

Cross Validation

- can use to choose a model, M , or to simply estimate the generalization error of a model: $E_{\mathcal{D} \sim \mathcal{P}^m, X, Y \sim \mathcal{P}} [l(Y, f_M(X))]$
- given our data, we can construct a pretty good estimator $\hat{\epsilon}_M$, of this error w/ the following algo:

For Min M ^{can be a grid of params}

For $k=1$ to K :

• Fit model w/ M on k ^{or whatever other loss you want}

$\hat{\epsilon}_M^{(k)} = \text{MSE}_M^{(k)}$

$\Rightarrow \hat{\epsilon}_M = \frac{1}{K} \sum_{k=1}^K \hat{\epsilon}_M^{(k)}$

$M^* = \arg \min_M (\hat{\epsilon}_M)$

should be repeated

→

• how to choose K ?

2 extremes:

LOOCV

- low bias

- training uses lots of data, so models are general and ~~do not have error~~ ~~are not~~ ~~biased~~ biased to the high end

2-Fold

- high bias

- training uses few data, so models aren't great and have error biased on the high end.

- high variance

- all models use essentially the same data and are thus highly correlated. Since we average the errors in the end, the correlation adds a lot to the variance

- low variance

- low for opposite reason

A medium K

(5 \rightarrow 10) is

usually preferred

- to hyperparameter tune and evaluate performance, we could do CV for the tuning, and test on a test set that wasn't used in tuning
- ^{for more robustness} as ¹ we could use nested CV (i.e., CV for hyperparam tuning, and for ~~the~~ generalization error)

For $l=1$ to L :

$$D' = D - \{l\}$$

For M in \mathcal{M} :

For $n=1$ to K :

$$D'' = D' - \{n\}$$

• train model w/ D''

$$\hat{\epsilon}_M^{(l,n)} = \text{MSE}_M^{(l,n)}$$

$$\hat{\epsilon}_M^{(l)} = \frac{1}{K} \sum_{n=1}^K \hat{\epsilon}_M^{(l,n)}$$

$$\hat{\epsilon}^{(l)} = \min_M (\hat{\epsilon}_M^{(l)})$$

$$\hat{\epsilon} = \frac{1}{L} \sum_{l=1}^L \hat{\epsilon}^{(l)}$$

- outer loop tests error
- inner loop does CV hyperparam tuning.
- After this, we use normal CV hyperparam tuning to pick tuned params.
- then we train on all data and deploy

Bias / Variance For Supervised learning (i.e., For Function estimation)

- Assuming in the real world that Y is generated from a deterministic function by adding noise: $Y = f(X) + \epsilon$, it is not difficult to show that:

$$E_{\mathcal{D} \sim P^m, \epsilon \sim \mathcal{N}(0, \sigma^2)} \left[(Y - \hat{f}_{\mathcal{D}}(X))^2 \right] = \text{Bias}^2(\hat{f}_{\mathcal{D}}(X)) + \text{Var}[\hat{f}_{\mathcal{D}}(X)] + \text{Var}[\epsilon]$$

" MSE
our estimator of Y

$$\text{w/ } \text{Bias}^2(\hat{f}_{\mathcal{D}}(X)) = \left(E[\hat{f}_{\mathcal{D}}(X)] - f(X) \right)^2; \text{Var}[\hat{f}_{\mathcal{D}}(X)] = E\left[(\hat{f}_{\mathcal{D}}(X) - E[\hat{f}_{\mathcal{D}}(X)])^2 \right]$$

