

## Model Selection

$$L_{\text{train}} = \frac{1}{n} \sum_{i=1}^n L(f(x_i^{\text{train}}), y_i^{\text{train}}) \leftarrow \text{Empirical Risk}$$

but we want to measure the error on a holdout set.

$$L_{\text{test}} = \frac{1}{n} \sum_{i=1}^n L(f(x_i^{\text{test}}), y_i^{\text{test}}) \leftarrow \text{Expected Risk.}$$

Expected Risk != Empirical Risk

Training Error is our empirical risk and test set error is  
our best approximation of expected risk.

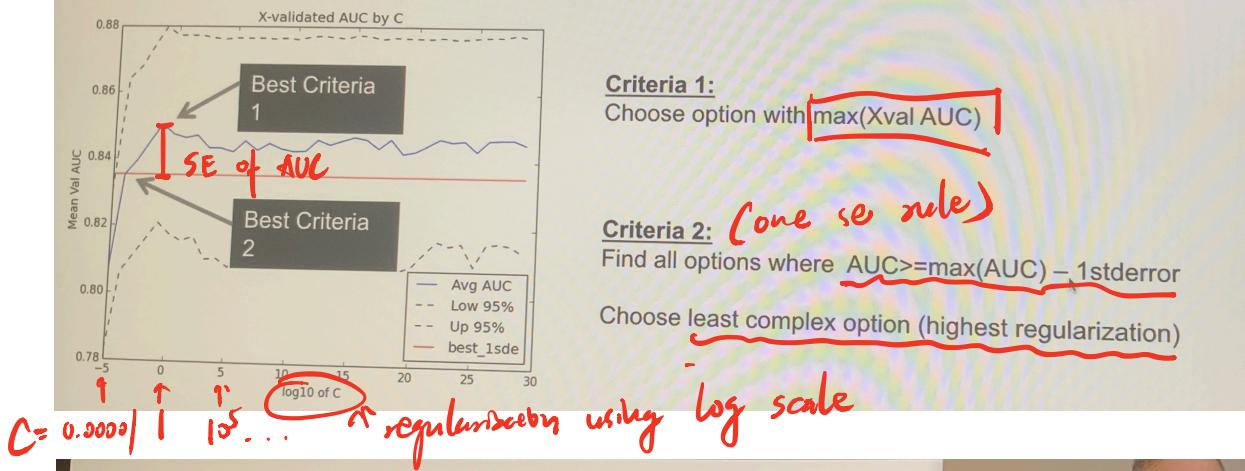
Example Routine:

Train Data. Validation Data

1. For each configuration  $c$  in the set  $C = \{\text{Algorithm} \times \text{Feature Set} \times \text{Hyper-parameters}\}$ 
  - find the function  $f_c$  such that:  $\hat{f}_c = \underset{f \in F}{\operatorname{argmin}} R_{\text{train}}$
  - With  $f_c$  estimated, get the validation loss:  $P_c^{\text{val}} = \frac{1}{n} \sum_{i=1}^n L(\hat{f}_c(x_i^{\text{val}}), y_i^{\text{val}})$
2. Choose the optimal configuration  $C_{\text{opt}}$  such that:  $C_{\text{opt}} = \underset{c \in C}{\operatorname{argmax}} P_c^{\text{val}}$
3. Define New train = Train + Validation data
4. Find the function  $f$  such that  $\hat{f} = \underset{f \in F}{\operatorname{argmin}} R^{\text{new train}}$
5. Estimate the test loss as:  $R^{\text{test}} = \frac{1}{n} \sum_{i=1}^n L(\hat{f}(x_i^{\text{test}}), y_i^{\text{test}})$

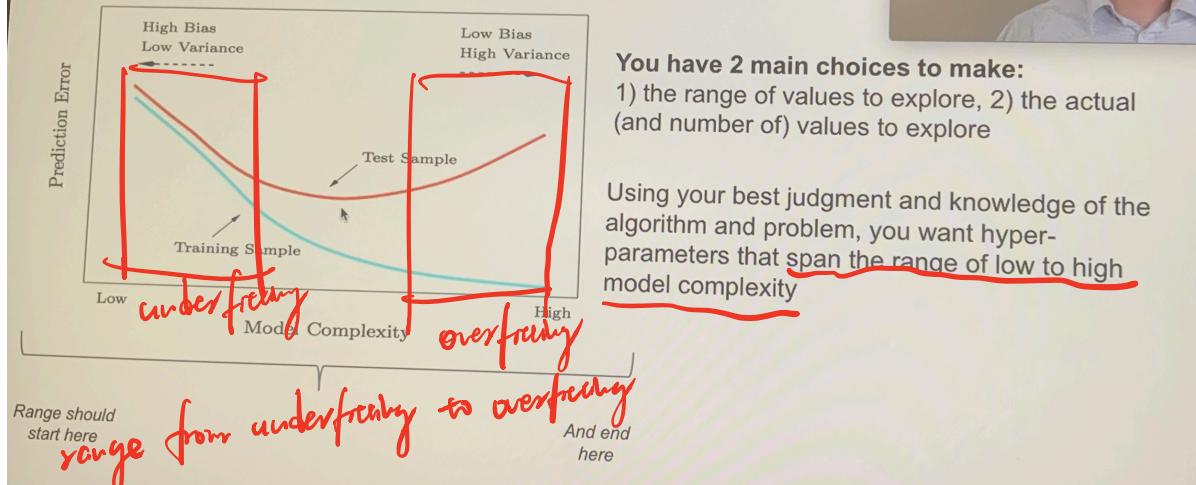
## RULE FOR CHOOSING BEST OPTION

When we zoom in we can see that statistically speaking, everything above  $C=1$  is essentially the same (i.e., strongly overlapping confidence intervals). There are two ways to choose  $C$  here.



## ON CHOOSING HYPER-PARAMETERS

Software will do most of the model training work, but the scientist still needs to instruct the software on which hyper-parameters to explore.



Validation loss metric does not have to be the same as the training loss.

- Sometimes the loss metric for an application is not very easy or even possible to directly minimize.
- We use logistic-loss to find  $f$ , but we need a loss metric that enables optimal ranking (such as AUC)
- We might want to minimize training loss, but maximize the validation loss (logistic loss vs AUC)  
*↑  
training loss      validation loss.*

Hyperparameter selection:

1. Experiment with parameters that span the range of low-to-high model complexity (from underfitting to overfitting).

Grid search,

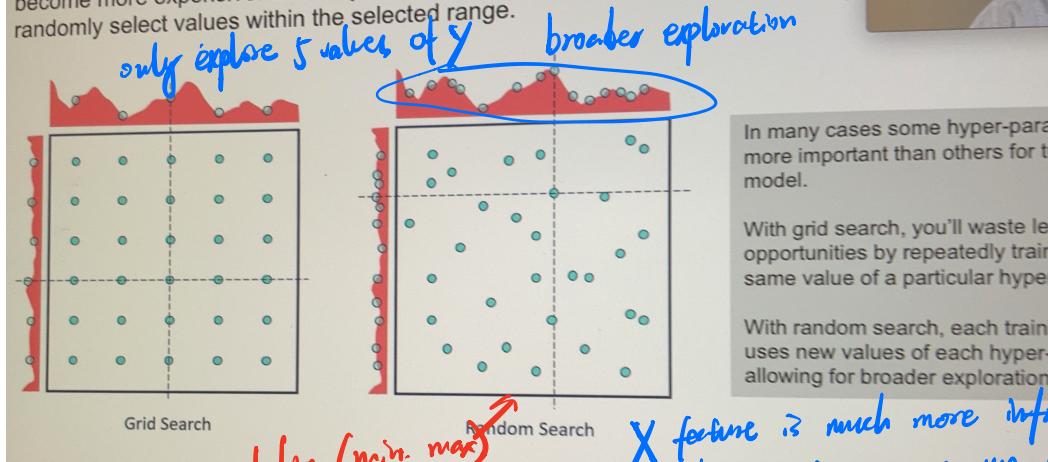
With grid search, you waste learning opportunities by repeatedly training on the same value of a particular hyper-parameter.

Random Search:

With Random search, each training iteration uses new values of each hyper-parameter, allowing for broader exploration of the space.

## ON CHOOSING HYPER-PARAMETERS

When optimizing multiple hyper-parameters, the search starts to become more expensive. One way to optimize learning per unit time is to randomly select values within the selected range.



In many cases some hyper-parameters are more important than others for the model.

With grid search, you'll waste lots of opportunities by repeatedly training the same value of a particular hyper-parameter.

With random search, each trial uses new values of each hyper-parameter, allowing for broader exploration.

defining (min, max)  
random sampling  
using uniform distribution

X feature is much more informative than Y feature, so we would like to explore X to a broader extent.

Piggy off the bee, score with simple model

(Random Forest or Logistic Regression) to get a baseline.

If there's a significant difference between the performance of Random Forest and Logistic Regression, it means there's a lot of non-linearity in your data that you need to capture.

Iterate towards better Model:

Starting from step 2, try new features and new algorithms. In Real life, simplicity is preferred. so don't use a model that is 5X complexity of your current model, but with 0.1X increase in metric.