

# Linear Regression

2/19/2013

$w$  world state

$x$  observed data

$$w_i \sim \mathcal{N}(\phi_0 + \phi_1 x_i, \sigma^2)$$

define

$$\phi = \begin{bmatrix} \phi_0 \\ \phi_1 \end{bmatrix} \quad x_i = \begin{bmatrix} 1 \\ x_i \end{bmatrix}$$

$$\theta = \{\phi_0, \phi_1, \sigma\}$$

$$P(w_i | x_i, \theta) = \mathcal{N}(x_i^T \phi, \sigma^2)$$

Assume  $(w_i, x_i)$  Training examples are independent  
 Therefore  $\Delta \times N$  matrix  $N \times 1$  vector  
 $X = \begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix}$   $W = \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{pmatrix}$   $I_n = N \times N$  Identity matrix

$$P(W | X, \theta) = \mathcal{N}(X^T \phi, \sigma^2 I_n)$$

$$= \frac{1}{(2\pi)^{n/2} |\det \sigma^2 I_n|^{1/2}} \exp \left\{ -\frac{1}{2} (W - X^T \phi)^T (\sigma^2 I)^{-1} (W - X^T \phi) \right\}$$

$$= (2\pi)^{-n/2} \sigma^{-n} \exp \left\{ -\frac{1}{2\sigma^2} (W - X^T \phi)^T (W - X^T \phi) \right\}$$

$$\log P = -\frac{n}{2} \log 2\pi - n \log \sigma - \frac{1}{2\sigma^2} (W - X^T \phi)^T (W - X^T \phi)$$

$$\underbrace{W^T W - 2\phi^T X W + \phi^T X X^T \phi}_{\text{call it "M"}}$$

$$\frac{\partial \log P}{\partial \phi} = -2XW + 2XX^T \phi = 0$$

$$XX^T \phi = XW \rightarrow \phi = (XX^T)^{-1} XW$$

$$\frac{\partial \log P}{\partial \sigma} = -\frac{n}{\sigma} - \frac{1}{2} M (-2\sigma^{-3}) = 0$$

$$\frac{M}{\sigma^3} = \frac{n}{\sigma} \quad \sigma^2 = \frac{M}{n} = \frac{(W - X^T \phi)^T (W - X^T \phi)}{n}$$



# Linear Regression Cont

2/21/13

Recall  $w_i$  = world state  $x_i$  = observed data

$$P(w_i | x_i, \phi, \sigma^2) = \mathcal{N}(x_i^T \phi, \sigma^2)$$

want to learn  $\phi, \sigma^2$  from training data pairs  $(x_i, w_i)$

$$X = \begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix} \quad \text{D x N matrix} \quad W = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} \quad \text{N x 1 vector} \quad I_n = \text{N x N Identity matrix}$$

MLE estimation

$$\Rightarrow \hat{\phi} = (X X^T)^{-1} X W \quad \hat{\sigma}^2 = \frac{(W - X^T \hat{\phi})^T (W - X^T \hat{\phi})}{n}$$

Predictive distribution: given a new measurement  $x^*$ , compute the posterior distribution over values of  $w$

$$P(w | x^*, \hat{\phi}, \hat{\sigma}^2) = \mathcal{N}(x^{*T} \hat{\phi}, \hat{\sigma}^2) \\ \propto \exp \left\{ -\frac{1}{2\hat{\sigma}^2} (w - x^{*T} \hat{\phi})^2 \right\}$$

We would like to take a Bayesian approach in most cases, to "regularize" the fitting (via ~~the~~ a prior on  $\phi$ ) and thus avoid overfitting. This is especially important as data dimension  $\Delta$  gets larger



Introduce a prior on  $\phi$

$$P(\phi) = \mathcal{N}(\phi | 0, S^2 \mathbf{I}_d)$$

zero mean  
variance  $S^2$  (in general a different variance than  $\sigma^2$ )

note:  $S^2$  should be relatively large, & we are initially uncertain about values of coefficients of  $\phi$ .

note: making this prior zero-mean influences estimates of  $\phi$  to have smaller ~~est~~ magnitude. This is what has the regularizing (smoothing) effects.

By the way, the Bayesian approach to regression we are deriving, using this prior, is equivalent to "Ridge regression".

Using Bayes Rule to compute a posterior over  $\phi$

$$P(\phi | x, w) = P(w | x, \phi) P(\phi) / P(w | x) \\ \propto \mathcal{N}(w | x^T \phi, \sigma^2 \mathbf{I}_N) \mathcal{N}(\phi | 0, S^2 \mathbf{I}_d)$$

$$\propto \exp \left\{ -\frac{1}{2\sigma^2} (w - x^T \phi)^T (w - x^T \phi) - \frac{1}{2S^2} \phi^T \mathbf{I}_d \phi \right\}$$

$$= \exp \left\{ -\frac{1}{2} \left[ \frac{1}{S^2} \phi^T \mathbf{I}_d \phi + \frac{w^T w}{\sigma^2} - \frac{2\phi^T x w}{\sigma^2} + \frac{1}{\sigma^2} \phi^T x x^T \phi \right] \right\}$$

$$= \exp \left\{ -\frac{1}{2} \left[ \phi^T \left( \frac{1}{S^2} \mathbf{I}_d + \frac{1}{\sigma^2} x x^T \right) \phi - 2\phi^T \left( \frac{1}{\sigma^2} x w \right) + \frac{1}{\sigma^2} w^T w \right] \right\}$$

Dimensions  
 $w = N \times 1$   
 $x_i = D \times 1$   
 $X = D \times N$   
 $\phi = D \times 1$   
 $\mathbf{I}_N = N \times N$   
 $\mathbf{I}_d = D \times D$



## Fun Fact

any distribution proportional to  $\exp\left\{-\frac{1}{2} [ax^2 - 2bx + c]\right\}$  for a scalar variable

is a Gaussian with mean  $\frac{b}{a}$  and variance  $\frac{1}{a}$  !

Proof: "completing the square" - notes will be posted on our web site.

matrix version if  $x$  is a vector variable

$$\exp\left\{-\frac{1}{2} [x^T A x - 2x^T b + c]\right\}$$

is a Gaussian with mean  $A^{-1}b$  and variance  $A^{-1}$   
 [note  $A$  is called a "precision matrix".]

So  $P(\phi | x, w) \propto \exp\left\{-\frac{1}{2} \phi^T A \phi - 2\phi^T b + c\right\}$

$A = \frac{1}{\sigma^2} I + \frac{1}{\sigma^2} x x^T$   
 $b = \frac{1}{\sigma^2} x w$   
 $c = \frac{1}{\sigma^2} w^T w$

is Gaussian with  
 mean  $\frac{1}{\sigma^2} A^{-1} x w$  and variance  $A^{-1}$

where  $A = \left( \frac{1}{\sigma^2} x x^T + \frac{1}{\sigma^2} I \right)$

We could now compute a MAP estimate for  $\phi$  to maximize this. It will of course be the mean of this Gaussian posterior

$$\hat{\phi}_{\text{map}} = \frac{1}{\sigma^2} A^{-1} x w$$



Alternatively, we could do full Bayesian inference for a new ~~test~~ measurement  $x^*$  to compute the distribution of values over predicted world state  $w$

Recall this is essentially computing the expected value of  $P(w | x^*, \phi)$  with respect to the posterior distribution  $P(\phi | x, w)$

$$\begin{aligned} P(w | x^*, x, w) &= \int_{\phi} P(w | x^*, \phi) P(\phi | x, w) d\phi \\ &= \int N(w | x^{*T} \phi, \sigma^2) N(\phi | \frac{1}{\sigma^2} A^{-1} X w, A^{-1}) d\phi \end{aligned}$$

It turns out that this is also a Gaussian,

$$= N(w | \underbrace{\frac{1}{\sigma^2} X^{*T} A^{-1} X w}_{\text{mean}}, \underbrace{X^{*T} A^{-1} X^* + \sigma^2}_{\text{variance}})$$