

# Lecture 13, Feb 26, 2024

## Linear Regression

- Consider a linear model  $Y = \mathbf{w}^T \mathbf{X} + Z$  where we have  $n$  noisy measurements  $y_i$  from  $n$  inputs  $\mathbf{x}_i$ 
  - Assume  $Z$  is some IID Gaussian random noise
  - Given these measurements, our goal is to find the best set of weights  $\mathbf{w}^T = [w_1 \ \dots \ w_D]$
  - Each weight  $w_j$  corresponds to the  $j$ th coefficient of  $\mathbf{x}$ , which has dimension  $D$

- Form the *design matrix*  $\begin{bmatrix} \mathbf{x}_1^T & y_1 \\ \vdots & \vdots \\ \mathbf{x}_n^T & y_n \end{bmatrix}$

- Consider the MLE  $\hat{\mathbf{w}}_{\text{ML}} = \underset{\mathbf{w} \in \mathbb{R}^D}{\text{argmax}} \log p((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n) | \mathbf{w})$

$$\begin{aligned} &= \underset{\mathbf{w} \in \mathbb{R}^D}{\text{argmax}} \log \prod_{i=1}^n p(\mathbf{x}_i, y_i | \mathbf{w}) \\ &= \underset{\mathbf{w} \in \mathbb{R}^D}{\text{argmax}} \sum_{i=1}^n \log \left( \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(y_i - \mathbf{w}^T \mathbf{x}_i)^2} \right) \\ &= \underset{\mathbf{w} \in \mathbb{R}^D}{\text{argmax}} - \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i)^2 \\ &= \underset{\mathbf{w} \in \mathbb{R}^D}{\text{argmin}} \mathbf{e}^T(\mathbf{w})\mathbf{e}(\mathbf{w}) \end{aligned}$$

- Where the error vector is  $\mathbf{e}(\mathbf{w}) = \begin{bmatrix} y_1 - \mathbf{w}^T \mathbf{x}_1 \\ \dots \\ y_n - \mathbf{w}^T \mathbf{x}_n \end{bmatrix} = \mathbf{y} - \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix} \mathbf{w}$

- This is now a *least squares regression problem*

- Let  $\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix}$  then we have  $\underset{\mathbf{w} \in \mathbb{R}^D}{\text{argmin}} (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w})$

- Expand:  $\frac{1}{2} \mathbf{w}^T \mathbf{X} \mathbf{X} \mathbf{w} - \mathbf{w} \mathbf{X}^T \mathbf{y} + \frac{1}{2} \mathbf{y}^T \mathbf{y}$  (note a factor of  $\frac{1}{2}$  was added)

- Derivative:  $\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y} = 0 \implies \mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y}$

- Therefore  $\hat{\mathbf{w}}_{\text{ML}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

- Another way to write this is  $\mathbf{X}^T (\mathbf{X} \mathbf{w} - \mathbf{y}) = 0$ , meaning we can interpret this as making the error vector orthogonal to all the input data

- This is the *normal equation*

- $\mathbf{X}^T \mathbf{X}$  is the *scatter matrix*

- This is an estimate of the covariance/correlation matrix of the data

- Regression can be performed in any general vector space, so our model can be nonlinear in  $\mathbf{x}$  (but still linear in  $\mathbf{w}$ )

- In general given any basis function  $\phi(\mathbf{x}_i)$  we can try to fit  $y_i = \mathbf{w}^T \phi(\mathbf{x}_i) + z_i$

- Let  $\mathbf{X} = \begin{bmatrix} \phi^T(\mathbf{x}_1) \\ \vdots \\ \phi^T(\mathbf{x}_n) \end{bmatrix}$  then  $\hat{\mathbf{w}}_{\text{ML}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

- e.g. we can work in the vector space of polynomials to perform polynomial regression, or the space of sinusoids for a Fourier series

\* For  $d$ -degree polynomial regression we'd have  $\phi^T(x_i) = [1 \ x_i \ x_i^2 \ \dots \ x_i^d]$

- Example: measuring the height of a cannonball  $h_i$  vs. time  $t_i$  for  $i = 1, \dots, n$

\* Use the model  $h_i = w_1 t_i + w_2 t_i^2 + z_i = \mathbf{w}^T \mathbf{x}_i + z_i$  where  $\mathbf{x}_i = \begin{bmatrix} t_i \\ t_i^2 \end{bmatrix}$

## Bayesian Regression – Regularization

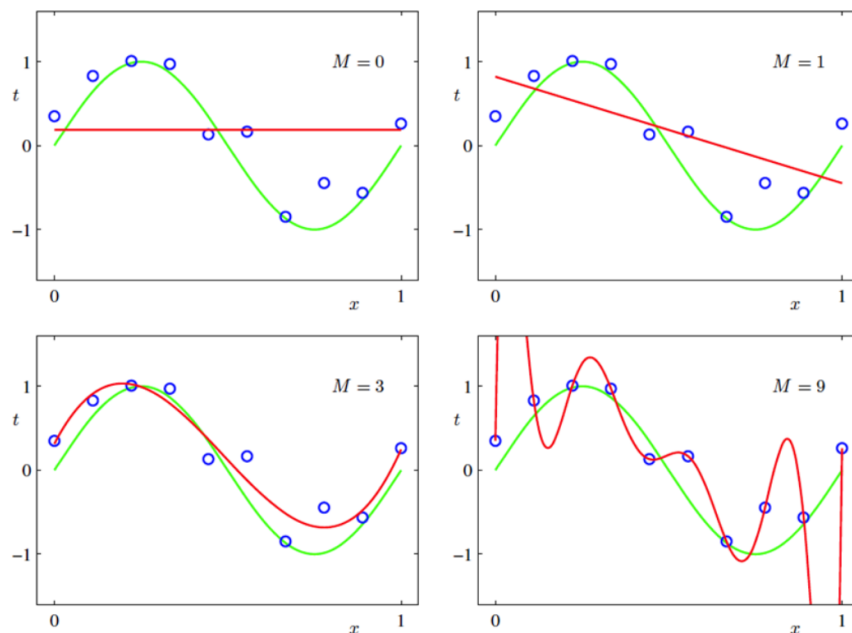


Figure 1: Polynomial regression for different degrees. Green is the underlying function we’re trying to approximate.

- If we make the model too complex, i.e. too high of a dimension for  $\phi$ , we will get overfitting
- Typically when the model overfits, we get very large weights that are not physically realistic for our system
  - To keep the weights down, we can use regularization
  - Here we show a way to derive the same result by instead assuming a prior on  $\mathbf{w}$
- Assume that each weight has a prior  $w_i \sim \mathcal{N}(0, \tau^2)$ ; now can find the MAP estimate
- $\hat{\mathbf{w}}_{\text{MAP}} = \underset{\mathbf{w} \in \mathbb{R}^D}{\text{argmax}} p((\phi(x_1), y_1), \dots, (\phi(x_n), y_n))p(\mathbf{w})$

$$= \underset{\mathbf{w} \in \mathbb{R}^D}{\text{argmax}} \sum_{i=1}^n \left[ \log \left( \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2} (y_i - \mathbf{w}^T \phi(\mathbf{x}_i))^2} \right) + \sum_{j=1}^D \log \frac{1}{\sqrt{2\pi}\tau} e^{-\frac{w_j^2}{2\tau^2}} \right]$$

$$= \underset{\mathbf{w} \in \mathbb{R}^D}{\text{argmin}} \sum_{i=1}^n (y_i - \mathbf{w}^T \phi(\mathbf{x}_i))^2 + \frac{\sigma^2}{\tau^2} \|\mathbf{w}\|^2$$

$$= \underset{\mathbf{w} \in \mathbb{R}^D}{\text{argmin}} \sum_{i=1}^n (y_i - \mathbf{w}^T \phi(\mathbf{x}_i))^2 + \lambda \|\mathbf{w}\|^2$$

- The first term is the same least squares term as before, but now we have an additional term that penalizes the norm of  $\mathbf{w}$ , effectively keeping the weights small

$$\text{– Let } \mathbf{e}(\mathbf{w}) = \begin{bmatrix} y_1 \\ \vdots \\ y_n \\ \mathbf{0} \end{bmatrix} - \begin{bmatrix} \phi^T(\mathbf{x}_1) \\ \vdots \\ \phi^T(\mathbf{x}_n) \\ -\sqrt{\lambda} \mathbf{1} \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_D \end{bmatrix} \text{ and } \tilde{\mathbf{X}} = \begin{bmatrix} \phi^T(\mathbf{x}_1) \\ \vdots \\ \phi^T(\mathbf{x}_n) \\ -\sqrt{\lambda} \mathbf{1} \end{bmatrix}$$

\* The error can again be written as  $\mathbf{e}^T(\mathbf{w})\mathbf{e}(\mathbf{w})$

- Using the same derivation as before,  $\hat{\mathbf{w}}_{\text{MAP}} = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \begin{bmatrix} \mathbf{y} \\ \mathbf{0} \end{bmatrix} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{1})^{-1} \mathbf{X}^T \mathbf{y}$
- Notice that this result is almost the same as the MLE solution, except with the addition of  $\lambda \mathbf{1}$
- This is known as *ridge regression*

- We can also solve this by writing it as a Gaussian system

$$- \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix} \begin{bmatrix} W_1 \\ \vdots \\ W_D \end{bmatrix} + \begin{bmatrix} Z_1 \\ \vdots \\ Z_n \end{bmatrix} \iff \mathbf{y} = \mathbf{X}\mathbf{w} + \mathbf{z}$$

$$- \hat{\mathbf{w}}_{\text{MAP}}(\mathbf{y}) = (\boldsymbol{\Sigma}_W^{-1} \mathbf{X}^T \boldsymbol{\Sigma}_Z^{-1} \mathbf{X})^{-1} \left( \mathbf{X}^T \frac{1}{\sigma^2} \mathbf{y} + \frac{1}{\tau^2} \cdot 0 \right) = (\mathbf{X}^T \mathbf{X} + \frac{\sigma^2}{\tau^2} \mathbf{1})^{-1} \mathbf{X}^T \mathbf{y}$$

$$- \text{This gives us the conditional precision } \boldsymbol{\Sigma}_{\mathbf{X}|\mathbf{Y}}^{-1} = \frac{1}{\tau^2} \mathbf{1} + \frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X}$$

$$- \text{So } \boldsymbol{\Sigma}_{\mathbf{X}|\mathbf{Y}} = \sigma^2 \left( \mathbf{X}^T \mathbf{X} + \frac{\sigma^2}{\tau^2} \mathbf{1} \right)^{-1}$$

\* Notice that  $\mathbf{X}^T \mathbf{X} + \frac{\sigma^2}{\tau^2} \mathbf{1}$  is related to the covariance

\* As we collect more data the added term becomes negligible

\*  $\mathbf{X}^T \mathbf{X}$  becomes bigger so the covariance shrinks