar     $\uparrow$  scalar     $\uparrow$  vector size  $M$  ✓

Step 3 set grad=0, solve for  $w, ML$

$$\vec{0} = \beta \sum_{n=1}^N (t_n - w^T \phi(x_n)) \phi(x_n)^T$$

$$\vec{0} = \sum_{n=1}^N t_n \phi(x_n)^T - w^T \sum_{n=1}^N \phi(x_n) \phi(x_n)^T$$

$\swarrow$   $M \times M$  matrix

$$w_{ML} = (\Phi^T \Phi)^{-1} \Phi^T t$$

$M \times M$      $M \times N$      $N \times 1$   
 $= M \times 1$  ✓

$$+ t \Phi^T = + w^T \Phi^T \Phi$$

$1 \times M$      $1 \times M$      $M \times M$      $\rightarrow$  transpose both sides     $\Phi^T t = \Phi^T \Phi w$   
 $M \times 1$      $M \times M$      $M \times 1$

Thus, the maximum likelihood estimator  $w_{ML}$  for parameter  $w$  is given by:

$$w_{ML} = (\Phi^T \Phi)^{-1} \Phi^T t$$

where  $t = \begin{bmatrix} t_1 \\ \vdots \\ t_N \end{bmatrix}$   $\Phi = \begin{bmatrix} 1 & \phi_1(x_1) & \dots & \phi_{M-1}(x_1) \\ \vdots & \vdots & \dots & \vdots \\ 1 & \phi_1(x_N) & \dots & \phi_{M-1}(x_N) \end{bmatrix}$

Only exists when inverse of  $\Phi^T \Phi$  exists, so that matrix must be full rank (rank  $M$ ).

Often, so long as #datapoints  $>$  # features, we'll be in good shape.  
 $N > M$

If inverse does not exist, can't estimate a unique  $w_{ML}$

What about ML estimate of  $\beta$ ? Our precision parameter?

Same process (Step 1, 2, 3) yields:

$$\beta_{ML}^{-1} = \frac{1}{\beta_{ML}} = \frac{1}{N} \sum_{n=1}^N (t_n - w_{ML}^T \phi(x_n))^2$$

Sum of squared error when we plug in ML estimate of  $w$

Looks similar to  $\sigma_{ML}^2$  for general  $\mathbb{R}^D$  Gaussians.

### Penalized ML Estimator

log likelihood

penalty term

Suppose we optimize

$$\max_w \alpha(w, \beta) + \lambda \sum_{m=1}^M w_m^2$$

sum of squares, also could write

$$\sum_m w_m^2 = w^T w$$

inner product of  $w$  with itself

with strength scalar  $\lambda > 0$  controlling strength of penalty

With this objective, we find

$$w^* = (\lambda \mathbf{I}_M + \Phi^T \Phi)^{-1} \Phi^T t$$

this is always rank  $M$  and always invertible

$M \times M$  with  $\lambda$  on diag. all zero off diagonal



# Towards a full probabilistic model for regression

Goal:

All unknown parameters (weight vector  $w \in \mathbb{R}^M$  precision  $\beta^{-1} > 0$ ) are treated probabilistically.

For now, we'll assume  $\beta^{-1} > 0$  is fixed known. Simpler.

Equivalently, we need to define a joint model

$$P(t, w | X, \beta) \\ = \underbrace{P(t | X, w, \beta)}_{\text{likelihood}} \cdot \underbrace{P(w)}_{\text{prior}}$$

Why? Given this joint, we can talk about posterior beliefs about parameters after seeing data

$$P(w | \{X_n, t_n\}_{n=1}^N, \beta)$$

- Can use the MAP estimate instead of ML (better with limited data)
- Can use samples from posterior to assess uncertainty

We can also use posterior to make predictions about new data

Use the predictive posterior.

$$P(t_* | X_*, \{X_n, t_n\}_{n=1}^N, \beta) = \int_w \underbrace{P(t_* | X_*, w, \beta)}_{\text{likelihood}} \underbrace{P(w | \{X_n, t_n\}_{n=1}^N, \beta)}_{\text{posterior}} dw$$

Key difference: Average over all  $w$  vectors, weighted by posterior density. Do NOT just commit 100% to one  $w$  vector.

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

# Posterior of Gaussian-Gaussian linear regression

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?

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

we've shown this is  
Gaussian:  $\mathcal{N}(m_N, S_N)$

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/  $\lambda = \frac{\alpha}{\beta}$