

## LECTURE 17 - BAYESIAN REGRESSION & ML, MAP CONTINUED...

### 1. REGRESSION, CONT.

- Look back our polynomial regression:

$$(1) \quad \min E^*(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (y(x_n, \mathbf{w}) - t_n)^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

This is equivalent to:

$$(2) \quad \max \prod_{n=1}^N \exp \left\{ -\frac{1}{2} (y(x_n, \mathbf{w}) - t_n)^2 \right\} \exp \left\{ -\frac{\lambda}{2} \|\mathbf{w}\|_2^2 \right\}$$

- As discussed, the first term is the Likelihood and the second term is the prior on the weights
- These are Gaussian distributions:

$$(3) \quad \mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2} \frac{(x - \mu)^2}{\sigma^2} \right\}$$

- $\sigma^2$  is the variance OR  $\frac{1}{\sigma^2}$  is the *precision*
- So, as  $\lambda$  gets big, variance gets smaller/tighter. As  $\lambda$  gets small, variance gets larger/wider.
- Previously, we used:

$$(4) \quad y = \sum_{j=0}^M w_j x^j$$

- We can extend this, make is more general and flexible:

$$(5) \quad y = \sum_{j=0}^M w_j \phi_j(\mathbf{x})$$

where  $\phi_j(\mathbf{x})$  is a *basis function*

- For example:
  - Basis function we were using previously:  $\phi_j(x) = x^j$  (for univariate  $x$ )
  - Linear Basis Function:  $\phi_j(\mathbf{x}) = x_j$
  - Radial Basis Function:  $\phi_j(\mathbf{x}) = \exp \left\{ -\frac{(x - \mu_j)^2}{2s_j^2} \right\}$
  - Sigmoidal Basis Function:  $\phi_j(\mathbf{x}) = \frac{1}{1 + \exp \left\{ \frac{\mathbf{x} - \mu_j}{s} \right\}}$

- As before:

$$(6) \quad t = y(\mathbf{x}, \mathbf{w}) + \epsilon$$

- However, now:

$$(7) \quad y = \mathbf{w}^T \Phi(\mathbf{x}) = [w_0, w_1, \dots, w_M][\phi_0(\mathbf{x}), \phi_1(\mathbf{x}), \dots, \phi_M(\mathbf{x})]^T$$

where  $\epsilon \sim \mathcal{N}(\cdot|0, \beta^{-1})$

$$(8) \quad p(t|\mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n|\mathbf{w}^T \Phi(\mathbf{x}_n), \beta^{-1})$$

- So, what is the “trick” to use to maximize this?

$$(9) \quad \mathcal{L} = \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta E(\mathbf{w})$$

$$(10) \quad \frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \beta \sum_{n=1}^N (t_n - \mathbf{w}^T \Phi(\mathbf{x}_n)) \Phi(\mathbf{x}_n)^T = 0$$

- This results in:

$$(11) \quad \mathbf{w}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

where

$$(12) \quad \Phi = [\Phi(x_1), \Phi(x_2), \dots]$$

- What would you do if you want to include a prior? get the MAP solution? If assuming zero-mean Gaussian noise, then Regularized Least Squares!

## 2. BAYESIAN LINEAR REGRESSION

- Recall:  $E_D(\mathbf{w}) + \lambda E_W(\mathbf{w})$  where  $\lambda$  is the trade-off regularization parameter
- A simple regularizer (and the one we used previously) is:  $E_W(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w}$
- If we assume zero-mean Gaussian noise:  $E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2$
- Then, the total error becomes:  $\frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$
- We can take the derivative, set it equal to zero and solve for the weights. When we do, we get:

$$(13) \quad \mathbf{w} = (\lambda \mathbf{I} + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

- Recall, we can interpret this as:

$$(14) \quad \min_{\mathbf{w}} E^* = \min_{\mathbf{w}} \{E_D(\mathbf{w}) + \lambda E_W(\mathbf{w})\}$$

$$(15) \quad = \max_{\mathbf{w}} \{-E_D(\mathbf{w}) - \lambda E_W(\mathbf{w})\}$$

$$(16) \quad = \max_{\mathbf{w}} \exp \{-E_D(\mathbf{w}) - \lambda E_W(\mathbf{w})\}$$

$$(17) \quad = \max_{\mathbf{w}} \exp \{-E_D(\mathbf{w})\} \exp \{-\lambda E_W(\mathbf{w})\}$$

$$(18) \quad \propto \max_{\mathbf{w}} \prod_{n=1}^N \mathcal{N}(t | \mathbf{w}^T \Phi(\mathbf{x}_n), \beta \mathbf{I}) \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0)$$

$$(19) \quad = \max_{\mathbf{w}} p(\mathbf{t} | \mathbf{w}, \mathbf{X}) p(\mathbf{w})$$

$$(20) \quad \propto \max_{\mathbf{w}} p(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N) = \mathcal{N}(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N)$$

where  $\mathbf{m}_N = \mathbf{S}_N (\mathbf{S}_0 \mathbf{m}_0 + \beta \Phi^T \mathbf{t})$  and  $\mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta \Phi^T \Phi$

- What happens with different values of  $\beta$  and  $\mathbf{S}_0$ ?
- To simplify, let us assume that  $\mathbf{S}_0 = \alpha^{-1} \mathbf{I}$  and  $\mathbf{m}_0 = \mathbf{0}$ , thus,  $\mathbf{m}_N = \beta \mathbf{S}_N \Phi^T \mathbf{t}$  and  $\mathbf{S}_N^{-1} = (\alpha^{-1} \mathbf{I})^{-1} + \beta \Phi^T \Phi = \alpha \mathbf{I} + \beta \Phi^T \Phi$
- This results in the following Log Posterior:

$$(21) \quad \ln p(\mathbf{w} | \mathbf{t}) = -\frac{\beta}{2} \sum_{n=1}^N (t_n - \mathbf{w}^T \Phi(\mathbf{x}_n))^2 - \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} + \text{const}$$

- Let us suppose we are dealing with 1-D data,  $\mathbf{X} = \{x_1, \dots, x_N\}$  and a linear form for  $y$ :  $y(x, \mathbf{w}) = w_0 + w_1 x$
- We are going to generate synthetic data from:  $t = -0.3 + 0.5x + \epsilon$  where  $\epsilon$  is from zero-mean Gaussian noise. The goal is to estimate the true values  $w_0 = -0.3$  and  $w_1 = 0.5$ .
- *Let us assume  $\beta = 25$  and  $\alpha = 2$  and run provided code example and step through it carefully.*