S&DS 242/542: Theory of Statistics Lecture 21: Models with covariates and response Predictive modeling

Predicting income from education and seniority (ISLR)

For n individuals, we observe their:

- Annual income
- Years of education
- Seniority

We wish to build a predictive model for income based on years of education and seniority.

Here income is the *output variable*, *response*, or *dependent variable*. We will denote these by Y_1, \ldots, Y_n .

Years of education and seniority are the *input variables*, *covariates*, *predictors*, or *independent variables*. We will denote these by X_1, \ldots, X_n . Each X_i is the tuple of predictors for the i^{th} individual.

We will call $(X_1, Y_1), \ldots, (X_n, Y_n)$ our training data.

Predicting income from education and seniority (ISLR)



Goal: Learn a function $f(\cdot)$ for which $Y \approx f(X)$

Prediction and inference

There are two types of reasons why we may wish to learn the relationship $Y \approx f(X)$.

Prediction:

• Given a new input X_{test} , predict the response Y_{test}

• $f(\cdot)$ is a black box, and we may not be interested in its form Inference:

Understand how X is associated to Y: What is the form of the relationship between each predictor and the response? Which predictors actually influence the response? How would perturbing a predictor affect the response?

• We may want $f(\cdot)$ to have an interpretable and simple form Both goals require learning a "true" relationship between X and Y, rather than some $f(\cdot)$ that only fits well on our training data.

Formulations of the data model

Suppose $(X_1, Y_1), \ldots, (X_n, Y_n), (X_{test}, Y_{test})$ are i.i.d., and let (X, Y) represent the distribution of a single observation.

There are different models we may posit for the data, expressing different assumptions about the distribution of (X, Y):

The distribution of (X, Y) is arbitrary, and we wish to learn the regression function

$$f(x) = \mathbb{E}[Y \mid X = x]$$

Writing $Y = f(X) + \varepsilon$, this means that conditional on X = x, ε has mean 0 but its distribution may otherwise depend on x.

The distribution of X is arbitrary, and

$$Y = f(X) + \varepsilon$$

where ε has mean 0 and is independent of X.

Formulations of the data model

The distribution of X is arbitrary, and

$$Y = f(X \mid \beta) + \varepsilon$$

where, in additional to the previous assumptions for ε , the function $f(\cdot | \beta)$ has a parametric form with parameter β .

For example, if $X = (x_1, ..., x_p)$ has p quantitative predictors, we may consider a linear model

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \ldots + \beta_p x_p + \varepsilon$$

with parameters $\beta = (\beta_0, \beta_1, \beta_2, \dots, \beta_p) \in \mathbb{R}^{p+1}$.

We will oftentimes make a distinction between the model for (X, Y) that we choose to fit, vs. the actual distribution of the data.

Irreducible vs. reducible prediction error

Suppose we predict Y by a function g(X). Fixing a test point x, we may be interested in the mean-squared prediction error

$$\mathbb{E}[(Y_{\text{test}} - g(x))^2 \mid X_{\text{test}} = x]$$

or its average over the distribution of X_{test} .

Suppose
$$Y = f(X) + \varepsilon$$
 where ε is mean-0 and independent of X.
We have:

$$\mathbb{E}\left[\left(Y_{test} - g(x)\right)^2 \middle| X_{test} = x\right) \\
 = \mathbb{E}\left[\left(\left\{I(x) + \varepsilon - g(x)\right)^2 \middle| X_{test} = x\right) \\
 = \mathbb{E}\left[\left(\left\{I(x) - g(x)\right)^2\right] + \mathbb{E}\left[\sum_{i=0}^{\infty} \left(I(x) - g(x)\right) + \varepsilon\right] + \mathbb{E}\left[\varepsilon^2\right] \\
 = \left(\left\{I(x) - g(x)\right)^2 + V_{ar}\left[\varepsilon\right] \right]$$

Irreducible vs. reducible prediction error

Thus

$$\mathbb{E}[(Y_{\text{test}} - g(x))^2 \mid X_{\text{test}} = x] = \underbrace{(f(x) - g(x))^2}_{\text{reducible error}} + \underbrace{\operatorname{Var}[\varepsilon]}_{\text{irreducible error}}$$

- We would incur the irreducible error even if we had perfect knowledge of the true function f(·). This error is non-zero if Y cannot be perfectly predicted given X, e.g. if there are unmeasured variables that also influence Y.
- ► The reducible error is the error that we may hope to minimize by learning an accurate prediction model g(·). In the best case, g(x) = f(x) and the reducible error is 0.

Bias-variance decomposition of reducible error

In practice, we estimate f(x) by a function $\hat{f}(x)$ that is random and depends on our training data $(X_1, Y_1), \ldots, (X_n, Y_n)$.

Still fixing the test point x, let us define the *bias* and *variance* of the prediction at x by

$$\mathsf{Bias} = \mathbb{E}[\hat{f}(x)] - f(x)$$
Variance $= \mathsf{Var}[\hat{f}(x)]$

where $\mathbb{E}[\cdot]$ and $Var[\cdot]$ are over $(X_1, Y_1), \ldots, (X_n, Y_n)$.

Then $\mathbb{E}[(\text{reducible error})^2] = \mathbb{E}\left[\left(\widehat{f}(x) - \widehat{f}(x)\right)^2\right]$ = $\mathbb{E}\left[\left(\widehat{f}(x) - \mathbb{E}\widehat{f}(y) + \mathbb{E}\widehat{f}(x) - \widehat{f}(x)\right)^2\right]$ = $\left(\widehat{f}(x) - \mathbb{E}\widehat{f}(x)\right)^2 + V_{n-1}\left[\widehat{f}(y)\right]$

Three-fold decomposition of prediction error

Thus the prediction mean-squared-error at a test point x is

$$\begin{split} \mathbb{E}[(Y_{\text{test}} - \hat{f}(x))^2 \mid X_{\text{test}} = x] \\ &= \underbrace{(\mathbb{E}[\hat{f}(x)] - f(x))^2}_{\text{squared bias}} + \underbrace{\operatorname{Var}[\hat{f}(x)]}_{\text{variance}} + \underbrace{\operatorname{Var}[\varepsilon]}_{\text{irreducible error}} \end{split}$$

- Variance refers to the variability of the our prediction f(x) under different possible realizations of our training data.
 Large variance may indicate that we have *overfitted* to the training data.
- Bias refers to a difference between the true regression function f(x) and our prediction f(x) averaged over the randomness of the training data. Large bias may indicate that the true f(x) is far from the class of functions represented by our model.

Bias-variance tradeoff in prediction error



Black curve: True regression function f(x)Orange, blue, green curves: Three different estimates $\hat{f}(x)$

Considerations for model complexity

The choice between fitting a simple model vs. a complex model may be determined by two trade-offs:

- Prediction vs. inference: To what extent do we care about the predictive accuracy of the model, and to what extent do we care about its interpretability and the meaning it encodes for the relationship between X and Y?
- Bias vs. variance: Even if our primary goal is predictive accuracy, what would be the right model complexity that balances bias and variance of our model predictions? This depends on the amount of available data and the complexity of the true regression function f(x).

Simple linear regression

Simple linear regression

Consider data $(X_1, Y_1), \ldots, (X_n, Y_n)$ where $X \in \mathbb{R}$ is a single quantitative predictor. The simple linear regression model is

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

with parameters β_0, β_1 , representing the slope and intercept of the regression line.



The population regression line

Restricting ourselves to this linear model, the parameters we wish to learn may be the *population least-squares coefficients*

$$eta_0^*,eta_1^* = rgmin_{eta_0,eta_1}\mathbb{E}[(eta_0+eta_1X-Y)^2]$$

These parameters minimize the prediction mean-squared-error over the true distribution of data.

These have an explicit form:

$$O = \frac{2}{3\rho_0} E\left[\left(\beta_0 + \beta_1 X - Y\right)^2\right]$$

= $E\left[Z\left(\rho_0 + \beta_1 X - Y\right)\right] = 2\rho_0 + 2\rho_1 EX - ZEY$
$$O = \frac{2}{3\rho_1} E\left[\left(\beta_0 + \beta_1 X - Y\right)^2\right]$$

= $E\left[ZX\left(\beta_0 + \beta_1 X - Y\right)\right] = 2\rho_0 EX + 2\rho_1 EE[X^2] - ZE[XY]_{13}$

The population regression line

=) β. EX + β. (EX)²= EX EY BOBX + BI ELX) = ELXY) =) $\beta_i^* = \frac{\mathbb{E}[XY] - \mathbb{E}X \cdot \mathbb{E}Y}{\mathbb{E}[X^2] - (\mathbb{E}X)^2} = \frac{Cov[X, Y]}{Var[X]}$

β. = ΙΕΥ- β. ΕΧ

Correlation and regression

Suppose X, Y both have mean 0 and variance 1, and their correlation is $Corr(X, Y) = \rho$. Then

$$\beta_0^* = 0, \quad \beta_1^* = \rho,$$

and the population least-squares regression line is $Y = \rho X$.

This is *not* symmetric between X and Y. If we were predicting X from Y, then the least-squares line would be $X = \rho Y$, i.e. $Y = \frac{X}{\rho}$.

Interpretation: Let (X, Y) = (height of father, height of son). Given that the father is 1 std. dev. taller than average, the son is only ρ std. dev. taller than average in expectation.

(Similarly, given that the son is 1 std. dev. taller than average, the father is only ρ std. dev. taller than average in expectation.)

Correlation and regression

Example: Suppose (X, Y) are bivariate normal w/ correlation 0.5



Black: Y = XRed: Least-squares regression line for predicting Y from X Blue: Least-squares regression line for predicting X from Y

Least squares estimation

We may fit β_0, β_1 by minimizing the average squared error on our training data,

$$(\hat{\beta}_0, \hat{\beta}_1) = \operatorname*{arg\,min}_{\beta_0, \beta_1} \frac{1}{n} \sum_{i=1}^n (\beta_0 + \beta_1 X_i - Y_i)^2$$

These are the **least-squares estimators** for β_0, β_1 .

The same derivation as for the population coefficients β_0^*, β_1^* shows that these estimators have the closed-form expressions

$$\hat{\beta}_{1} = \frac{\sum_{i=1}^{n} (X_{i} - \bar{X})(Y_{i} - \bar{Y})}{\sum_{i=1}^{n} (X_{i} - \bar{X})^{2}} \qquad \hat{\beta}_{0} = \bar{Y} - \hat{\beta}_{1}\bar{X}$$

A maximum likelihood interpretation of least squares

Suppose that $X_1, \ldots, X_n \stackrel{IID}{\sim} f_X$, and

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$$

where $\varepsilon_1, \ldots, \varepsilon_n \stackrel{IID}{\sim} \mathcal{N}(0, \sigma^2)$ are normal errors independent of X_1, \ldots, X_n . Then $(\hat{\beta}_0, \hat{\beta}_1)$ is the MLE for (β_0, β_1) : Joint dusty of (X,Y): fx(x): f= e - tor (p-1p.x-y)2 $= \int_{n} \left(\beta_{i}, \beta_{i} \right) = \sum_{i=1}^{n} \int_{y} \left[f_{x}(x_{i}) \int_{z_{x}} e^{-\frac{1}{2\sigma^{2}}} \left(\beta_{i} t_{i} \beta_{i} x_{i} - y_{i} \right)^{2} \right)$ = Z - 1 ((" + p, x, - y)" + log [fx (x:) - for] Marin's $l_n(\beta, \beta) \hookrightarrow M_n'm's f_n \sum_{i=1}^{n} (\beta, f_i, \chi, -\chi)^2$

Unbiasedness in a correctly specified model

Theorem

Suppose $(X_1, Y_1), \ldots, (X_n, Y_n)$ are *i.i.d.* and distributed according to a true model

$$Y = \beta_0^* + \beta_1^* X + \varepsilon$$

where ε has mean 0 and is independent of X. Then $\mathbb{E}[\hat{\beta}_0] = \beta_0^*$ and $\mathbb{E}[\hat{\beta}_1] = \beta_1^*$, i.e. $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1)$ is unbiased for $\beta^* = (\beta_0^*, \beta_1^*)$.

Proof:

$$\hat{\beta}_{i} = \frac{\tilde{z}_{i}(X_{i} - \bar{X})(Y_{i} - \bar{Y})}{\tilde{z}_{i}(X_{i} - \bar{X})^{2}} \qquad \text{Dende } X = (X_{y}, X_{y})$$

$$E[\hat{\beta}_{i} | X = x] = \frac{\tilde{z}_{i}(x_{i} - \bar{x})E[Y_{i} - \bar{Y} | X = x]}{\tilde{z}_{i}(x_{i} - \bar{x})E[Y_{i} - \bar{Y} | X = x]}$$

Unbiasedness in a correctly specified model

We have:
$$\mathbb{E}[Y_{\cdot}|X=x] = \beta_{0}^{*} + \beta_{1}^{*} \times \mathbb{I}$$

 $\mathbb{E}[\overline{Y}|X=x] = \beta_{0}^{*} + \beta_{1}^{*} \times \mathbb{I}$
 $\Rightarrow \mathbb{E}[\widehat{\beta}_{1}|X=x] = \frac{\widehat{z}_{1}^{*}(x_{1}-\overline{x})\cdot(x_{2}-\overline{x})\cdot\beta_{1}^{*}}{\widehat{z}_{1}^{*}(x_{2}-\overline{x})^{*}} = \beta_{1}^{*}$
 $\widehat{\beta}_{0}^{*} = \overline{Y} - \widehat{\beta}_{1}\cdot\overline{X} \Rightarrow \mathbb{E}[\widehat{\beta}_{0}|X=x]$
 $= \mathbb{E}[\overline{Y}|X=x] - \mathbb{E}[\widehat{\beta}_{1}|X=x]\cdot\overline{x}$
 $= \mathbb{E}[\overline{Y}|X=x] - \mathbb{E}[\widehat{\beta}_{1}|X=x]\cdot\overline{x}$
 $= (\beta_{0}^{*} + \beta_{1}^{*}\overline{x}) - (\beta_{1}^{*}\overline{x} = \beta_{0}^{*} \Rightarrow)\mathbb{E}[\widehat{\beta}_{0}] = \beta_{0}^{*}.$

(In fact, this only required the assumption $\mathbb{E}[\varepsilon \mid X = x] = 0$ for each $x \in \mathbb{R}$, i.e. ε has mean-0 conditional on X.)

Consistency for the population coefficients

Theorem Suppose $(X_1, Y_1), \ldots, (X_n, Y_n)$ are *i.i.d.*, and let $(\beta_0^*, \beta_1^*) = \arg \min \mathbb{E}[(\beta_0 + \beta_1 X - Y)^2]$ β_0, β_1 Then $(\hat{\beta}_0, \hat{\beta}_1) \to (\beta_0^*, \beta_1^*)$ in probability as $n \to \infty$, i.e. $\hat{\beta}$ is consistent for β^* . Proof: $\widehat{\beta}_{1} = \frac{1}{4} \sum_{i=1}^{2} (X_{i} - \overline{X})(Y_{i} - \overline{Y}) = \frac{1}{4} \sum_{i=1}^{2} X_{i} Y_{i} - \overline{X} \overline{Y}$ $\frac{1}{4} \sum_{i=1}^{2} (X_{i} - \overline{X})^{2} = \frac{1}{4} \sum_{i=1}^{2} X_{i} Y_{i} - \overline{X} \overline{Y}$ As no 20, by UN $\cdot \downarrow \bar{z} \times \gamma \rightarrow E[X\gamma]$ ·XJEX $\cdot \neq \tilde{z} \chi^2 \rightarrow E[\chi']$ ·YAEY

Consistency for the population coefficients

$$\Rightarrow \hat{\beta}_i \rightarrow \frac{IE[XY] - IEX \cdot IEY}{E[XY] - (IEX)^2} = \beta_i^*$$

$$\hat{\beta}_{o} = \overline{\nabla} - \hat{\beta}_{i} \cdot \overline{X}$$

$$\longrightarrow EY - \beta_{i}^{+} \cdot EX = \beta_{o}^{*}$$

Decomposition of prediction error

Recall our decomposition of prediction mean-squared-error at a fixed test point *x*:

$$\begin{split} \mathbb{E}[(Y_{\text{test}} - \hat{f}(x))^2 \mid X_{\text{test}} = x] \\ &= \underbrace{(\mathbb{E}[\hat{f}(x)] - f(x))^2}_{\text{squared bias}} + \underbrace{\operatorname{Var}[\hat{f}(x)]}_{\text{variance}} + \underbrace{\operatorname{Var}[\varepsilon]}_{\text{irreducible error}} \end{split}$$

Suppose the true model is $Y = \beta_0^* + \beta_1^* X + \varepsilon$. Then

- The true regression function is $f(x) = \beta_0^* + \beta_1^* x$.
- ▶ The estimated regression function is $\hat{f}(x) = \hat{\beta}_0 + \hat{\beta}_1 x$. Since $\hat{\beta}$ is unbiased, $\mathbb{E}[\hat{f}(x)] = f(x)$ so the squared bias is 0.

As n→∞, since β̂→ β* in probability, we have also f̂(x) → f(x) in probability. Then we expect Var[f̂(x)] → 0.
 Assuming that the irreducible error Var[ε] > 0 is fixed for all n, the prediction error is dominated by this irreducible error for large n.

Decomposition of prediction error

More generally, suppose $Y = f(X) + \varepsilon$ where ε is mean-0 and independent of X. Let

$$(\beta_0^*, \beta_1^*) = \operatorname*{arg\,min}_{\beta_0, \beta_1} \mathbb{E}[(\beta_0 + \beta_1 X - Y)^2]$$

Then, as $n \to \infty$,

Since $\hat{\beta} \to \beta^*$ in probability, also $\hat{f}(x) \to \beta_0^* + \beta_1^* x$ in probability. Then we expect $\mathbb{E}[\hat{f}(x)] \to \beta_0^* + \beta_1^* x$, so the squared bias approaches $(\beta_0^* + \beta_1^* x - f(x))^2$.

• We still expect
$$Var[\hat{f}(x)] \rightarrow 0$$
.

Thus for large n, the prediction error is approximately

$$\mathbb{E}[(Y_{\text{test}} - \hat{f}(x))^2 \mid X_{\text{test}} = x] \approx (\beta_0^* + \beta_1^* x - f(x))^2 + \mathsf{Var}[\varepsilon]$$

the sum of the irreducible error and the squared approximation error of f(x) by the best linear predictor.