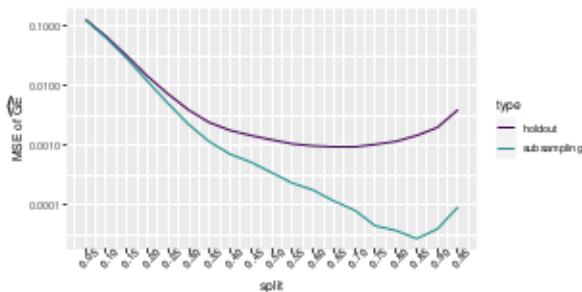


# BIAS-VARIANCE ANALYSIS FOR SUBSAMPLING



- MSE of  $\widehat{GE}$  strictly better for SS
- Smaller var of SS enables to use larger  $s$  for optimal choice
- The optimal split rate now is a higher  $s \approx 0.8$ .
- Beyond  $s = 0.8$ : MSE goes up because var doesn't go down as much as we want due to increasing overlap in trainsets (see later)

## DEDICATED TESTSET SCENARIO - ANALYSIS

- Goal: estimate  $GE(\hat{f}) = E[L(y, \hat{f}(\mathbf{x}))]$  via

$$\widehat{GE}(\hat{f}) = \frac{1}{m} \sum_{(\mathbf{x}, y) \in \mathcal{D}_{\text{test}}} L(y, \hat{f}(\mathbf{x}))$$

Here, only  $(\mathbf{x}, y)$  are random, they are i.i.d. fresh test samples

- This is: average over i.i.d  $L(y, \hat{f}(\mathbf{x}))$ , so directly know  $E$  and var.  
And can use CLT to approx distrib of  $\widehat{GE}(\hat{f})$  with Gaussian.
- $E[\widehat{GE}(\hat{f})] = E[L(y, \hat{f}(\mathbf{x}))] = GE(\hat{f})$
- $V[\widehat{GE}(\hat{f})] = \frac{1}{m} V[L(y, \hat{f}(\mathbf{x}))]$
- So  $\widehat{GE}(\hat{f})$  is unbiased estimator of  $GE(\hat{f})$ , var decreases linearly in testset size, have an approx of full distrib (can do NHST, CIs, etc.)
- NB: Gaussian may work less well for e.g. 0-1 loss, with  $E$  close to 0, can use binomial or other special approaches for other losses



## PESSIMISTIC BIAS IN RESAMPLING

- Estim  $GE(\mathcal{I}, n)$  (surrogate for  $GE(\hat{f})$ ) when  $\hat{f}$  is fit on full  $\mathcal{D}$ , with  $|\mathcal{D}| = n$  via resampling based estim  $\widehat{GE}(\mathcal{I}, n_{\text{train}})$

$$\widehat{GE}(\mathcal{I}, \mathcal{J}, \rho, \lambda) = \text{agr} \left( \rho \left( \mathbf{y}_{\mathcal{J}_{\text{test},1}}, \mathbf{F}_{\mathcal{J}_{\text{test},1}; \mathcal{I}(\mathcal{D}_{\text{train},1}, \lambda)} \right), \right. \\ \vdots \\ \left. \rho \left( \mathbf{y}_{\mathcal{J}_{\text{test},B}}, \mathbf{F}_{\mathcal{J}_{\text{test},B}; \mathcal{I}(\mathcal{D}_{\text{train},B}, \lambda)} \right) \right),$$

- Let's assume agr is avg and  $\rho$  is loss-based, so  $\rho_L$
- The  $\rho$  are simple holdout estims. So:

$$\mathbf{E}[\widehat{GE}(\mathcal{I}, \mathcal{J}, \rho, \lambda)] \approx \mathbf{E}[\rho \left( \mathbf{y}_{\mathcal{J}_{\text{test}}}, \mathbf{F}_{\mathcal{J}_{\text{test}}; \mathcal{I}(\mathcal{D}_{\text{train}}, \lambda)} \right)]$$

- NB1: In above, as always for  $GE(\mathcal{I})$ , both  $\mathcal{D}_{\text{train}}$  and  $\mathcal{D}_{\text{test}}$  (and so  $\mathbf{x} \in \mathcal{D}_{\text{test}}$ ) are random vars, and we take E over them
- NB2: Need  $\approx$  as maybe not all train/test sets in resampling of exactly same size



$$\mathbf{E}[\widehat{GE}(\mathcal{I}, \mathcal{J}, \rho, \lambda)] \approx \mathbf{E}[\rho(\mathbf{y}_{\mathcal{J}_{\text{test}}}, \mathbf{F}_{\mathcal{J}_{\text{test}}; \mathcal{I}(\mathcal{D}_{\text{train}}; \lambda)})] =$$

$$\mathbf{E} \left[ \frac{1}{m} \sum_{\substack{(\mathbf{x}, y) \in \mathcal{D}_{\text{test}} \\ (\mathbf{x}, y) \in \mathcal{D}_{\text{test}}} L(y, \mathcal{I}(\mathcal{D}_{\text{train}})(\mathbf{x})) \right] \equiv GE(\mathcal{I}, n_{\text{train}})$$

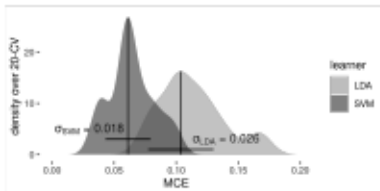


- So when we use  $\widehat{GE}(\mathcal{I}, \mathcal{J}, \rho, \lambda)$  to estimate  $GE(\mathcal{I}, n)$ , our expected value is nearly correct, it's  $GE(\mathcal{I}, n_{\text{train}})$
- But fitting  $\mathcal{I}$  on less data ( $n_{\text{train}}$  vs full  $n$ ) usually results in model with worse perf, hence estimator is pessimistically biased
- Bias the stronger, the smaller our training splits in resampling.

# NO INDEPENDENCE OF CV RESULTS

- Similar analysis as before holds for CV
- Might be tempted to report distribution or SD of individual CV split perf values, e.g. to test if perf of 2 learners is significantly different
- But  $k$  CV splits are not independent

A t-test on the difference of the mean GE estimators yields a highly significant p-value of  $\approx 7.9 \cdot 10^{-5}$  on the 95% level.

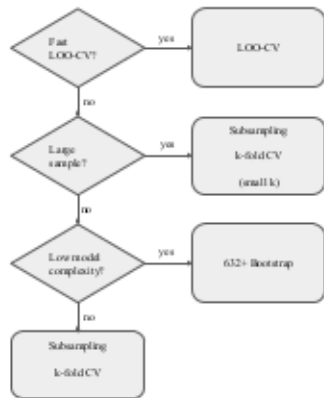


LDA vs SVM on a pass classification problem, performance estimation via 20-CV w.r.t. MCE.

# NO INDEPENDENCE OF CV RESULTS

- $\mathbb{V}[\widehat{GE}]$  of CV is a difficult combination of
  - average variance as we estim on finite trainsets
  - covar from test errors, as models result from overlapping trainsets
  - covar due to the dependence of trainsets and test obs appear in trainsets
- Naively using the empirical var of  $k$  individual  $\widehat{GE}$ s (as on slide before) yields biased estimator of  $\mathbb{V}[\widehat{GE}]$ . Usually this underestimates the true var!
- Worse: there is no unbiased estimator of  $\mathbb{V}[\widehat{GE}]$  [Bengio, 2004]
- Take into account when comparing learners by NHST
- Somewhat difficult topic, we leave it with the warning here

# SHORT GUIDELINE



- 5-CV or 10-CV have become standard.
- Do not use hold-out, CV with few folds, or SS with small split rate for small  $n$ . Can bias estim and have large var.
- For small  $n$ , e.g.  $n < 200$ , use LOO or, probably better, repeated CV.
- For some models, fast tricks for LOO exist
- With  $n = 100,000$ , can have "hidden" small-sample size, e.g. one class very small
- SS usually better than bootstrapping. Repeated obs can cause problems in training, especially in nested setups where the "training" set is split up again.