

SPR Day 07

Probabilistic Linear Regression

Reading: Bishop Sec 3.1

we'll cover

Linear Basis Function models
Maximum likelihood
"Least Squares"
closed form solution
Penalized maximum likelihood

Reading: Bishop Sec 3.3

we'll cover

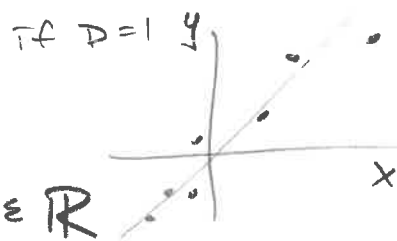
Gaussian-Gaussian model
Posterior for weights
MAP estimator

next time

predictive posterior

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

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

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

Define
$$\phi(x_n) = \begin{bmatrix} 1 & \phi_1(x_n) & \phi_2(x_n) & \dots & \phi_{M-1}(x_n) \end{bmatrix}$$

$\phi_M(x_n)$ 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:

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

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}

β^{-1} is $\text{Var}[t_n]$

β is precision
of t_n

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

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

Taking log of both sides and simplifying

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

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

View log likelihood as function of w, β , 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, β

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

by chain rule

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

Step 3 set grad=0, solve for w, β

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

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

\swarrow $M \times M$ matrix

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

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

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

$1 \times M$ $1 \times M$ $M \times M$ \rightarrow transpose both sides

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

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

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

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 β ? 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 σ_{ML}^2 for general N Gaussians.

Penalty term goal: avoid weight coefficients with large values

Penalized ML Estimator

Suppose we optimize

$$\max_w \mathcal{L}(w, \beta) - \lambda \sum_{m=1}^M w_m^2$$

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

With this objective, we find

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

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



penalty term

sum of squares, also could write

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

inner product of w with itself

this is always rank M and always invertible

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

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

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

We can also make good predictions about new data

Use the predictive posterior.

$$P(t_* | X_*, \{X_n, t_n\}_{n=1}^N) = \iint_{w, \beta} p(t_* | X_*, w, \beta) p(w | \{X_n, t_n\}, \beta) dw d\beta$$

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:

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

$D \times 1$
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 Σ is $D \times D$

has an inverse

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

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

Marginal

$$\begin{aligned} P(x_A) &= \int P(x_A, x_B) dx_B \\ &= N(x_A \mid \mu_A, \Sigma_{AA}) \end{aligned}$$

Conditional

$$\begin{aligned} P(x_A \mid x_B) &= \frac{P(x_A, x_B)}{P(x_B)} \\ &= N(x_A \mid \mu_A - \Delta_{AA}^{-1} \Delta_{AB} (x_B - \mu_B), \Delta_{AA}^{-1}) \end{aligned}$$

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