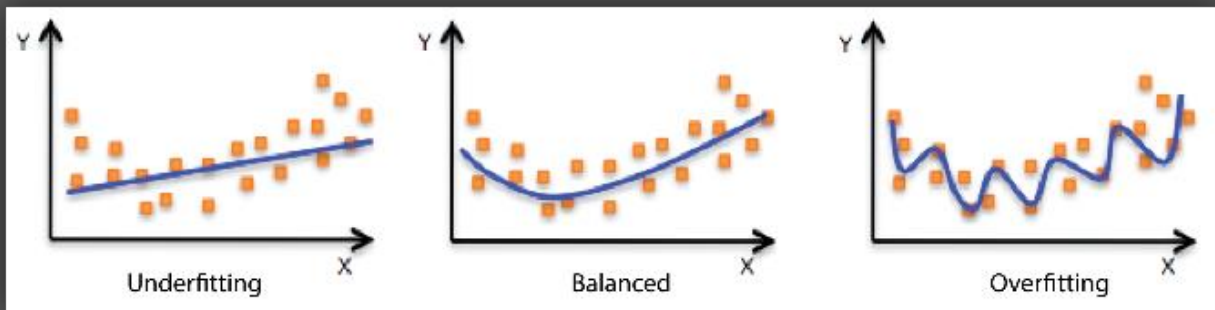




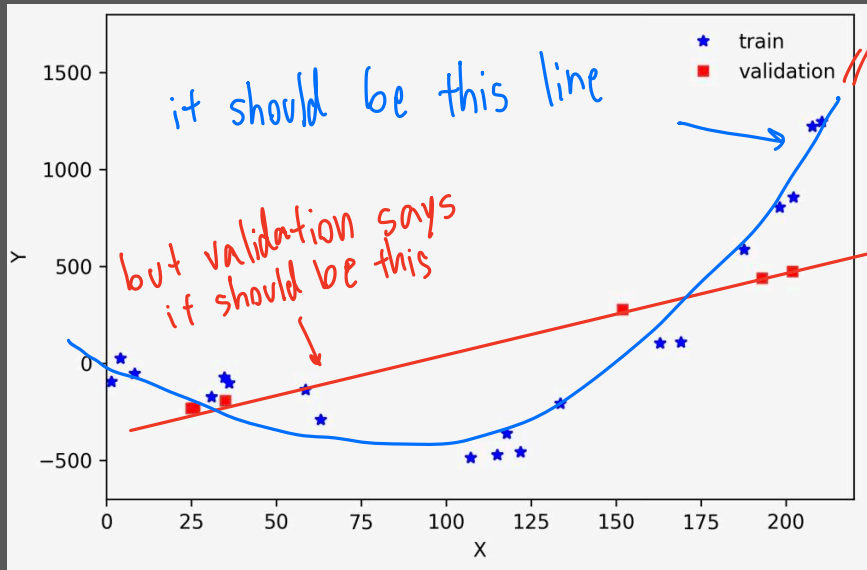
Discussion CAEsaaaaaar



# Recap



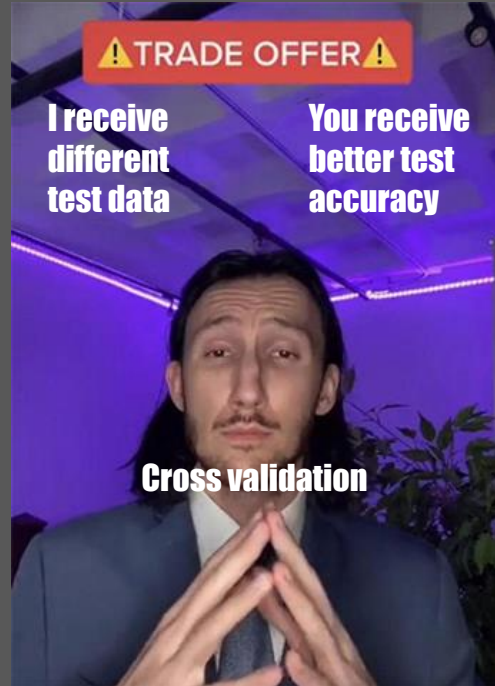
**Motivation:** we have validation dataset to measure how well the model is. But what if the validation dataset is poorly-chosen?



// for now, we  
can think  
validation is  
similar to  
test data //

← result  
in poor  
score

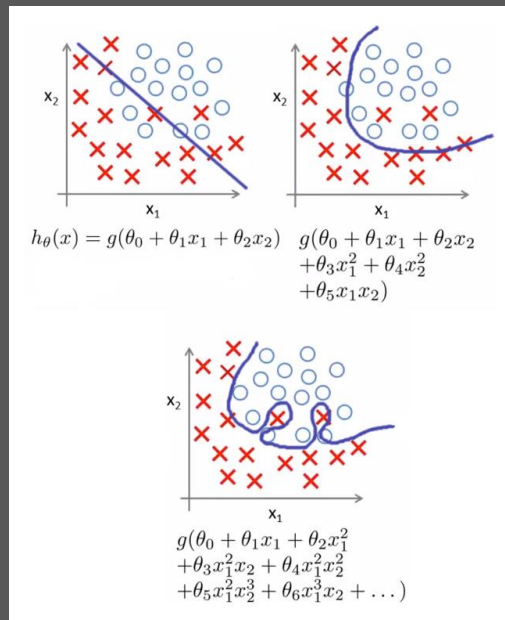
**Solution:** repeat the trials,  
change validation dataset, <sup>(test)</sup>  
average accuracy across all  
trials



**Recap:** We can make the model more complex to capture non-linear data

**Problem:** What is the right degree complexity?

Or what features should be dropped?



**Solution:**

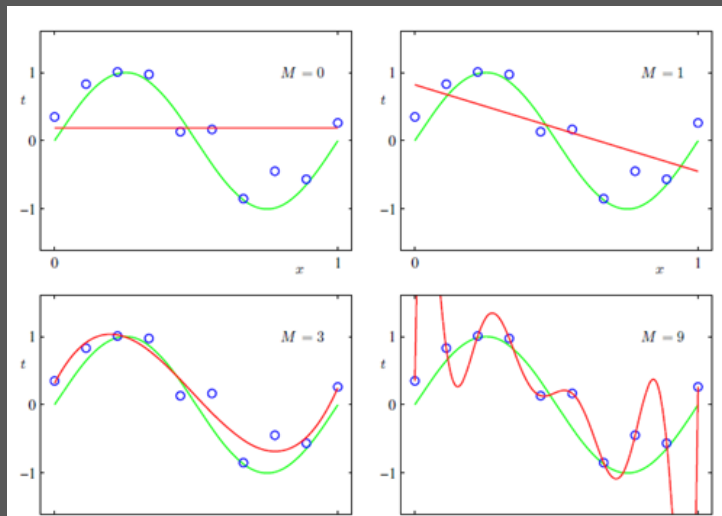
Re



gura

rization

# Regularization



$$y = ax + b$$



$$y = a_1x + a_2x^2 + a_3x^3 + b \rightarrow$$

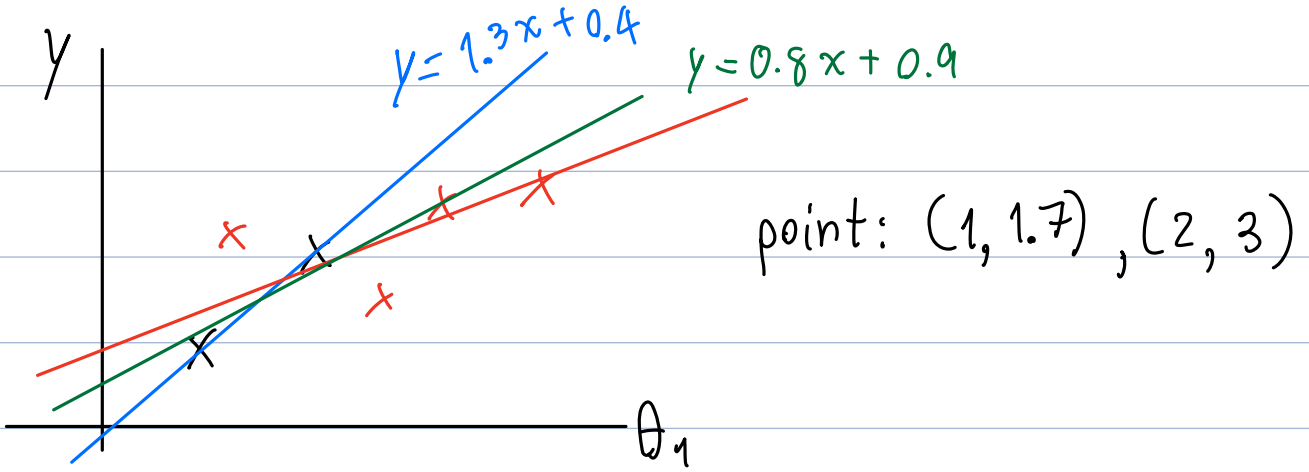
$$y = a_1x + a_2x^2 + \dots + a_qx^q + b$$



reminder : 
$$\text{Cost}(\Theta) = \frac{1}{2m} \sum_i \left[ \text{predicted}^{(i)} - \text{observed}^{(i)} \right]^2$$







Modify  $\text{Cost}(\theta_1)$  to be =  $\sum_i [\text{predicted}_{\theta_1}^{(i)} - \text{observed}_0^{(i)}]^2$   
 $+ \lambda \times (\text{slope})^2$   
 $\lambda = 1$  penalty term

Ex. Blue:  $0 + (1)(1.3)^2 = 1.69$

Green:  $\left[ (1.7 - 1.5)^2 + (2.5 - 3)^2 \right] + 1(0.8)^2$   
 $= 0.04 + 0.25 + 0.64$   
 $= 0.93$

lower cost function value

Expand from this 2D example to 3D, each feature has its own slope  $a_i$ , then in general

Ridge regression:  $\text{Cost}(a) = \frac{1}{2m} \sum_i (\text{pred}_a^{(i)} - \text{obs}^{(i)})^2$   
 $+ \lambda (a_1^2 + a_2^2 + \dots + a_n^2)$

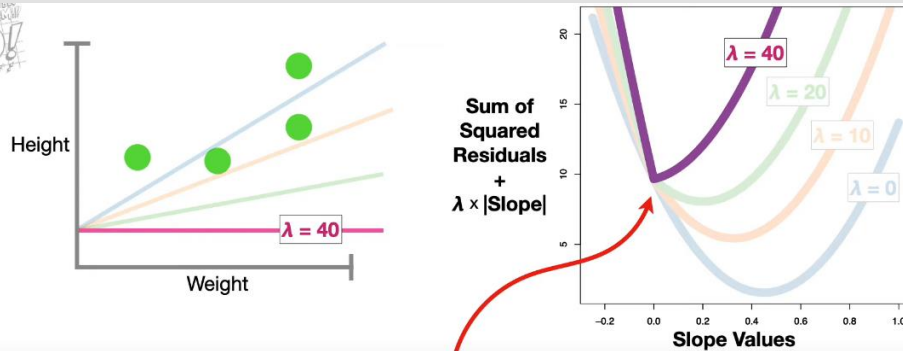
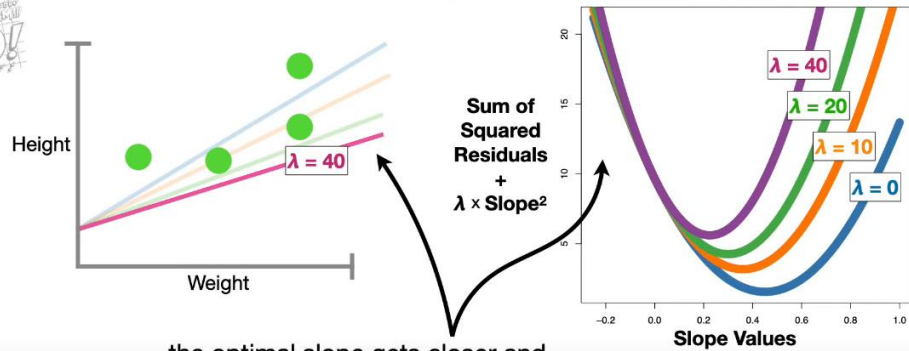
Key takeaway: We find a way to make the model underfits, so we can get higher testing accuracy even if training accuracy is low

Note: there is another type of regularization, which is called **Lasso**

The difference is that the penalty term, we use **absolute** instead of squaring the parameters

Ridge	Lasso
Squared the parameters	Take absolute of the parameters
Parameters get close to zero	Parameters can reach zero
Better when we believe every parameters are useful	Can exclude useless parameters

$$\text{Cost}(\theta) = \frac{1}{2} \sum_{i=1}^n (\text{predicted}_{\theta}(x_i) - y_i)^2 + \lambda \sum_i |\theta_i|$$



Now the lowest point in the **purple curve**, aka, the optimal slope given the **Absolute Value Penalty** when  $\lambda = 40$ , is 0.

Statquest Youtube video:  
Ridge vs Lasso Regression,  
Visualized!!!

[https://www.youtube.com/watch?v=Xm2C\\_gTAI8c](https://www.youtube.com/watch?v=Xm2C_gTAI8c)

But you know, I learned something today



- We use cross validation to average models across all trials, instead of accidentally pick the invalid test data
- We use ridge/lasso regression to lower training accuracy, but get higher test accuracy