

Modeling

Raguvir Kunani

Discussion 6

February 28, 2020

Big Picture View of Modeling

We model processes that we don't fully understand or processes that inherently have some randomness to them (e.g. how long will it take for you to walk to class).

What does modeling mean?

- sampling* { 1. Identify a target variable that you want to predict
- 2. Gather some observational data about that target variable
- stats + ML* { 3. Propose some relationship between the columns in your observational data and the target variable
- 4. Use ~ machine learning ~ to see how well that relationship actually predicts the target variable
- 5. Repeat steps 3 and 4 for multiple different relationships
- 6. Choose the relationship that best predicts the target variable

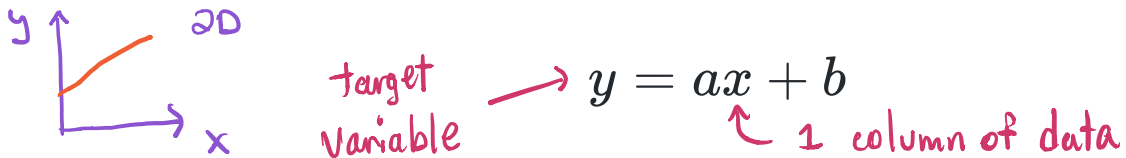
Foreword about Modeling

This class from now on is pretty much only going to be about modeling. But, we will quickly move past the abstract view of modeling and dive into the nitty gritty of the math behind modeling.

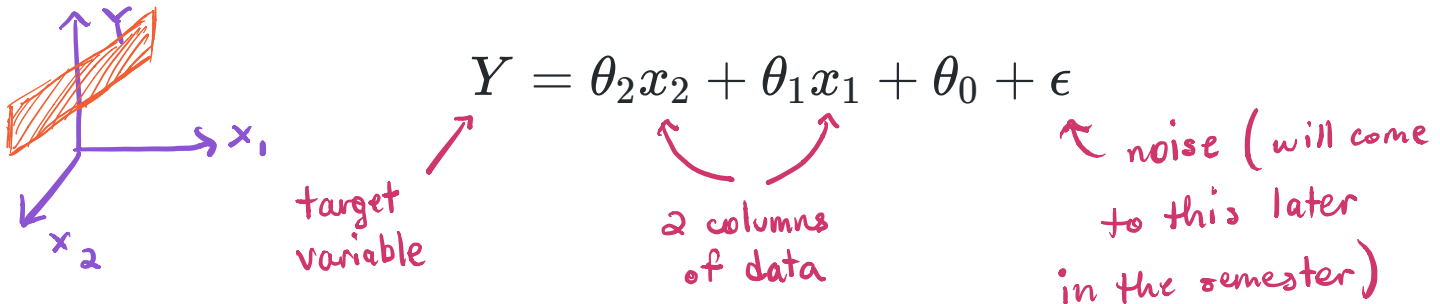
Although the math may get overwhelming sometimes, **don't lose sight of the big picture**. All the math we will do from now on fits somewhere into the previous 6 steps of modeling.

The Math of Modeling: The Model

After you've identified a variable you want to predict and gathered some data about that target variable, step 3 of modeling data says to propose a relationship between the columns of your data and the target variable. In formal terms, this is called **proposing a model**. For example,



is a model you've seen before. In DS 100, we'll graduate to models like:



The Math of Modeling: The Loss Function

After proposing a model, step 4 says to figure out how well your model actually predicts the target variable. For that, we will need a **loss function** that quantifies how *bad* your model is at predicting the target variable.

loss on 1 pt $\rightarrow l(\theta_2, \theta_1, \theta_0) = (y - (\underbrace{\theta_2 x_2 + \theta_1 x_1 + \theta_0}_{\text{model's prediction}}))^2$

loss on entire dataset $\rightarrow L(\theta_2, \theta_1, \theta_0) = \frac{1}{n} \sum_{i=1}^n (\underbrace{y_i - (\theta_2 x_{2,i} + \theta_1 x_{1,i} + \theta_0)}_{l(\theta_2, \theta_1, \theta_0)})^2$

This is the first example of where it is easy to get lost in the math.

Remember the goal of the loss function and that will help you understand the math more easily!

The Math of Modeling: Minimizing the Loss Function

At this point we have a model $\hat{y} = \theta_2 x_2 + \theta_1 x_1 + \theta_0$, but we still need to actually choose some values for θ_2 , θ_1 and θ_0 .

Logically, we want to choose the θ_2 , θ_1 and θ_0 that *minimizes* the loss function, since the loss function tells us how *bad* our model is.

Minimizing a function involves calculus, so again it is easy to get lost in the math here. Keep the big picture in mind!

Goals of Modeling: Interpretability vs. Accuracy

Ideally, we want our models to make predictions about their target variables as accurately as possible. But, in the real world, we also care about how interpretable our model is.

For examples, there is a model called COMPAS that predicts how likely a defendant is to re-offend. Judges can use this model to help them make legal decisions. In this scenario, it is of high importance that the model is interpretable, because the judge must be able to explain their decision.

This is a preview of a topic later in the semester (bias-variance tradeoff).

Worksheet!

Feedback Form

This *anonymous* form is for me to learn what I can do to ensure you all get the most of discussion. This form will be open all semester, and I'll be checking it regularly. Be as ruthless as you want, I promise my feelings won't get hurt.

Feedback Form: tinyurl.com/raguvirTAfeedback