

$$④ \text{ weighted} = 1/2 \sum_{i=1}^n t_i (\omega^T \phi(x_i) - y_i)^2 = 1/2 \sum_{i=1}^n t_i (\omega^T \phi(x_i) - 2\omega^T \phi(x_i) y_i + y_i^2)$$

$$\begin{aligned} & \frac{1}{2} (\Phi \omega - y)^T T (\Phi \omega - y) \quad , \text{ with } T \rightarrow \mathbb{R}^{n \times n} \\ & = \frac{1}{2} ((\Phi \omega)^T - y^T) (T (\Phi \omega - y)) \\ & = \frac{1}{2} (\omega^T \Phi^T T \Phi \omega - \omega^T \Phi^T T y - y^T T \Phi \omega + y^T T y) \end{aligned}$$

$$\begin{aligned} \text{By parts:} \\ ① \nabla_{\omega} (\omega^T \underbrace{\Phi^T T \Phi}_{Z} \omega) \quad Z &= \Phi^T T \Phi \\ &= \nabla_{\omega} (\omega^T Z \omega) \\ &= (Z + Z^T) \cdot \omega \\ &= (\Phi^T T \Phi + (\Phi^T T \Phi)^T) \cdot \omega \\ &= (\Phi^T T \Phi + \Phi^T T^T \Phi) \cdot \omega \\ \text{Since } T^T \text{ is a diagonal matrix } T^T &= T \\ &= 2 \cdot (\Phi^T T \Phi) \cdot \omega \\ ② \nabla_{\omega} (\omega^T \underbrace{\Phi^T T}_{Z} y) \quad Z &= \Phi^T T y \\ &= \nabla_{\omega} (\omega^T Z) \\ &= \Phi^T T y \\ ③ \nabla_{\omega} (y^T T \Phi \omega) \quad ④ \nabla_{\omega} (y^T T y) &= 0 \\ &= \nabla_{\omega} ((\Phi^T T y)^T \omega) \\ &= \Phi^T T y \\ &= \Phi^T T y \end{aligned}$$

$$\begin{aligned} \nabla_{\omega} E_{\text{weighted}} &= \frac{1}{2} [2 \cdot (\Phi^T T \Phi) \cdot \omega - \Phi^T T y \cdot 2] \\ &= (\Phi^T T \Phi) \cdot \omega - \Phi^T T y \\ \nabla_{\omega} E_{\text{weighted}} = 0 &\Leftrightarrow (\Phi^T T \Phi) \cdot \omega = \Phi^T T y \\ &\Leftrightarrow \omega^* = (\Phi^T T \Phi)^{-1} \cdot \Phi^T T y \end{aligned}$$

- 1) If we attribute weights to each datapoint our model gains flexibility. For example, we can assign less weight to outliers that follow no perceivable pattern, thus minimizing the noise in our data.
- 2) Repeated values bring no value to the data and thus can be ignored in the learning process by assigning them a 0 weighting factor.

$$⑤ \quad \hat{\Phi} = \begin{pmatrix} \frac{1}{\sqrt{\lambda}} \mathbf{I}_m \\ \mathbf{0} \end{pmatrix}, \quad \hat{\mathbf{y}} = \begin{pmatrix} \mathbf{y} \\ 0 \end{pmatrix}, \quad \ell_{LS}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^m (\mathbf{w}^T \phi_i(x_i) - y_i)^2$$

$$\Rightarrow \sum_{i=1}^m (\mathbf{w}^T \phi_i(x_i) - y_i)^2 = [\mathbf{w}^T \hat{\Phi} - \hat{\mathbf{y}}]^T [\mathbf{w}^T \hat{\Phi} - \hat{\mathbf{y}}] \quad \left| \quad \begin{aligned} \hat{\Phi} &= \begin{bmatrix} \mathbf{I} & \sqrt{\lambda} \mathbf{I} \end{bmatrix} \\ \hat{\mathbf{y}} &= \begin{pmatrix} \mathbf{y} \\ 0 \end{pmatrix} \end{aligned} \right.$$

$$\text{ridge} = [\mathbf{w}^T \hat{\Phi} - \hat{\mathbf{y}}]^T [\mathbf{w}^T \hat{\Phi} - \hat{\mathbf{y}}]$$

$$\Rightarrow \mathbf{w}^T \hat{\Phi} - \hat{\mathbf{y}} = \mathbf{w}^T \begin{bmatrix} \mathbf{I} & \sqrt{\lambda} \mathbf{I} \end{bmatrix} - \hat{\mathbf{y}}$$

$$\Rightarrow \mathbf{w} \begin{bmatrix} \mathbf{I} & \sqrt{\lambda} \mathbf{I} \end{bmatrix} - \hat{\mathbf{y}} = \begin{bmatrix} w_1 \\ \vdots \\ w_m \\ w_1 \sqrt{\lambda} \\ \vdots \\ w_m \sqrt{\lambda} \end{bmatrix} - \begin{bmatrix} y_1 \\ \vdots \\ y_m \\ 0 \\ \vdots \\ 0 \end{bmatrix} = \begin{bmatrix} w_1 - y_1 \\ \vdots \\ w_m - y_m \\ w_1 \sqrt{\lambda} \\ \vdots \\ w_m \sqrt{\lambda} \end{bmatrix} \quad \left\{ \right. \}_{m}$$

$$\Rightarrow \frac{1}{2} \left([\mathbf{w} \begin{bmatrix} \mathbf{I} & \sqrt{\lambda} \mathbf{I} \end{bmatrix} - \hat{\mathbf{y}}]^T [\mathbf{w} \begin{bmatrix} \mathbf{I} & \sqrt{\lambda} \mathbf{I} \end{bmatrix} - \hat{\mathbf{y}}] \right) = \frac{1}{2} \left(\sum_{i=1}^m (\mathbf{w}^T \phi_i(x_i) - y_i)^2 + \sum_{i=1}^m (w_i \sqrt{\lambda})^2 \right)$$

$$= \frac{1}{2} \sum_{i=1}^m (\mathbf{w}^T \phi_i(x_i) - y_i)^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 //$$

⑥

John is allowing his polynomial regression model to go up to degree 50, thus losing all generality associated with allowing some error % to occur when using lower degrees.

By trying to minimize the error mindlessly, he happens to be overfitting his model to the training data and is learning a function that has a lot of variance in order to match all data points to the predicted values.

This could be solved by splitting the data into a training and validation set. Such that the minimization of the error function would occur in unseen data and actually be learning a model that encapsulates the behaviour of his data. Also, we could use K-fold cross validation and learn what degree has the least error on average regarding all folds and thus create the most generalizable model.

⑦

$$\text{likelihood} \rightarrow p(y_i | \Phi, w, \beta) = \prod_{i=1}^N \mathcal{N}(y_i | w^T \phi(x_i), \beta^{-1})$$

$$\text{conjugate prior} \rightarrow p(w, \beta) = \mathcal{N}(w | m_0, \beta^{-1} S_0) \text{Gamma}(\beta | a_0, b_0)$$

$$\frac{1}{\sigma \sqrt{2\pi}} e^{-1/2 \left(\frac{x-\mu}{\sigma} \right)^2}$$

$$\mathcal{N}(y_i | w^T \phi(x_i), \beta^{-1}) \rightarrow \prod_{i=1}^N \frac{\beta}{\sqrt{2\pi}} e^{-1/2 (y_i - w^T \phi(x_i) \beta)^2} \propto e^{-1/2 \left[\sum_{i=1}^N (y_i - w^T \phi(x_i) \beta)^2 \right]}$$

$$\mathcal{N}(w | m_0, \beta^{-1} S_0) \rightarrow \frac{\beta}{S_0 \sqrt{2\pi}} e^{-1/2 \left(\frac{w - m_0}{S_0} \beta \right)^2}$$

$$\text{Gamma}(\beta | a_0, b_0) \rightarrow \frac{\Gamma(a_0 + b_0)}{\Gamma(a_0) \Gamma(b_0)} \beta^{a_0-1} (1-\beta)^{b_0-1}$$

$$\begin{aligned} \rightarrow p(w, \beta | y) &= \prod_{i=1}^N \mathcal{N}(y_i | w^T \phi(x_i), \beta^{-1}) \times \mathcal{N}(w | m_0, \beta^{-1} S_0) \text{Gamma}(\beta | a_0, b_0) \\ &= \left(\frac{\beta}{\sqrt{2\pi}} \right)^N e^{-1/2 \left[\sum_{i=1}^N (y_i - w^T \phi(x_i) \beta)^2 \right]} \left(\frac{\beta}{S_0 \sqrt{2\pi}} \right) e^{-1/2 \left(\frac{w - m_0}{S_0} \beta \right)^2} \frac{\Gamma(a_0 + b_0)}{\Gamma(a_0) \Gamma(b_0)} \beta^{a_0-1} (1-\beta)^{b_0-1} \end{aligned}$$

$$= e^{-\frac{\beta^2}{2} \left(\sum_{i=1}^N (y_i - w^T \phi(x_i))^2 + \left(\frac{w - m_0}{S_0} \right)^2 \right)} \beta^{a_0-1} (1-\beta)^{b_0-1} \frac{1}{S_0 \sqrt{2\pi}} \frac{\Gamma(a_0 + b_0)}{\Gamma(a_0) \Gamma(b_0)}$$

$$= e^{-\frac{\beta^2}{2} \left(y^T y - 2y^T \Phi w + w^T \Phi^T \Phi w + \left(\frac{w - m_0}{S_0} \right)^2 \right)} \beta^{a_0-1} (1-\beta)^{b_0-1} \frac{1}{S_0 \sqrt{2\pi}} \frac{\Gamma(a_0 + b_0)}{\Gamma(a_0) \Gamma(b_0)}$$

$$= e^{-\frac{\beta^2}{2} \left(y^T y - 2y^T \Phi w + w^T \Phi^T \Phi w + \left(\frac{w - m_0}{S_0} \right)^2 \right)} \beta^{a_0+N} (1-\beta)^{b_0-1} \frac{1}{S_0 \sqrt{2\pi}} \frac{\Gamma(a_0 + b_0)}{\Gamma(a_0) \Gamma(b_0)}$$

$$\ln \rightarrow -\frac{\beta^2}{2} \left(y^T y - 2y^T \Phi w + w^T \Phi^T \Phi w + \left(\frac{w - m_0}{S_0} \right)^2 \right) + (a_0 + N) \ln \beta + (b_0 - 1) \ln (1-\beta) - \ln(S_0 \sqrt{2\pi}) - \ln \left(\frac{\Gamma(a_0 + b_0)}{\Gamma(a_0) \Gamma(b_0)} \right)$$

$$\bullet \mathcal{N}(w | m_N, \beta^{-1} S_N) \text{Gamma}(\beta | a_N, b_N) = \frac{\beta}{S_N \sqrt{2\pi}} e^{-1/2 \left(\frac{w - m_N}{S_N} \beta \right)^2} \cdot \frac{\Gamma(a_N + b_N)}{\Gamma(a_N) \Gamma(b_N)} \beta^{a_N-1} (1-\beta)^{b_N-1}$$

$$\ln \rightarrow \ln \left(\frac{\beta}{S_N \sqrt{2\pi}} \right) - 1/2 \left(\frac{w - m_N}{S_N} \beta \right)^2 + \ln \left(\frac{\Gamma(a_N + b_N)}{\Gamma(a_N) \Gamma(b_N)} \right) + (a_N - 1) \ln \beta + (b_N - 1) \ln (1-\beta) =$$

$$= -\ln(S_N \sqrt{2\pi}) - 1/2 \left(\frac{w - m_N}{S_N} \beta \right)^2 + \ln \left(\frac{\Gamma(a_N + b_N)}{\Gamma(a_N) \Gamma(b_N)} \right) + (a_N) \ln \beta + (b_N - 1) \ln (1-\beta)$$

$$\bullet \frac{1}{2} \left(y^T y - 2y^T \Phi w + w^T \Phi^T \Phi w + \left(\frac{w - m_0}{s_0} \right)^2 \right) + (a_0 + N) \ln \beta + (b_0 - 1) \ln(1 - \beta) - \ln(s_0 \cdot (\sqrt{2\pi})^{N+1}) - \ln \left(\frac{\gamma(a_0 + b_0)}{\gamma(a_0) \gamma(b_0)} \right)$$

$$\bullet -\ln(s_N \sqrt{2\pi}) - \frac{1}{2} \left(\frac{(w - m_N)^T \beta}{s_N} \right)^2 + \ln \left(\frac{\gamma(a_N + b_N)}{\gamma(a_N) \gamma(b_N)} \right) + (a_N) \ln(\beta) + (b_N - 1) \ln(1 - \beta)$$

$$\Rightarrow \begin{cases} s_N = s_0 (\sqrt{2\pi})^N \\ m_N = * \text{ (found on the next page)} \\ a_N = a_0 + N \\ b_N = b_0 \end{cases}$$

② we have $E_{\text{ridge}}(w) = \frac{1}{2} \sum_{i=1}^N (w^T \phi(x_i) - y_i)^2 + \frac{\lambda}{2} w^T w$

$$\Rightarrow \nabla_{\text{ridge}} = \nabla_{\text{LS}} + \nabla_w \left(\frac{\lambda}{2} \|w\|_2^2 \right) \stackrel{\text{from lecture notes}}{=} \left(X^T X w - X^T y \right) + \lambda w$$

$$\bullet \|w\|_2^2 = \left(\sum_{i=1}^N w_i^2 \right)^{1/2} = \sum_{i=1}^N w_i^2$$

$$\Rightarrow \frac{\partial}{\partial w_k} \|w\|_2^2 = \frac{\partial}{\partial w_k} \sum_{i=1}^N w_i^2 = \sum_{i=1}^N \frac{\partial}{\partial w_k} w_i^2 = 2w_k$$

$$\Rightarrow \nabla_w (w^T w) = 2w //$$

$$\stackrel{\text{max}}{\Rightarrow} X^T X w - X^T y + \lambda w = 0 \Leftrightarrow (X^T X + \lambda I) w = X^T y \Leftrightarrow w = \frac{X^T y}{X^T X + \lambda I} \quad \left(\Leftrightarrow w = X^T y (X^T X + \lambda I)^{-1} \right)$$

assuming $X^T X + \lambda I \neq 0$
 $\Rightarrow \det \neq 0$

Since the amount of basis functions behaves like the degrees of freedom in the standard polynomial regression. If we were to have less data than basis functions, we would be encountering some form of overfitting in our models which is observed by the growing norm of the weight vector.

With regularization, since the norm is taken into account in the minimization portion of the training, it is expected that the level of overfitting should be drastically reduced.

$$(*) \sum_{i=1}^N (y_i - w^T \phi(x_i))^2 + \left(\frac{w - m_0}{z} \right)^2 = \left(\frac{w - m_N}{s_N} \right)^2 = \left(\frac{w - m_N}{s_0 \sqrt{2N}} \right)^2 \quad (=)$$

$$\Leftrightarrow \left(\sum_{i=1}^N (y_i - w^T \phi(x_i))^2 + \left(\frac{w - m_0}{z} \right)^2 \right) s_0^2 (2N)^N = (w - m_N)^2 \Leftrightarrow$$

$$\Leftrightarrow m_N = w - \left[s_0^2 2N^N \left(\sum_{i=1}^N (y_i - w^T \phi(x_i))^2 + \left(\frac{w - m_0}{s_0} \right)^2 \right) \right]^{1/2}$$

\Rightarrow note that this is likely wrong since m_N shouldn't depend on w

9)

a) prediction is $f(X) = w^T X$

\rightarrow find (w_{new}) such that $w_{\text{new}}^T X_{\text{new}} = w^T X \Leftrightarrow$

$$\Leftrightarrow w_{\text{new}}^T aX = w^T X \Leftrightarrow$$

$$\Leftrightarrow w_{\text{new}}^T = \frac{1}{a} w^T \Leftrightarrow w_{\text{new}} = 1/a \cdot w$$

b) from problem 8 we have that $w = X^T y (X^T X + \lambda I)^{-1}$ in the original ridge regression problem

$$\Rightarrow X_{\text{new}}^T y (X_{\text{new}}^T X_{\text{new}} + \lambda_{\text{new}} I)^{-1} = X^T y (X^T X + \lambda I)^{-1} \Leftrightarrow$$

$$\Leftrightarrow a X^T y (a^2 X^T X + X_{\text{new}}^T I)^{-1} = X^T y (X^T X + \lambda I)^{-1} \Leftrightarrow$$

$$\Leftrightarrow a (a^2 X^T X + \lambda_{\text{new}} I)^{-1} = (X^T X + \lambda I)^{-1} \Leftrightarrow$$

$$\Leftrightarrow a (X^T X + \lambda I) = a^2 X^T X + \lambda_{\text{new}} I \Leftrightarrow$$

$$\Leftrightarrow \lambda_{\text{new}} = (a - a^2) X^T X + \lambda I$$