

Summary of Week 7

- Categorical predictors
- Interactions
- residual plots

this week:

- other nonlinear models
- overfitting bias variance trade-off
- a bit of statistical learning theory

Statistical learning (ch 2 ISLP)

Consider general problem of predicting Y based on X w/ model

$$E[Y|X] = f(X), \quad Y|X \sim \text{Dist depending on } X$$

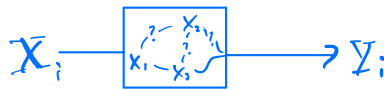
Prediction Vs. inference

Prediction: want to know $Y|X$ for X which are not in our data set

Example: predict CO2 in the future



Inference: want to understand relationships between variables
Example: what % of CO2 variation is driven by seasonal trend.



inform each other

Goal: Understand issues that emerge when adding more complexity to models

(focus on prediction for now)

Especially to work w/ 1 predictor and complex Y vs. X relationship than multiple predictors w/ simple relationship between each X and Y (as we have before)

Polynomial Regression Suppose $X \in \mathbb{R}$

$$Y = \beta_0 + \sum_{j=1}^k \beta_j X^j + \varepsilon \iff Y|X \sim \mathcal{N}\left(\sum_{j=0}^k \beta_j X^j, \sigma^2\right)$$

Can define predictors $X_j = X^j$ and estimate β_j using least squares

Let's review how this works:

$$\text{Cov}(Y, X^m) = \text{Cov}\left(\sum_{j=0}^k \beta_j X^j + \varepsilon, X^m\right) = \sum_{j=0}^k \beta_j \text{Cov}(X^j, X^m)$$

Given data $D = \{(x_1, y_1), \dots, (x_N, y_N)\}$ can approximate

$$\text{Cov}(X^j, X^m) \approx \frac{1}{N} \sum_{i=1}^N (x_i^j - \bar{x}^j)(x_i^m - \bar{x}^m)$$

$$\text{Cov}(Y, X^m) \approx \frac{1}{N} \sum_{i=1}^N (y_i - \bar{y})(x_i^m - \bar{x}^m)$$

and solve for β . This gives same result as

$$\hat{\beta} = \begin{bmatrix} \hat{\beta}_0 \\ \vdots \\ \hat{\beta}_k \end{bmatrix} = \underset{\hat{\beta}}{\text{argmin}} \sum_{i=1}^N \left(y_i - \sum_{j=0}^k \beta_j x_i^j\right)^2$$

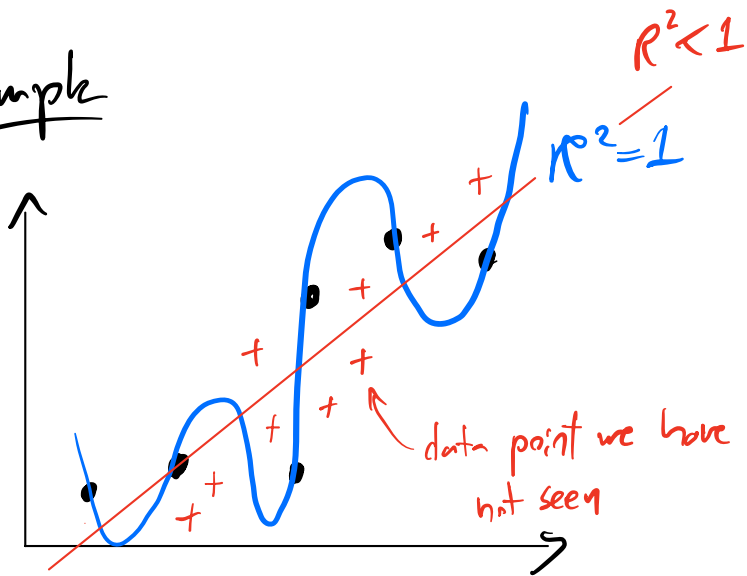
Notation: let $\hat{y}(x; D) = \sum_{j=0}^k \hat{\beta}_j x^j$ when coefficients are computed from data set D

$$\text{and } \text{RSS} = \sum_{i=1}^N (y_i - \hat{y}(x_i; D))^2, \text{ recall } \frac{\text{RSS}}{N} \rightarrow \hat{\sigma}_\varepsilon^2$$

How big should we make K ?

previously we used $R^2 = 1 - \frac{RSS}{\hat{\sigma}_y^2}$, but w/ $K=N-1$ we can always make $R^2=1$

Example



R^2 is computed based on the data we used to fit the coefficients

\Rightarrow doesn't tell us about ability to generalize!

Cross Validation

Break up into two sets

$$D = D^{\text{train}} \cup D^{\text{test}}$$

treat as seen data (green arrow pointing to D^{train})
treat as unseen data (red arrow pointing to D^{test})

② fit the model using only data in D^{train}

③ for each X_i^{test} in D^{test} compute $(X_i^{\text{test}})^j$ for polynomial regression

$$\hat{y}(X_i^{\text{test}}, D^{\text{train}}) = \sum_{j=1}^K \hat{\beta}_j X_{j,i}^{\text{test}}$$

④ look at $\overline{\mathcal{E}_{\text{test}}^2} = \frac{1}{N_{\text{test}}} \text{RSS}_{\text{test}} = \frac{1}{N_{\text{test}}} \sum_{i=1}^{N_{\text{test}}} (\mathcal{Y}_i - \hat{y}(X_i^{\text{test}}, D^{\text{train}}))^2$

[eq. 2.6 in ISLP]

Now we can use $\overline{\mathcal{E}_{\text{test}}^2}$ as measure of how well model predicts new data points!

Also note that we can define

$$\bar{\mathcal{E}}_{\text{train}}^2 = \frac{1}{N_{\text{test}}} \sum_{i=1}^{N_{\text{test}}} (Y_i - \hat{y}(X_i^{\text{train}}, D^{\text{train}}))^2$$

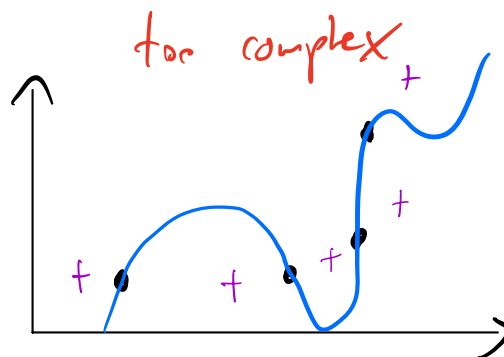
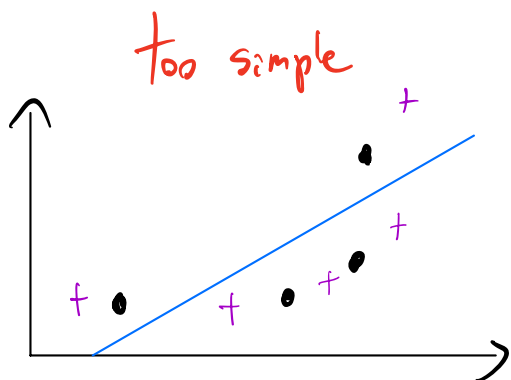
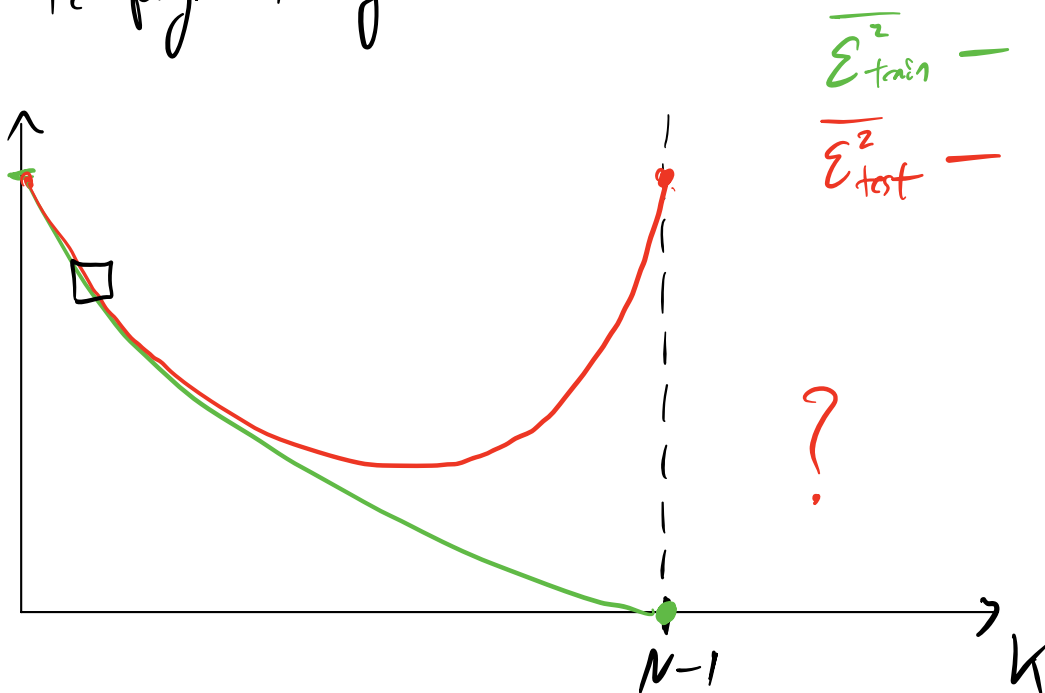
[eq. 2.5 in ISLP]

and notice

$$R^2 = R_{\text{train}}^2 = 1 - \frac{\bar{\mathcal{E}}_{\text{train}}^2}{\hat{\sigma}_{Y^{\text{train}}}^2}$$

We will assume $\hat{\sigma}_Y^2$ is close for test and training points so we can just work w/ $\bar{\mathcal{E}}^2$ s

Now let's look at behavior of $\bar{\mathcal{E}}_{\text{train}}^2$, $\bar{\mathcal{E}}_{\text{test}}^2$ as function of K for the polynomial regression model



Bias Variance Decomposition

let θ be any quantity we would like to estimate
and $\hat{\theta}$ be an estimator of it

in discussion above
 $\hat{\theta} = \hat{y}$, $\theta = E[y|x]$

$$\text{let } \text{MSE}_{\hat{\theta}} = E[(\theta - \hat{\theta})^2]$$

Because adding a constant does not change the variance

$$\begin{aligned} \text{Var}(\hat{\theta}) &= \text{Var}(\hat{\theta} - \theta) = E[(\hat{\theta} - \theta)^2] - E[\hat{\theta} - \theta]^2 \\ &= \text{MSE}_{\hat{\theta}} - \underbrace{E[\hat{\theta} - \theta]^2}_{= \text{Bias}_{\hat{\theta}}^2} \end{aligned}$$

Bias-Variance decomposition

$$\text{MSE}_{\hat{\theta}} = \underbrace{\text{Var}(\hat{\theta})}_{\text{Variance}} + \underbrace{\text{Bias}_{\hat{\theta}}^2}_{\text{Bias}^2}$$

large Bias means result will differ systematically from truth

large Variance means less confident in result / will depend more on specific data set

Example Laplace Rule of Succession (see HW5)

$$\hat{p} = \frac{S+1}{N+2} \quad (\text{compare to } \hat{p} = \frac{S}{N})$$

Now to relate back to $\overline{\Sigma}_{\text{test}}^2$ vs $\overline{\Sigma}_{\text{train}}^2$ picture:

Suppose $y = f(x) = \mathbb{E}[Y|X=x]$ is true relationship between x and avg. of Y
 Then $\hat{\Theta} = \hat{y}(x, D)$ is estimator of $\Theta = f(x)$

$$\overline{\Sigma}_{\text{test}}^2 = \frac{1}{N_{\text{test}}} \sum_{i=1}^{N_{\text{test}}} (Y_i - \hat{y}(X_i, D))^2$$

$$\begin{aligned} \approx \mathbb{E}[(Y - \hat{y})^2] &= \mathbb{E}[(f(x) + \varepsilon - \hat{y}(x, D))^2] \\ &= \mathbb{E}[(f(x) - \hat{y}(x, D))^2 + 2f(x)\varepsilon - 2\varepsilon\hat{y} + \varepsilon^2] \\ &= \text{MSE}_{\hat{y}} + \mathbb{E}[\varepsilon^2] = \text{MSE}_{\hat{y}} + \sigma_{\varepsilon}^2 \end{aligned}$$

↑
 assumes N_{test}
 large

Hence $\overline{\Sigma}_{\text{test}}^2$ almost approximates MSE of \hat{y} , but since we don't know the true relationship there is an extra σ_{ε}^2
 This extra term is called irreducible error

Note when $K \approx N-1$ this relationship breaks down, but it still gives us an idea of why the curve looks the way it does

