

Need each point's leverage statistic. For simple linear regression: leverage "given by"

$$h_i = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

(meaning one variable.)

For multiple linear regression, it's probably something like

$$h_i = \frac{1}{n} + \frac{(x_i - \bar{x})^{12}}{\sum_{i=1}^n (x_i - \bar{x})^{12}}$$

- $h_i \in [0, 1]$, average leverage is always $(p+1)/n$. So if some point has leverage significantly higher than $(p+1)/n$, we should be suspicious of high leverage.
- ★ • Outliers with high leverage can really screw up a model. (Plot standardized residuals vs. leverage to see.)

2/19/23 6. Collinearity

Day 8

- Collinearity refers to the situation in which two or more predictor variables are closely related to one another.
- It can be difficult to separate the individual effects of collinear variables on the response.
- Collinearity increases the uncertainty of the coefficient estimates!
- Collinearity \Rightarrow increased standard error $\Rightarrow t = \hat{\beta}_j / SE(\hat{\beta}_j) \downarrow$
 \Rightarrow failure to detect $\beta_j \neq 0 \Rightarrow$ power of hypothesis test (probability of correctly detecting non-zero coefficient) decreases.

- You can detect collinearity between 2 variables by detecting inspecting the correlation matrix of the predictors.
- But, there is multicollinearity: when 3 or more variables are correlated even if all individual pairs are not.

Instead, we have a variance inflation factor (VIF). VIF is the variance of $\hat{\beta}_j$ when fitting the full model divided by the variance of $\hat{\beta}_j$ if fit on its own.

$1 \leq \text{VIF}$

$\text{VIF} \geq 5$ or 10 is a problem

An expression for the VIF is

$$\text{VIF}(\hat{\beta}_j) = \frac{1}{1 - R^2_{X_j|X_{-j}}} \quad \text{on. col.}$$

where $R^2_{X_j|X_{-j}}$ is the R^2 from a regression of X_j onto all of X_{-j} the other predictors.

If $R^2_{X_j|X_{-j}}$ is close to 1, then collinearity is present, so VIF will be large.

Solutions:

- 1) Drop one of the problematic variables
- 2) Combine the collinear predictors into a new variable
(you could sum them or something)

So, for the residuals, you could test the hypothesis that the best fit line to the residuals is 0, and if you can't reject, then proceed.

3.5: Comparison of Linear Regression w/ KNN Regression

KNN regression is a ~~nonparametric~~ nonparametric regression algorithm.

Given a value K and a prediction point x_0 , KNN reg finds the K nearest predictor points (x_i) and then estimates

$$\hat{f}(x_0) = \frac{1}{K} \sum_{x_i \in N(x_0)} y_i \quad \begin{matrix} \text{Choose } K \text{ according to bias} \\ \text{-variance tradeoff} \end{matrix}$$

- Parametric forms outperform nonparametric if the parametric form is close to the true f.
- Note that in higher dimensions, KNNeg often performs worse than linear regression.
- This is because higher dimension effectively results in a decrease in sample size, the curse of dimensionality, where neighbor every point becomes more distant from its neighbors

Go back and do the lab.

Ch. 4: Classification

• "Qualitative" = "Categorical"

This chapter discusses 3 classifiers: logistic regression, linear discriminant analysis, and KNN

Later chapters discuss generalized additive models, trees, random forests, boosting, and support vector machines.

4.1: An Overview of Classification

Training observations $\{(x_i, y_i)\}_{i=1}^n$

4.1 Why not linear regression?

ex) Classify medical condition

$$y = \begin{cases} 1 & \text{if stroke} \\ 2 & \text{if drug overdose} \\ 3 & \text{if epileptic seizure} \end{cases}$$

4.3: Logistic Regression

- Models the probability that Y belongs to a particular category

4.3.1: The Logistic Model

$$(4.2) \quad p(X) = \frac{e^{\beta_0 + \beta_1 X}}{1 + e^{\beta_0 + \beta_1 X}}, \text{ the logistic function is } \frac{e^x}{1 + e^x}.$$

To fit the model, we use maximum likelihood.

One can manipulate (4.2) to get

$$(4.3) \quad \frac{p(X)}{1 - p(X)} = e^{\beta_0 + \beta_1 X}$$

odds (Note that odds are better for horse betting).

Note that an increase β_1 are and in X will increase the odds by a factor of e^{β_1} .

4.3.2: Estimating the Regression Coefficients

• Maximum likelihood has good statistical properties, so this is the preferred way of estimating the coefficients.

• Basic intuition: choose β_0, β_1 so that $\hat{p}(x_i)$ matches the observed class for x_i as closely as possible for all i .

• Likelihood function

$$l(\beta_0, \beta_1) = \prod_{i:y_i=1} p(x_i) \prod_{i:y_i=0} (1 - p(x_i))$$

Software can do it.

4.3.3: Making Predictions

- Just plug in the predictor values
- You can use dummy variables the same way as with linear regression.

4.3.4: Multiple Logistic Regression

$$\log\left(\frac{p(X)}{1-p(X)}\right) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$$

$$\Leftrightarrow p(X) = \frac{e^{\sum_{i=0}^p \beta_i X_i}}{1 + e^{\sum_{i=0}^p \beta_i X_i}}$$

Use ML method to fit model.

Note that confounding occurs when other predictors are relevant, and so the confounding variable needs to be included in the model.

4.3.5: Logistic Regression for >2 response classes

Could model

$$P(Y=1|X) \text{ & } P(Y=2|X) \text{ and then}$$

$$P(Y=3|X) = 1 - P(Y=1|X) - P(Y=2|X)$$

But ~~linear~~ discriminant analysis is more popular for this problem.

4.4: Linear Discriminant Analysis

- Approach: Model the distribution of the predictors X in each of the response classes. Then, use Bayes' theorem to flip these around into estimates for $P(Y=k|X=x)$.

Why use another method?

1) When classes are well separated, the parameter estimates in logistic regression are unstable. Not so for LDA

2) If n is small and the distn of X is approx normal in each class, then LDA again more stable.

LDA "works better for multiple classes."

4.4.1) Using Bayes' Theorem for Classification

We want to classify an obs into one of $K \geq 2$ classes.

Let π_k represent the overall prior probability that a randomly chosen observation comes from the k -th class. The prior is the probability that the observation comes from class k .

Let $f_k(x) = P(X=x | Y=k)$ denote the density function of X for an observation coming from class k .

Then Bayes' theorem states that

$$P(Y=k | X=x) = \frac{\pi_k f_k(x)}{\sum_{l=1}^K \pi_l f_l(x)}$$

- π_k is the fraction of training data belonging to class k
- f_k is harder to estimate unless you assume a simple form.

4.4.2: Linear Discriminant Analysis for $p=1$

• We have one predictor.

• We want to estimate $p_k(x)$ and class the function into the most likely class.

Assumptions

- 1) $f_k(x)$ is a normal density: $f_k(x) = [\sqrt{2\pi}\sigma_k]^{-1} \exp(-\frac{1}{2\sigma_k^2}(x-\mu_k)^2)$
- 2) For now, $\sigma_1^2 = \dots = \sigma_K^2 =: \sigma^2$

If we plug this into the estimator, we get

$$p_k(x) = \frac{\pi_k \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(x-\mu_k)^2\right)}{\sum_{l=1}^K \pi_l \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(x-\mu_l)^2\right)}$$

$$\log(p_k(x)) = \log(\pi_k \exp)$$

$$= \log(\pi_k) - \frac{(x-\mu_k)^2}{2\sigma^2}$$

Taking the log of both sides and rearranging, we can show that this is equivalent to assigning x to the class for which

$$f_k(x) := x \frac{\mu_k}{\sigma^2} - \frac{\mu_k^2}{2\sigma^2} + \log(\pi_k) \quad (\text{go through } x)$$

is largest

LDA estimates the Bayes classifier w/

$$\hat{\pi}_k = n_k/n$$

$$\hat{\mu}_k = \frac{1}{n_k} \sum_{i:y_i=k} x_i \quad \hat{\sigma}^2 = \frac{1}{n-K} \sum_{k=1}^K \sum_{i:y_i=k} (x_i - \hat{\mu}_k)^2$$

So, based on the scratch work, this is the simplification that makes the discriminant function linear.

4.4.3: Linear Discriminant Analysis for p>1

Assume $X = (X_1, \dots, X_p)$ is drawn from a multivariate Gaussian distribution with class specific mean vectors and a common covariance matrix.

The Multivariate Gaussian distribution

• Each predictor is μ -distributed, w/ some correlation.

We write this as $X \sim N(\mu, \Sigma)$, where Σ is the covariance matrix.

$$f(x) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

In the case of