

# Lecture 6, Feb 2, 2024

## Generalized Linear Models (GLMs)

### Definition

*Generalized Linear Models:* A GLM is given by

$$\hat{f}(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{i=1}^{M-1} w_i \phi_i(\mathbf{x})$$

where  $\mathbf{w}$  is a set of undetermined weights and  $\phi_i: \mathbb{R}^D \mapsto \mathbb{R}$  are a set of known basis functions.

- The models may be nonlinear in the inputs  $\mathbf{x}$ , but still linear in the weights  $\mathbf{w}$ , which makes it still possible to use linear techniques
- To construct a GLM we need to select the appropriate basis functions, and formulate a strategy to estimate the weights

- Let  $\phi_0(\mathbf{x}) = 1$  (the bias term) and  $\phi(\mathbf{x}) = \begin{bmatrix} \phi_0(\mathbf{x}) \\ \phi_1(\mathbf{x}) \\ \vdots \\ \phi_{M-1}(\mathbf{x}) \end{bmatrix} \in \mathbb{R}^M$

- Then if we define the weight vector  $\mathbf{w} = [w_0 \ \dots \ w_m]^T$ , we can write  $\hat{f}(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \phi(\mathbf{x})$
- We are using  $\phi$  to map from the input space  $\mathcal{X}$  to the *feature space*  $\mathcal{F}$ , and performing linear regression in the feature space

- Let the vector of training targets  $\mathbf{y} = [y^{(1)} \ \dots \ y^{(N)}]^T \in \mathbb{R}^N$  and  $\Phi \in \mathbb{R}^{N \times M}$  where the  $i$ th row contains  $\phi(\mathbf{x}^{(i)}) \in \mathbb{R}^M$

- Then  $\hat{\mathbf{y}} = \Phi \mathbf{w}$

- Use the  $l_2$  loss function  $\hat{\mathbf{w}} = \underset{\mathbf{w} \in \mathbb{R}^M}{\operatorname{argmin}} \|\mathbf{y} - \hat{\mathbf{y}}\|_2^2$

- Again the loss function can be written as  $(\mathbf{y} - \Phi \mathbf{w})^T (\mathbf{y} - \Phi \mathbf{w})$

- We can use the same techniques for the linear model, but instead of  $D + 1$  weights we have  $M$  weights

- We essentially replace  $\mathbf{X} \in \mathbb{R}^{N \times (D+1)}$  with  $\Phi \in \mathbb{R}^{N \times M}$

- Derivation:

- $\mathcal{L}(\mathbf{w}) = \mathbf{y}^T \mathbf{y} + \mathbf{w}^T \Phi^T \Phi \mathbf{w} - 2\mathbf{y}^T \Phi \mathbf{w}$

- $\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = (\Phi^T \Phi + (\Phi^T \Phi)^T) \mathbf{w} - 2\Phi^T \mathbf{y} = 2\Phi^T \Phi \mathbf{w} - 2\Phi^T \mathbf{y} = 0$

- Therefore  $\Phi^T \Phi \mathbf{w} = \Phi^T \mathbf{y}$

- We can use the same techniques as linear models to solve the normal equations:

- Cholesky factorization

- \* Avoid this because the condition number is squared

- (Economy) QR factorization

- \* Use this only if  $\Phi$  is not rank-deficient

- (Economy) SVD, or truncated SVD if  $\Phi$  is rank-deficient

- \* Slowest, but the most stable

## Polynomial Regression

- The basis functions are the univariate polynomials  $\{1, x_i, x_i^2, \dots, x_i^p\}$  up to order  $p$

- If  $D = 1$ , then the basis functions are  $\phi_i = x^i$  so  $\hat{f}(x) = w_0 + w_1 x + \dots + w_p x^p$

- Note taking tensor products of higher-order univariate polynomials is not a good idea since we will generate  $p^D$  basis functions

- For arbitrary  $D$  we would have  $\phi(\mathbf{x}) = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_D \\ x_1x_2 \\ \vdots \\ x_{D-1}x_D \\ \vdots \\ x_1x_2 \dots x_D \end{bmatrix}$
- We can circumvent this with the *kernel trick* covered later

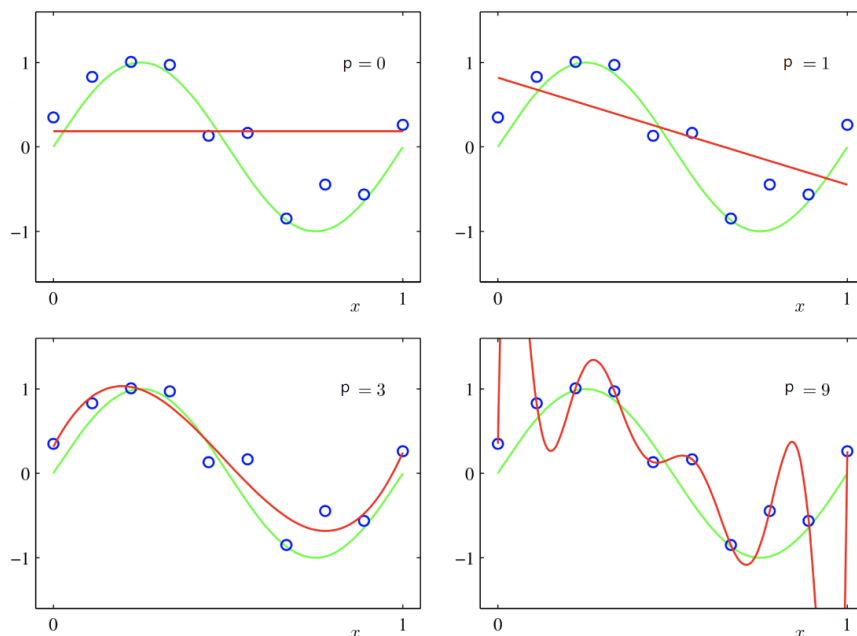


Figure 1: 1D polynomial regression for various values of  $p$ . (Note  $M = p + 1$ .)

- Results for various values of  $p$  are shown above
  - Notice that for smaller values of  $p$  the model doesn't fit well since it doesn't have enough complexity (underfitting)
  - But for large values of  $p$ , the polynomial matches the training points perfectly but approximates the underlying function poorly
- To prevent overfitting, we need to restrict the number of features  $M$  (which in this case restricts the degree of the polynomial)
  - Increasing the number of data points also helps but we often can't just get more data
  - What if we know the underlying model is complex but we don't have enough data points?
  - How do we deal with noise?

## Regularization

- One pattern we may notice is that when the model is overfitting ( $M$  too large), the weights start becoming very big in magnitude
- *Regularization* tries to keep the magnitudes of the weights reasonably small, as a way to prevent overfitting
- To keep the weight small, we can introduce the norm of the weights to the loss function, so the model is penalized for having weights that are too large

## Definition

*Ridge Regression Method:* Choose the weights as

$$\hat{\mathbf{w}} = \underset{\mathbf{w} \in \mathbb{R}^M}{\operatorname{argmin}} \|\mathbf{y} - \Phi \mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$$

where  $\lambda$  is the *regularization parameter*.

- Note  $w_0$  is often excluded from the regularization term
- The regularized loss function is also quadratic in  $\mathbf{w}$ , so we can use the same steps as before
  - Expanded loss:  $\mathbf{y}^T \mathbf{y} + \mathbf{w}^T \Phi^T \Phi \mathbf{w} - 2\mathbf{w}^T \Phi^T \mathbf{y} + \lambda \mathbf{w}^T \mathbf{w}$
  - $\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = 2\Phi^T \Phi \mathbf{w} - 2\Phi^T \mathbf{y} + 2\lambda \mathbf{w} = 0$
  - Rearrange:  $\Phi^T \Phi \mathbf{w} + \lambda \mathbf{w} = \Phi^T \mathbf{y} \implies (\Phi^T \Phi + \lambda \mathbf{1}) \mathbf{w} = \Phi^T \mathbf{y}$ 
    - \* Therefore  $l_2$  regularization is equivalent to adding a small positive perturbation to the diagonal of  $\Phi^T \Phi$
    - \* We saw this in a previous lecture – this also helps with ill-conditioning
    - \* If  $\lambda$  is sufficiently large we can avoid ill-conditioning completely
- Using SVD:  $\Phi = U \Sigma V^T$ 
  - $((U \Sigma V^T)^T U \Sigma V^T + \lambda \mathbf{1}) \mathbf{w} = (U \Sigma V^T)^T \mathbf{y}$
  - Simply to get  $V(\Sigma^T \Sigma + \lambda \mathbf{1}) V^T \mathbf{w} = V \Sigma^T U^T \mathbf{y}$
  - Multiply each side by  $V^T$  to get  $(\Sigma^T \Sigma + \lambda \mathbf{1}) V^T \mathbf{w} = \Sigma^T U^T \mathbf{y}$
  - Therefore  $\mathbf{w} = V(\Sigma^T \Sigma + \lambda \mathbf{1})^{-1} \Sigma^T U^T \mathbf{y}$
  - This can be rewritten as  $\hat{\mathbf{w}}(\lambda) = \sum_{i=1}^M v_i \frac{\sigma_i u_i^T \mathbf{y}}{\sigma_i^2 + \lambda}$ 
    - \*  $v_i, u_i$  are the  $i$ th columns of  $V$  and  $U$
    - \* If  $0 \approx \sigma_i \ll \lambda$  this goes to 0
    - \* If  $\sigma_i \gg \lambda$  this goes to the original unregularized solution
- The regularization has almost no impact on the contributions of large singular values but zeros out the contribution of smaller singular values

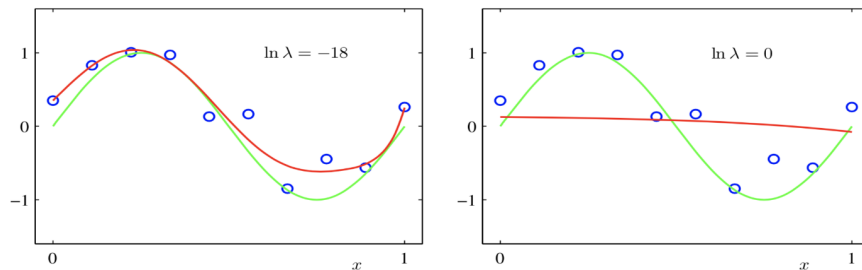


Figure 2: The same polynomial regression from above for  $p = 9$ , with different values of  $\lambda$ .

- $\lambda$  is an important hyperparameter
  - Notice that with a reasonable value of  $\lambda$  we have a pretty good model even at  $p = 9$
  - However if the regularization is too extreme, the model will underfit as the loss is too focused on minimizing  $\|\mathbf{w}\|_2^2$
- To estimate  $\lambda$  we can again use  $\nu$ -fold cross-validation just like we chose  $k$  for  $k$ -NN
  - If the training dataset is small we can use leave-one-out cross-validation (i.e.  $\nu = 1$ )
  - There are fast algorithms for calculating this
  - Use cross-validation to select the best value of  $\lambda$ , then retrain the model on all the data using this new value