Linear Regression as Prediction

36-462/36-662, Spring 2022

20 January 2022

Context

- We want to find patterns in our data which will let us predict one variable from another (or from more than one variable)
- The oldest and most widely-used prediction method is linear regression
- We're going to review linear regression, as a prediction method
- Without a lot of the mythology from 401
- But with an eye to more powerful methods

Optimal prediction in general

- We want to predict a *numerical* random variable Y
- We want a one-number ("point") prediction, not a range or a distribution
- We say how good our guess is by expected squared error

Optimal prediction in general (cont'd.)

What's the best *constant* guess for a random variable Y?

$$\mu = \operatorname*{argmin}_{m \in \mathbb{R}} \mathbb{E}\left[(Y - m)^2 \right]$$
(1)

$$= \operatorname{argmin}_{m} \operatorname{Var}\left[(Y-m)\right] + \left(\mathbb{E}\left[Y-m\right]\right)^{2}$$
(2)

$$= \operatorname{argmin} \operatorname{Var} [Y] + (\mathbb{E} [Y] - m)^2$$
(3)

$$= \operatorname*{argmin}_{m} \left(\mathbb{E}\left[Y\right] - m \right)^2 \tag{4}$$

$$= \mathbb{E}[Y] \tag{5}$$

(Because: $\mathbb{E}[Z^2] = \operatorname{Var}[Z] + (\mathbb{E}[Z])^2$, always)

Optimal prediction in general (cont'd.)

What's the best *function* of X to guess for Y?

$$\mu = \underset{m:\mathcal{X} \to \mathbb{R}}{\operatorname{argmin}} \mathbb{E}\left[(Y - m(X))^2 \right]$$
(6)

$$= \operatorname{argmin}_{m} \mathbb{E}\left[\mathbb{E}\left[(Y - m(X))^{2} | X\right]\right]$$
(7)

For each x, best m(x) is $\mathbb{E}[Y|X = x]$ (by previous slide)

$$\mu(x) = \mathbb{E}\left[Y|X=x\right]$$

Optimal prediction in general (cont'd.)

Learning arbitrary functions is hard! Who knows what the right function might be? What if we *decide* to make our predictions linear?

Optimal linear prediction with univariate predictor

Our prediction will be of the form

$$m(x) = a + bx$$

and we want the best a, b

Optimal linear prediction, univariate case

$$(\alpha, \beta) = \underset{a,b}{\operatorname{argmin}} \mathbb{E}\left[(Y - (a + bX))^2 \right]$$

Expand out that expectation, then take derivatives and set them to 0

The intercept

$$\mathbb{E}\left[(Y - (a + bX))^2\right] = \mathbb{E}\left[Y^2\right] - 2\mathbb{E}\left[Y(a + bX)\right] + \mathbb{E}\left[(a + bX)^2\right]$$
(8)

$$= \mathbb{E}\left[Y^2\right] - 2a\mathbb{E}\left[Y\right] - 2b\mathbb{E}\left[YX\right] + \tag{9}$$

$$a^{2} + 2ab\mathbb{E}\left[X\right] + b^{2}\mathbb{E}\left[X^{2}\right] \tag{10}$$

$$\frac{\partial}{\partial a} \mathbb{E}\left[(Y - (a + bX))^2 \right] \Big|_{a=\alpha,b=\beta} = -2\mathbb{E}\left[Y \right] + 2\alpha + 2\beta \mathbb{E}\left[X \right] = 0$$
(11)

$$\alpha = \mathbb{E}[Y] - \beta \mathbb{E}[X] \tag{12}$$

 \therefore optimal linear predictor $m(X) = \alpha + \beta X$ looks like

$$m(X) = \alpha + \beta X \tag{13}$$

$$= \mathbb{E}[Y] - \beta \mathbb{E}[X] + \beta X \tag{14}$$

$$= \mathbb{E}[Y] + \beta(X - \mathbb{E}[X]) \tag{15}$$

The optimal linear predictor only cares about how far X is from its expectation $\mathbb{E}[X]$ And when $X = \mathbb{E}[X]$, we will always predict $\mathbb{E}[Y]$

The slope

$$\frac{\partial}{\partial b} \mathbb{E}\left[(Y - (a + bX))^2 \right] \Big|_{a=\alpha, b=\beta} = -2\mathbb{E}\left[YX \right] + 2\alpha \mathbb{E}\left[X \right] + 2\beta \mathbb{E}\left[X^2 \right] = 0$$
(16)

$$0 = -\mathbb{E}[YX] + (\mathbb{E}[Y] - \beta \mathbb{E}[X])\mathbb{E}[X] + \beta \mathbb{E}[X^2]$$
(17)

$$0 = \mathbb{E}[Y]\mathbb{E}[X] - \mathbb{E}[YX] + \beta(\mathbb{E}[X^2] - \mathbb{E}[X]^2)$$
(18)

$$0 = -\operatorname{Cov}[Y, X] + \beta \operatorname{Var}[X]$$

$$(19)$$

$$\beta = \frac{\operatorname{Cov}[Y, X]}{\operatorname{Var}[X]}$$
(20)

- If we replace X with $X' = X \mathbb{E}[X]$, β doesn't change
- If we replace Y with $Y' = Y \mathbb{E}[Y]$, β doesn't change
- \therefore centering the variables doesn't change the slope

The optimal linear predictor of Y from X

The optimal linear predictor of Y from a single X is *always*

$$\alpha + \beta X = \mathbb{E}\left[Y\right] + \left(\frac{\operatorname{Cov}\left[X,Y\right]}{\operatorname{Var}\left[X\right]}\right)\left(X - \mathbb{E}\left[X\right]\right)$$

What did we not assume?

- That the true relationship between Y and X is linear
- That anything is Gaussian
- That anything has constant variance
- That anything is independent or even uncorrelated

NONE OF THAT MATTERS for the optimal linear predictor

The prediction errors average out to zero

$$\mathbb{E}[Y - m(X)] = \mathbb{E}[Y - (\mathbb{E}[Y] + \beta(X - \mathbb{E}[X]))]$$
(21)
$$= \mathbb{E}[Y] - \mathbb{E}[Y] - \beta(\mathbb{E}[X] - \mathbb{E}[X]) = 0$$
(22)

- If they didn't average to zero, we'd adjust the coefficients until they did
- Important: In general, $\mathbb{E}[Y m(X)|X] \neq 0$

The prediction errors are uncorrelated with X

$$\operatorname{Cov}\left[X, Y - m(X)\right] = \mathbb{E}\left[X(Y - m(X))\right] \text{ (by previous slide)}$$
(23)

$$= \mathbb{E}\left[X(Y - \mathbb{E}[Y] - \frac{\operatorname{Cov}[Y, X]}{\operatorname{Var}[X]}(X - \mathbb{E}[X]))\right]$$
(24)

$$= \mathbb{E}\left[XY - X\mathbb{E}\left[Y\right] - \frac{\operatorname{Cov}\left[Y, X\right]}{\operatorname{Var}\left[X\right]}(X^2) + \frac{\operatorname{Cov}\left[Y, X\right]}{\operatorname{Var}\left[X\right]}(X\mathbb{E}\left[X\right])\right]$$
(25)

$$= \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y] - \frac{\operatorname{Cov}[Y,X]}{\operatorname{Var}[X]}\mathbb{E}[X^{2}] + \frac{\operatorname{Cov}[Y,X]}{\operatorname{Var}[X]}(\mathbb{E}[X])^{2}$$
(26)

$$= \operatorname{Cov} [X, Y] - \frac{\operatorname{Cov} [Y, X]}{\operatorname{Var} [X]} (\operatorname{Var} [X])$$
(27)

$$= 0$$
 (28)

• If the errors weren't uncorrelated with X, we'd adjust the coefficients until they were

The prediction errors are uncorrelated with X

Alternate take:

$$\operatorname{Cov} [X, Y - m(X)] = \operatorname{Cov} [X, Y] - \operatorname{Cov} [X, \alpha + \beta X]$$
(29)

$$= \operatorname{Cov}[Y, X] - \operatorname{Cov}[X, \beta X]$$
(30)

$$= \operatorname{Cov}[Y, X] - \beta \operatorname{Cov}[X, X]$$
(31)

$$= \operatorname{Cov}[Y, X] - \beta \operatorname{Var}[X]$$
(32)

$$= \operatorname{Cov}[Y, X] - \operatorname{Cov}[Y, X] = 0 \tag{33}$$

How big are the prediction errors?

$$\operatorname{Var}\left[Y - m(X)\right] = \operatorname{Var}\left[Y - \alpha - \beta X\right]$$
(34)

$$= \operatorname{Var}\left[Y - \beta X\right] \tag{35}$$

(36)

After-class exercise: Reduce this to an expression involving only $\operatorname{Var}[Y]$, $\operatorname{Var}[X]$ and $\operatorname{Cov}[Y, X]$; if you get the right answer you should see that it's $< \operatorname{Var}[Y]$ unless $\operatorname{Cov}[Y, X] = 0$

 \Rightarrow Optimal linear predictor is almost always better than nothing...

Multivariate case

- We try to predict Y from a whole bunch of variables
- Bundle those predictor variables into \vec{X}
- We need a vector of slope coefficients $\vec{\beta}$
- Solution:

$$m(\vec{X}) = \alpha + \vec{\beta} \cdot \vec{X} = \mathbb{E}\left[Y\right] + \operatorname{Var}\left[\vec{X}\right]^{-1} \operatorname{Cov}\left[\vec{X}, Y\right] (\vec{X} - \mathbb{E}\left[\vec{X}\right])$$

and

$$\operatorname{Var}\left[Y - m(\vec{X})\right] = \operatorname{Var}\left[Y\right] - \operatorname{Cov}\left[Y, \vec{X}\right]^{T} \operatorname{Var}\left[\vec{X}\right]^{-1} \operatorname{Cov}\left[Y, \vec{X}\right]$$

(Gory details in the back-up slides)

What we don't assume, again

- That the linear predictor is correct
- That anything is Gaussian
- Anything about the distributions of Y or \vec{X}

Estimation: Data, not the full distribution

- Var [X], Cov [Y, X], $\mathbb{E}\left[(Y (a + bX))^2\right]$ all involve the *true* distribution, which we don't know
- : We can't just calculate β
- What we do know is $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ - That is, *n* pairs of predictor and outcome

Estimation: Ordinary Least Squares (OLS)

Set up the mean squared error, and minimize it:

$$MSE(a,b) = \frac{1}{n} \sum_{i=1}^{n} (y_i - (a + bx_i))^2$$
(37)

$$(\hat{\alpha}, \hat{\beta}) \equiv \operatorname*{argmin}_{a,b} MSE(a, b)$$
 (38)

Do the calculus:

$$\frac{\partial MSE}{\partial a} = \frac{1}{n} \sum_{i=1}^{n} -2(y_i - (a + bx_i))$$
(39)

$$\frac{\partial MSE}{\partial b} = \frac{1}{n} \sum_{i=1}^{n} -2(y_i - (a + bx_i))x_i$$

$$\tag{40}$$

$$\hat{\alpha} = \overline{y} - \hat{\beta}\overline{x} \tag{41}$$

$$\hat{\beta} = \frac{\overline{yx} - \overline{yx}}{\overline{x^2}} = \frac{\operatorname{Cov}\left[Y, X\right]}{\widehat{\operatorname{Var}\left[X\right]}}$$
(42)

(with \overline{x} = sample mean of x, etc.)

Optimum vs. estimate (I)

• Optimal linear model *versus* Linear model estimated by OLS:

$$\alpha + \beta x = \mathbb{E}[Y] + \left(\frac{\operatorname{Cov}[X,Y]}{\operatorname{Var}[X]}\right)(x - \mathbb{E}[X])$$
(43)

$$\hat{\alpha} + \hat{\beta}x = \overline{y} + \frac{\widetilde{\operatorname{Cov}}[Y, \overline{X}]}{\widetilde{\operatorname{Var}}[X]}(x - \overline{x})$$
(44)

• These are not the same, but we can hope they're close, and increasingly close with more data

When does OLS/plug-in work?

- "Work" = converge on the optimal linear predictor
- Jointly sufficient conditions:
- 1. Sample means converge on expectation values
- 2. Sample covariances converge on true covariance
- 3. Sample variances converge on true, invertible variance
- Then by continuity OLS coefficients converge on true β
- None of this requires that the linear model is right, that anything is Gaussian, etc.

Optimum vs. estimate (II)

• A little more subtle: for every fixed (a, b), by the law of large numbers,

$$MSE(a,b) \to \mathbb{E}\left[(Y - (a + bX))^2 \right]$$

- Mean squared error \rightarrow expected squared error
- This is *almost* enough to ensure that minimize MSE will converge on the true optimum
- We'll come back to this, because this approach generalizes better to other models and loss functions

What do the estimates look like?

- Bundle all the regressor values into an $n \times (p+1)$ matrix **x**, with an extra column of 1s
- Bundle all the regressand values into an $n \times 1$ matrix **y**
- Then the $(p+1) \times 1$ matrix of least-squares coefficients $\hat{\beta}$ is

$$\hat{\beta} = (\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{y}$$

(as you learned in linear models)

• This is really

$$\hat{\beta} = \left(\frac{1}{n}\mathbf{x}^T\mathbf{x}\right)^{-1}\left(\frac{1}{n}\mathbf{x}^T\mathbf{y}\right)$$
(45)

= (sample variance matrix of regressors)⁻¹(sample covariances between regressors and respo(46))

What do the predictions look like?

• The prediction at an arbitrary point $\vec{x} = (x_1, x_2, \dots, x_p)$ is

$$\hat{m}(\vec{x}) = \begin{bmatrix} 1 & x_1 & x_2 & \dots & x_p \end{bmatrix} \hat{\beta} = \begin{bmatrix} 1 & x_1 & x_2 & \dots & x_p \end{bmatrix} ((\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{y})$$

• The $n \times 1$ matrix of fitted values, at the training points, is

$$\mathbf{\hat{m}} = \mathbf{x}(\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{y}$$

• The matrix $\mathbf{h} \equiv \mathbf{x} (\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T$ is the weight matrix, influence matrix or hat matrix and says how much y_i matters to the prediction of y_i

Fitted values and other predictions are weighted sums of the observations

$$\hat{m}(\vec{x}) = \vec{x}\hat{\beta} = \vec{x}(\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{y}$$
(47)

$$\hat{\mathbf{m}} = \mathbf{x} (\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{y}$$
(48)

- Every prediction we make is linear in y, it's a weighted sum of all the observed values of the response
- How much weight do we put on y_i when we're predicting at \vec{x} ?

Explicit form of the weights for OLS

• For univariate OLS, if we're trying to predict at the point x and built the model with data $(x_1, y_1), \ldots, (x_n, y_n)$, we end up with:

$$\hat{m}(x) = \bar{y} + \hat{\beta}_1(x - \bar{x}) \tag{49}$$

$$= \sum_{i=1}^{n} \frac{1}{n} \left(1 + \frac{(x-\overline{x})(x_i-\overline{x})}{\hat{\sigma}_x^2} \right) y_i \tag{50}$$

- We give more weight to y_i if x_i is far from \overline{x}
- and we give more weight to y_i if x is far from \overline{x}
- Weights don't care about $x x_i$, which is stupid, but also required to get a straight line
- Same idea, with more algebra, for multiple regressors

Generalizing: linear smoothers

• Linear regression is a special case of a linear smoother

$$\widehat{\mu}(\vec{x}) = \sum_{i=1}^{n} w(\vec{x}, \vec{x}_i) y_i$$

- Notice: Linear in the y's, not in \vec{x}
- General idea: predict at \vec{x} by finding similar \vec{x}_i 's and averaging their y_i 's
- Different choices of w are different ideas about "similar" and about weighted averaging
- Most prediction methods used in practice are linear smoothers
 - Nearest neighbors
 - Trees and forests
 - Kernels
 - Neural networks, a.k.a. "deep learning", a.k.a. "artificial intelligence"

What about the rest of your linear models course?

- Say you want to test whether $\beta_{37} = 1$, or to give a confidence set for (α, β_{12})
- That's hard to do *explicitly* without more assumptions
- Generally, you need to know the sampling distribution of $\hat{\beta}$

What about the rest of your linear models course? (cont'd)

- The usual ("401") assumptions give nice formulas for the sampling distribution of parameter estimators, and so for parameter inference:
- 1. The true regression function is exactly linear.

- 2. $Y = \alpha + \vec{X} \cdot \vec{\beta} + \epsilon$ where ϵ is *independent* of x.
- 3. ϵ is independent across observations.
- 4. $\epsilon \sim \mathcal{N}(0, \sigma^2)$.
- Assuming (1)–(4), the usual formulas (embedded in R) are exactly right

 Gaussian distribution for β̂, centered at true β
- Assuming (1)–(3), non-Gaussian noise, and the influence of any one observation $\rightarrow 0$, then the usual formulas hold as $n \rightarrow \infty$
 - Get back the Gaussian distribution by a (complicated) central limit theorem
- Do not use those formulas without checking those assumptions

The most important assumption to check

• If the **true** regression function is *exactly* linear, then

$$\mathbb{E}\left[Y - (\alpha + \vec{X}\beta) | \vec{X} = \vec{x}\right] = 0$$

for all \vec{x}

• So

$$\mathbb{E}\left[Y-(\hat{\alpha}+\vec{X}\hat{\beta})|\vec{X}=\vec{x}\right]\approx 0$$

• The residuals should have conditional mean 0 everywhere if the linear model is right – This is something we can check using nonlinear regression methods

Summing up

- We can always *decide* to use a linear predictor, $m(\vec{X}) = \alpha + \vec{\beta} \cdot \vec{X}$
- The optimal linear predictor of Y from \vec{X} always takes the same form:

$$m(Y) = \mathbb{E}[Y] + \operatorname{Var}\left[\vec{X}\right]^{-1} \operatorname{Cov}\left[Y, \vec{X}\right] (\vec{X} - \mathbb{E}\left[\vec{X}\right])$$

- Doing linear prediction requires finding those covariances
- We usually estimate those covariances (implicitly) by ordinary least squares
- OLS converges pretty robustly to the optimal linear predictor (not to the truth)
- The usual "401" assumptions are needed for *parameter* inference, not point prediction
- Once we use OLS, all our predictions are weighted sums of the observations
- The weights for linear regression are weird and implausible
- We're going to explore other ways of weighting
- We're also going to explore minimize other measures of prediction error, over other classes of predictors

A final thought

When you're fundraising, it's AI When you're hiring, it's ML When you're implementing, it's linear regression

• Baron Schwartz, 15 November 2017¹

 $^{^1}$ Since-deleted tweet, but see e.g. [https://twitter.com/bc238dev/status/1225150435729666048]. I've often seen the last line quoted as "it's logistic regression", which fits with computer science's emphasis on classification rather than regression, but so far as I can work out that's a later mutation.

Backup: Further reading

- See Berk (2008), chapter 1, and Hastie, Tibshirani, and Friedman (2009), sections 2.3.1 and 2.6, and chapter 3 (especially through section 3.4), for more on linear regression as a prediction method
- If you want more from the point of view taken here, see Shalizi (n.d.), chapter 2; if you want a *lot* more, see Shalizi (2015)
- If you need to do inference about the coefficients but the usual assumptions don't hold, your best bet is the bootstrap (covered in 402, and in some excellent textbooks like Davison and Hinkley (1997)); some of the complicated procedures developed in econometrics for "robust standard errors" in linear models turn out to be equivalent to simple bootstraps (Buja et al. 2014)
 - If you need prediction intervals or predictive distributions from your linear model, there are also ways of doing that using the bootstrap — see Davison and Hinkley (1997) again, or again Berk (2008), chapter 1
- The view that the optimal linear model is just $\operatorname{Var}\left[\vec{X}\right]^{-1}\operatorname{Cov}\left[Y,\vec{X}\right]$ and all we need is for sample covariance to converge, goes back to Wiener (1949)
 - A notoriously hard-to-read book, but that's because Wiener was trying to explain this in a way which made mathematical sense where instead of regressing Y on a fixed-length vector \vec{X} , you're regressing Y on a continuous function, and framing the solution in a way which could be implemented using 1940-vintage analog electronics...

Backup: The optimal regression function

• We set up the problem

$$\mu = \underset{m:\mathcal{X} \to \mathbb{R}}{\operatorname{argmin}} \mathbb{E}\left[\mathbb{E}\left[(Y - m(X))^2 | X\right]\right]$$
(51)

$$= \operatorname{argmin}_{m:\mathcal{X} \mapsto \mathbb{R}} \int dx p(x) \int dy p(y|x) (y - m(x))^2$$
(52)

- For each x, we can minimize the *inner* integral by setting $m(x) = \mathbb{E}[Y|X = x]$
- We just pasted together all of those pointwise minimizers to get μ
- We could do this *because* we were minimizing over all possible functions, so minimizing at x_1 doesn't interfere with minimizing at x_2
- This doesn't work if we limit ourselves to a particular set of functions \mathcal{M}

- Unless of course $\mu \in \mathcal{M}$, i.e., that model is well-specified

Backup: Gory details for multivariate predictors

$$m(\vec{X}) = a + \vec{b} \cdot \vec{X} \tag{53}$$

$$(\alpha, \vec{\beta}) = \underset{a, \vec{b}}{\operatorname{argmin}} \mathbb{E}\left[(Y - (a + \vec{b} \cdot \vec{X}))^2 \right]$$
(54)

$$\mathbb{E}\left[(Y - (a + \vec{b} \cdot \vec{X}))^2 \right] = \mathbb{E}\left[Y^2 \right] + a^2 + \mathbb{E}\left[(\vec{b} \cdot \vec{X})^2 \right] -2\mathbb{E}\left[Y(\vec{b} \cdot \vec{X}) \right] - 2\mathbb{E}\left[Ya \right] + 2\mathbb{E}\left[a\vec{b} \cdot \vec{X} \right]$$
(55)

$$= \mathbb{E}\left[Y^{2}\right] + a^{2} + \vec{b} \cdot \mathbb{E}\left[\vec{X} \otimes \vec{X}\right] b$$

$$(56)$$

$$-2a\mathbb{E}\left[Y\right] - 2\vec{b} \cdot \mathbb{E}\left[Y\vec{X}\right] + 2a\vec{b} \cdot \mathbb{E}\left[\vec{X}\right]$$

$$\tag{57}$$

Backup: Gory details: the intercept

Take derivative w.r.t. a, set to 0:

$$0 = -2\mathbb{E}[Y] + 2\beta\mathbb{E}\left[\vec{X}\right] + 2\alpha \tag{58}$$

$$\alpha = \mathbb{E}[Y] - \vec{\beta} \cdot \mathbb{E}\left[\vec{X}\right]$$
(59)

(60)

just like when X was univariate

Backup: Gory details: the slopes

$$-2\mathbb{E}\left[Y\vec{X}\right] + 2\mathbb{E}\left[\vec{X}\otimes\vec{X}\right]\beta + 2\alpha\mathbb{E}\left[\vec{X}\right] = 0$$
(61)

$$\mathbb{E}\left[Y\vec{X}\right] - \alpha \mathbb{E}\left[\vec{X}\right] = \mathbb{E}\left[\vec{X} \otimes \vec{X}\right]\beta$$

$$\mathbb{E}\left[Y\vec{X}\right] - \left(\mathbb{E}\left[Y\right] - \vec{\beta} \cdot \mathbb{E}\left[\vec{X}\right]\right)\mathbb{E}\left[\vec{X}\right] - \mathbb{E}\left[\vec{X} \otimes \vec{X}\right]\beta$$
(62)
(63)

$$\mathbb{E}\left[Y\vec{X}\right] - (\mathbb{E}\left[Y\right] - \beta \cdot \mathbb{E}\left[\vec{X}\right])\mathbb{E}\left[\vec{X}\right] = \mathbb{E}\left[\vec{X} \otimes \vec{X}\right]\beta$$

$$Cov\left[Y, \vec{X}\right] = Var\left[\vec{X}\right]\beta$$
(63)
(64)

$$[X] = \operatorname{Var} [X] \beta \tag{64}$$

$$\beta = (\operatorname{Var}\left[\vec{X}\right])^{-1} \operatorname{Cov}\left[Y, \vec{X}\right]$$
(65)

Reduces to $\operatorname{Cov}[Y, X] / \operatorname{Var}[X]$ when X is univariate

Backup: Gory details: the PCA view

The factor of $\mathrm{Var}\left[\vec{X}\right]^{-1}$ rotates and scales \vec{X} to uncorrelated, unit-variance variables

$$\operatorname{Var}\left[\vec{X}\right] = \mathbf{w} \mathbf{\Lambda} \mathbf{w}^{T}$$
(66)

$$\operatorname{Var}\left[\vec{X}\right]^{-1} = \mathbf{w} \mathbf{\Lambda}^{-1} \mathbf{w}^{T}$$
(67)

$$\operatorname{Var}\left[\vec{X}\right]^{-1} = (\mathbf{w}\Lambda^{-1/2})(\mathbf{w}\Lambda^{-1/2})^{T}$$
(68)

$$= \operatorname{Var}\left[\vec{X}\right]^{-1/2} \left(\operatorname{Var}\left[\vec{X}\right]^{-1/2}\right)^{T}$$
(69)

$$\vec{U} \equiv \vec{X} \operatorname{Var} \left[\vec{X} \right]^{-1/2} \tag{70}$$

$$\operatorname{Var}\left[\vec{U}\right] = \mathbf{I} \tag{71}$$

$$\vec{X} \cdot \vec{\beta} = \vec{X} \cdot \operatorname{Var}\left[\vec{X}\right]^{-1} \operatorname{Cov}\left[\vec{X}, Y\right]$$
(72)

$$= \vec{X} \operatorname{Var} \left[\vec{X} \right]^{-1/2} \left(\operatorname{Var} \left[\vec{X} \right]^{-1/2} \right)^{T} \operatorname{Cov} \left[\vec{X}, Y \right]$$
(73)

$$= \vec{U} \operatorname{Cov} \left[\vec{U}, Y \right] \tag{74}$$

(75)

Backup: Square root of a matrix

- A square matrix **d** is a square root of **c** when $\mathbf{c} = \mathbf{d}\mathbf{d}^T$
- If there are any square roots, there are many square roots
 - Pick any **orthogonal** matrix $\mathbf{o}^T = \mathbf{o}^{-1}$
 - $(\mathbf{do})(\mathbf{do})^T = \mathbf{dd}^T$
 - Just like every real number has two square roots...
- If ${\bf c}$ is diagonal, define ${\bf c}^{1/2}$ as the diagonal matrix of square roots
- If $\mathbf{c} = \mathbf{w} \mathbf{\Lambda} \mathbf{w}^T$, one square root is $\mathbf{w} \mathbf{\Lambda}^{1/2}$

Backup/Aside: R^2 is useless

- $R^2 \equiv \text{Var}\left[\text{fittedvalues}\right]/\text{Var}\left[Y\right]$
- Suppose we know $Y = \beta X + \epsilon$, with $\operatorname{Var}[\epsilon] = \sigma^2$ and ϵ independent of X
- Then $R^2 = \frac{\beta^2 \operatorname{Var}[X]}{\beta^2 \operatorname{Var}[X] + \sigma^2}$
- Consequence 1: R^2 can be arbitrarily close to 0, even for the completely correct model (shrink Var [X] and/or grow σ^2)
- Consequence 2: If we make Var[X] bigger, and the linear model is right, R^2 will grow $\uparrow 1$ (unless $\beta = 0$); R^2 measures the range of the regressor, not goodness of fit
- Suppose the linear model is wrong but we use it anyway; $R^2 = \beta^2 \operatorname{Var}[X] / \operatorname{Var}[Y]$
- Consequence 3: \mathbb{R}^2 can be arbitrarily close to 1 even if every single assumption of the linear model is badly, obviously broken
- Conclusion: R^2 tells us nothing about how good or bad our model is that we didn't already know from the mean squared error; you'd be better off forgetting about it
- More: Shalizi (2015), chapter 10

References

Berk, Richard A. 2008. Statistical Learning from a Regression Perspective. New York: Springer-Verlag.

Buja, Andreas, Richard Berk, Lawrence Brown, Edward George, Emil Pitkin, Mikhail Traskin, Linda Zhao, and Kai Zhang. 2014. "Models as Approximations, Part I: A Conspiracy of Nonlinearity and Random Regressors in Linear Regression." arxiv:1404.1578. http://arxiv.org/abs/1404.1578.

Davison, A. C., and D. V. Hinkley. 1997. *Bootstrap Methods and Their Applications*. Cambridge, England: Cambridge University Press. https://doi.org/10.1017/CBO9780511802843.

Hastie, Trevor, Robert Tibshirani, and Jerome Friedman. 2009. The Elements of Statistical Learning: Data Mining, Inference, and Prediction. Second. Berlin: Springer. http://www-stat.stanford.edu/~tibs/ ElemStatLearn/.

Shalizi, Cosma Rohilla. 2015. "The Truth About Linear Regression." Online Manuscript. http:///www.stat. cmu.edu/~cshalizi/TALR.

. n.d. Advanced Data Analysis from an Elementary Point of View. Cambridge, England: Cambridge University Press. http://www.stat.cmu.edu/~cshalizi/ADAfaEPoV.

Wiener, Norbert. 1949. Extrapolation, Interpolation, and Smoothing of Stationary Time Series: With Engineering Applications. Cambridge, Massachusetts: The Technology Press of the Massachusetts Institute of Technology.