

Sample Size Considerations for Raman Spectroscopic Cell Identification

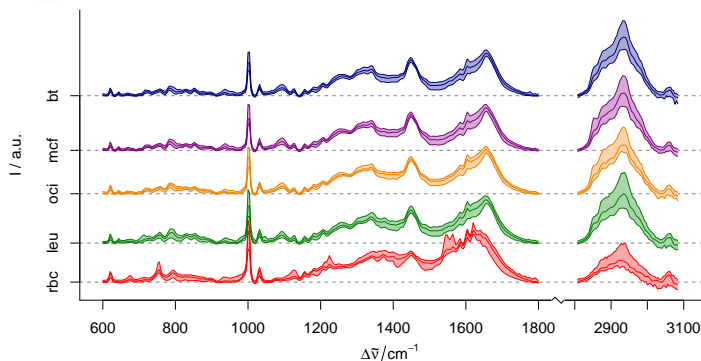
Claudia Beleites¹ (Claudia.Beleites@ipht-jena.de),
Ute Neugebauer^{1,2}, Thomas Bocklitz³, Christoph Krafft¹,
and Jürgen Popp^{1,3}

¹Institute of Photonic Technology, Jena/Germany

²Center for Sepsis Control and Care, University Hospital Jena, Jena/Germany

³Institute of Physical Chemistry and Abbe Center of Photonics,
Friedrich-Schiller-University Jena/Germany

Chemometrics in Analytical Chemistry 2012



rbc	normal red blood cells	5 donors	372 spectra
leu	normal leukocytes	5 donors	569 spectra
oci	acute myelotic leukemia cell line OCI-AML	5 batches	518 spectra
mcf	breast cancer cell line MCF-7	5 batches	558 spectra
bt	breast cancer cell line BT-20	5 batches	532 spectra
total			2549 spectra

- Samples are ever too few...

- Samples are ever too few...
- But: how many samples (per class) do we **really** need?

- Samples are ever too few...
- But: how many samples (per class) do we **really** need?
- ...to train a good classifier:
⇒ **learning curve** model performance $p = f(n_{\text{train}})$

- Samples are ever too few...
- But: how many samples (per class) do we **really** need?
- ...to train a good classifier:
⇒ **learning curve** model performance $p = f(n_{\text{train}})$
- ...to **precisely** measure the classifier's performance:
⇒ **confidence interval for test results**: $\sigma^2(\hat{p}) = \frac{p(1-p)}{n_{\text{test}}}$

Take home message I

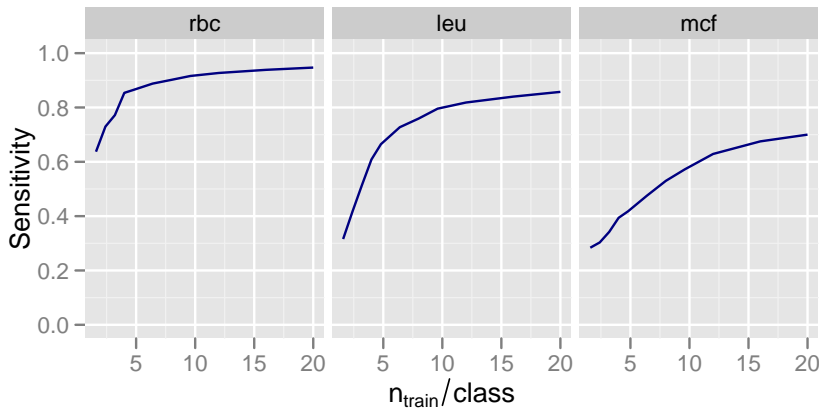
Training a good classifier is not enough:
performance must be **demonstrated**, too.

- PLS-LDA, 25 latent variables (for $n_{\text{train}} / \text{class} < 10$: $\frac{1}{2} n_{\text{train}}$)
- 50× iterated 5-fold cross validation
- 100 growing “small” data sets
- “large” test with 320 – 520 spectra / class
 - ↪ 95 % confidence interval: $p = 0.5 \pm 0.055$
 - ≈ 1 : 9 split for “small” : “large” sets

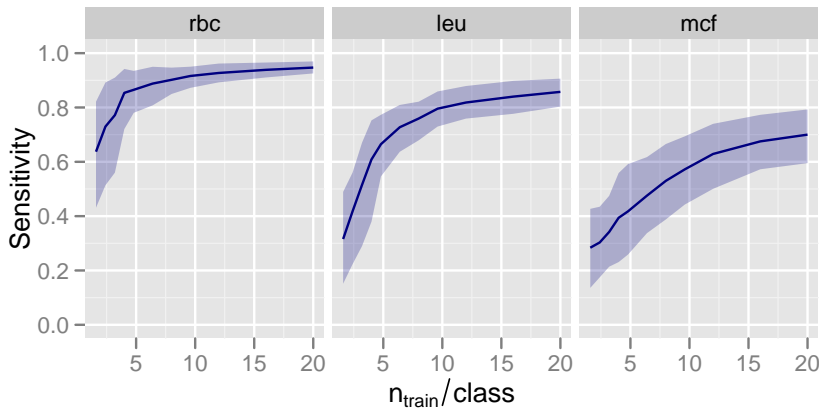
- Learning curve:
- true performance p of model trained with n_{train} samples:

$$\text{error}^2 = \underbrace{\text{Bayes-error}^2 + \text{bias}^2(n_{\text{train}})}_{\text{learning curve}} + \underbrace{\text{var}(n_{\text{train}})}_{\text{model instability}}$$

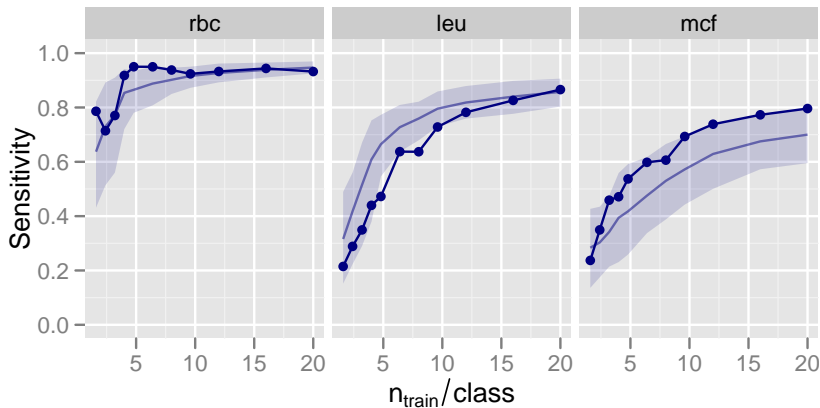
- observed performance $\hat{p} =$ true performance p
 - + systematic test error(n , method)
 - + random test error(n_{test})



- Confidence band: 5th – 95th percentile of observations
- 100 repetitions
- tested with large test set



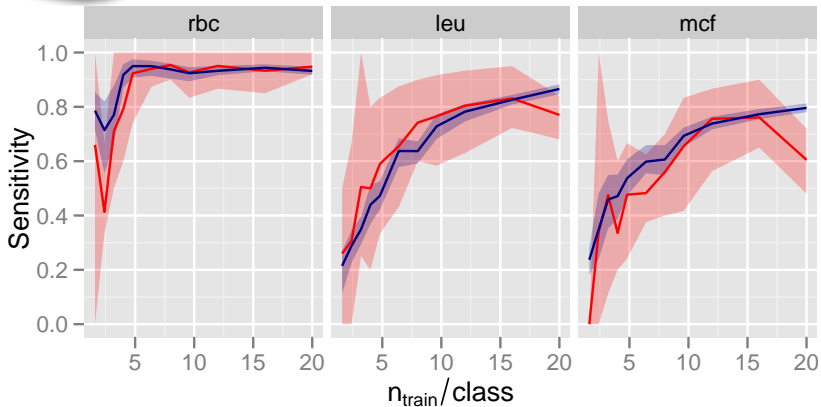
- Confidence band: 5th – 95th percentile of observations
- 100 repetitions
- tested with large test set



- Confidence band: 5th – 95th percentile of observations
- 100 repetitions
- tested with large test set

Take home message II

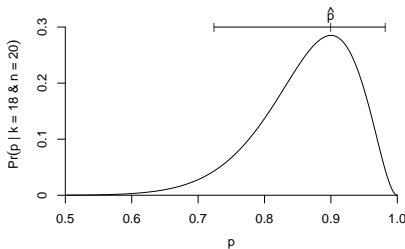
The learning curve of **particular** small data set can differ substantially from average performance for data set of given size.



- Confidence band: 5th – 95th percentile of observations
- single, growing data set: iteration no. 17
- blue: tested with large test set
- red: 50× 5-fold cross validation

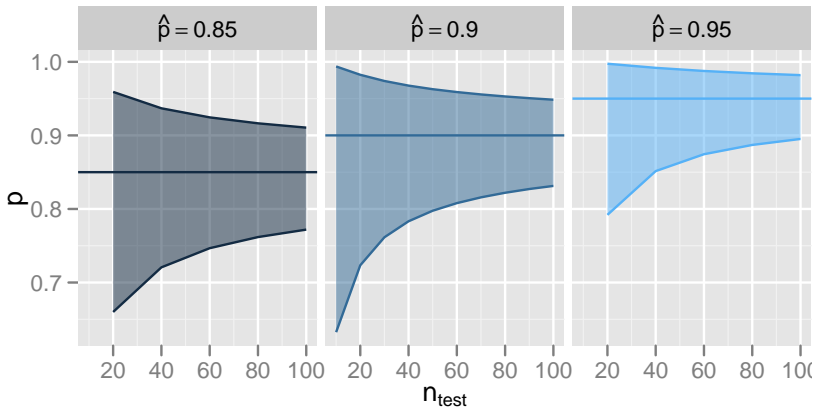
Take home message III

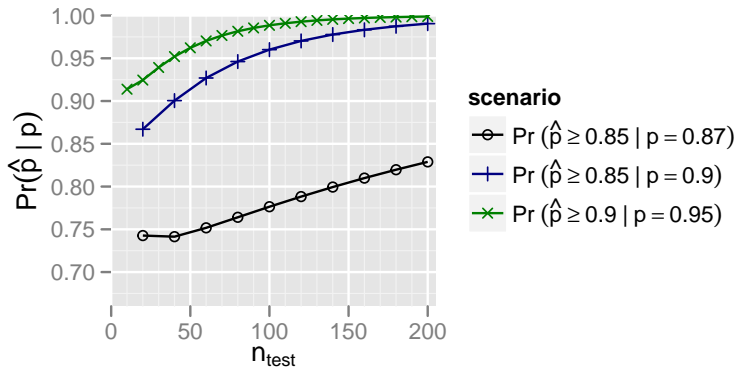
For estimating the learning curve of **particular** small data set, performance estimation uncertainty is **huge**.



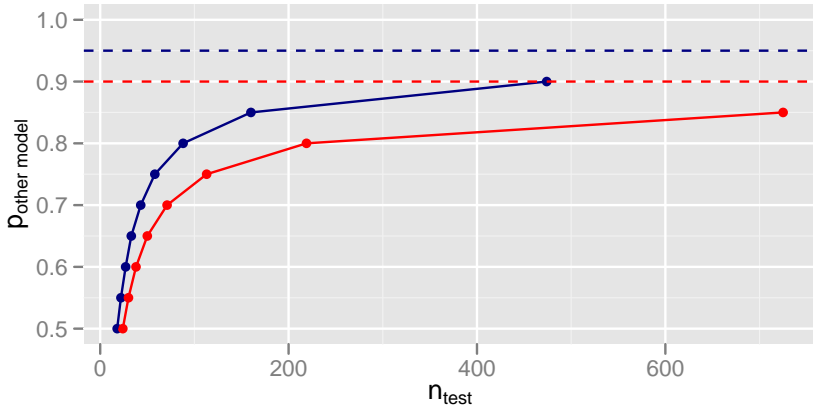
- Classifier performance: proportions
- Statistical description: Bernoulli trial
- Uncertainty on proportion: $\text{var}(\hat{p}) = \frac{p(1-p)}{n_{\text{test}}}$

~> Estimate necessary n_{test}





iphtena Proving Advantage over Other Model



$\alpha = 5 \%, \beta = 20 \%$

from: Fleiss "Statistical Methods for Rates and Proportions"

- More powerful tests available for **paired** test

Summary



- Learning curve: check variance as well as expected performance
- Performance of data set of size n vs. particular data set
- Learning curve is difficult to measure from small sample set:
Uncertainty dominated by **testing**.
- Calculating necessary test sample size
 - Confidence interval width $p \leq \Delta p$
 - Observe $\hat{p} \geq x$
 - Show advantage over model with \hat{p}_A
- Necessary n_{test} often \gg necessary n_{train}