

# Pattern Recognition (PR)

Prof. Dr.-Ing. Andreas Maier  
Pattern Recognition Lab (CS 5), Friedrich-Alexander-Universität Erlangen-Nürnberg  
Winter Term 2020/21



This is a printable version of the slides of the lecture

**Pattern Recognition (PR)**  
*Winter term 2020/21*  
*Friedrich-Alexander University of Erlangen-Nuremberg.*

These slides are released under Creative Commons License Attribution CC BY 4.0.

Please feel free to reuse any of the figures and slides, as long as you keep a reference to the source of these slides at <https://lme.tf.fau.de/teaching/> acknowledging the authors Niemann, Hornegger, Hahn, Steidl, Nöth, Seitz, Rodriguez, and Maier.

Erlangen, October 28, 2020  
Prof. Dr.-Ing. Andreas Maier

# Model Assessment



## No Free Lunch

- In the past lectures, we have come across many learning algorithms and classification techniques.
- They have properties such as
  - low computational complexity
  - incorporation of prior knowledge
  - linearity / non-linearity
  - optimality with respect to certain cost functions, etc.
- Some compute smooth decision boundaries, some compute rather non-smooth decision boundaries.

We really have to ask:

Are there any reasons to favor one algorithm over another?

## No Free Lunch (cont.)

### Theorem

Given a cost function  $f \in \mathcal{F}$ , an algorithm  $A$  and costs  $c_m$  for a specific sample that is iterated on  $m$  times.

The performance of an algorithm is the conditional probability  $P(c_m|f, m, A)$ .

The *No Free Lunch Theorem* states that for any two algorithms  $A_1$  and  $A_2$ :

$$\sum_f P(c_m|f, m, A_1) = \sum_f P(c_m|f, m, A_2)$$

## No Free Lunch (cont.)

Consequences for classification methods:

- If no prior assumptions about the problem are made, there is **NO** overall superior or inferior classification method!
- We should be skeptical regarding studies that demonstrate the overall superiority of a particular method.
- We have to focus on the aspects that matter most for the classification problem:
  - prior information
  - data distribution
  - amount of training data
  - cost functions

## Off-Training Set Error

Off-training set error:

- Specifies the error on samples that are not contained within the training set.
- For large training data sets, the off-training set is necessarily small.
- Used to compare general classification performance of algorithms.

## Off-Training Set Error (cont.)

- Consider a two-class problem with training data set  $\mathcal{D}$  consisting of patterns  $\mathbf{x}_i$  and labels  $y_i = \pm 1$ .
- $y_i$  is generated by an unknown target function:  $F(\mathbf{x}_i) = y_i$ .
- The expected off-training set classification error for the  $k$ -th learning algorithm is:

$$E_k\{e|F, n\} = \sum_{\mathbf{x} \notin \mathcal{D}} p(\mathbf{x}) [1 - \delta(F(\mathbf{x}), h(\mathbf{x}))] p_k(h(\mathbf{x})|\mathcal{D})$$

where  $e$  is the error and  $h(\mathbf{x})$  the hypothesis on the data.

## Off-Training Set Error (cont.)

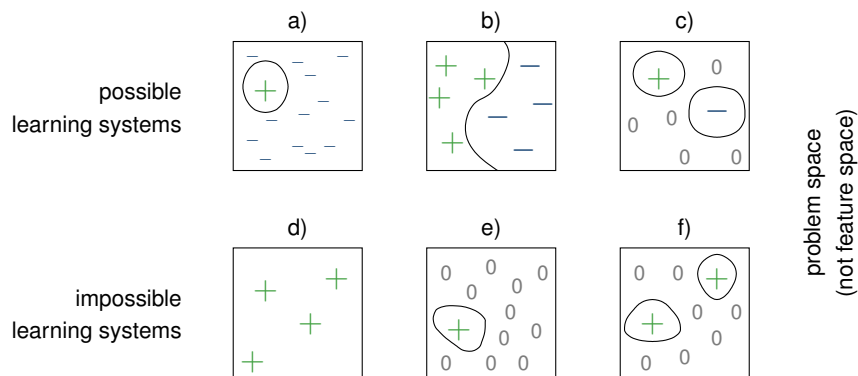


Fig.: Each square represents all possible classification problems. +/− indicates better/worse generalization than the average (adapted from Duda, Hart).

## Bias and Variance

- The *No Free Lunch Theorem* states that there is no general best classifier.
- But we have to assess the quality of a learning algorithm in terms of the **alignment** to the classification problem.
- This can be achieved using the **bias-variance** relation.

### Bias:

- The bias measures the accuracy or quality of the match:  
high bias means poor match.

### Variance:

- The variance measures the precision of specificity for the match:  
high variance implies a weak match.

## Bias and Variance for Regression

The bias-variance relation is very demonstrative in the context of regression:

- Let  $g(\mathbf{x}; \mathcal{D})$  be the regression function.
- The mean-square deviation from the true function  $F(\mathbf{x})$  is:

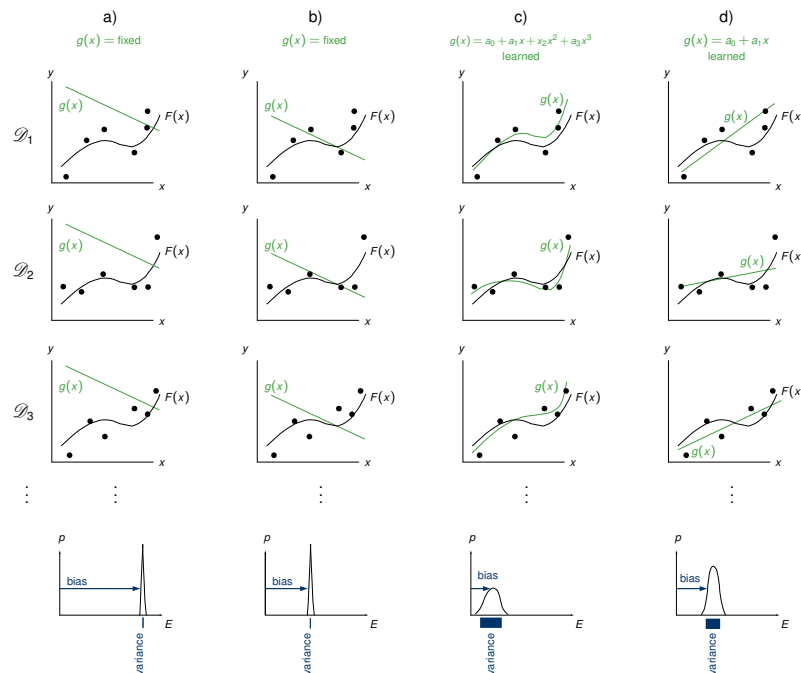
$$\begin{aligned}
 E_{\mathcal{D}} \left\{ \left( g(\mathbf{x}; \mathcal{D}) - F(\mathbf{x}) \right)^2 \right\} \\
 = \underbrace{E_{\mathcal{D}} \left\{ g(\mathbf{x}; \mathcal{D}) - F(\mathbf{x}) \right\}^2}_{(\text{bias})^2} + \underbrace{E_{\mathcal{D}} \left\{ \left( g(\mathbf{x}; \mathcal{D}) - E_{\mathcal{D}} \{ g(\mathbf{x}; \mathcal{D}) \} \right)^2 \right\}}_{\text{variance}}
 \end{aligned}$$

## Bias and Variance for Regression (cont.)

Bias-Variance Trade-Off:

- Methods with **high flexibility** to adapt to the training data
  - generally have low bias
  - but yield high variance.
- Methods with **few parameters** and less degrees of freedom
  - tend to have a high bias, as they may not fit the data well.
  - However, this does not change a lot between different data sets, so these methods generally have low variance.
- Unfortunately, we can virtually never get both zero bias and zero variance!
- We need to have as **much prior information** about the problem as possible to reduce both values.

## Bias and Variance for Regression (cont.)



## Bias and Variance for Classification

Assuming a two-class classification problem:

- In a two-class problem, the target function changes to:

$$F(\mathbf{x}) = p(y = 1|\mathbf{x}) = 1 - p(y = -1|\mathbf{x})$$

- We cannot compare  $g(\mathbf{x}; \mathcal{D})$  and  $F(\mathbf{x})$  based on the mean-square error as in regression.
- For simplicity, let us assume **identical priors**:  $p_1 = p_2 = 0.5$ 
  - The Bayes discriminant  $y_B$  has the threshold 0.5.
  - The Bayes decision boundary is the set of points for which  $F(\mathbf{x}) = 0.5$ .

## Bias and Variance for Classification (cont.)

### Boundary error

- $p(g(\mathbf{x}; \mathcal{D}))$  is the pdf of obtaining a particular estimate of the discriminant given  $\mathcal{D}$ .
- Because of random variations in the training set, the boundary error will depend upon  $p(g(\mathbf{x}; \mathcal{D}))$ .

$$p(g(\mathbf{x}; \mathcal{D}) \neq y_B) = \begin{cases} \int_{0.5}^{\infty} p(g(\mathbf{x}; \mathcal{D})) dg & \text{if } F(\mathbf{x}) < 0.5 \\ \int_{-\infty}^{0.5} p(g(\mathbf{x}; \mathcal{D})) dg & \text{if } F(\mathbf{x}) \geq 0.5 \end{cases}$$

## Bias and Variance for Classification (cont.)

- Convenient assumption that  $p(g(\mathbf{x}; \mathcal{D}))$  is a Gaussian:

$$p(g(\mathbf{x}; \mathcal{D}) \neq y_B) = \Phi \left[ \underbrace{\text{sgn} \left( F(\mathbf{x}) - \frac{1}{2} \right) \cdot \left( E_{\mathcal{D}} \{ g(\mathbf{x}; \mathcal{D}) \} - \frac{1}{2} \right)}_{\text{boundary bias}} \cdot \underbrace{\text{var} (g(\mathbf{x}; \mathcal{D}))^{-1/2}}_{\text{variance}} \right]$$

where  $\Phi$  is a nonlinear function:

$$\Phi(t) = \frac{1}{\sqrt{2\pi}} \int_t^{\infty} e^{-\frac{1}{2}u^2} du$$

- $p(g(\mathbf{x}; \mathcal{D}) \neq y_B)$  represents the incorrect estimation of the Bayes boundary.



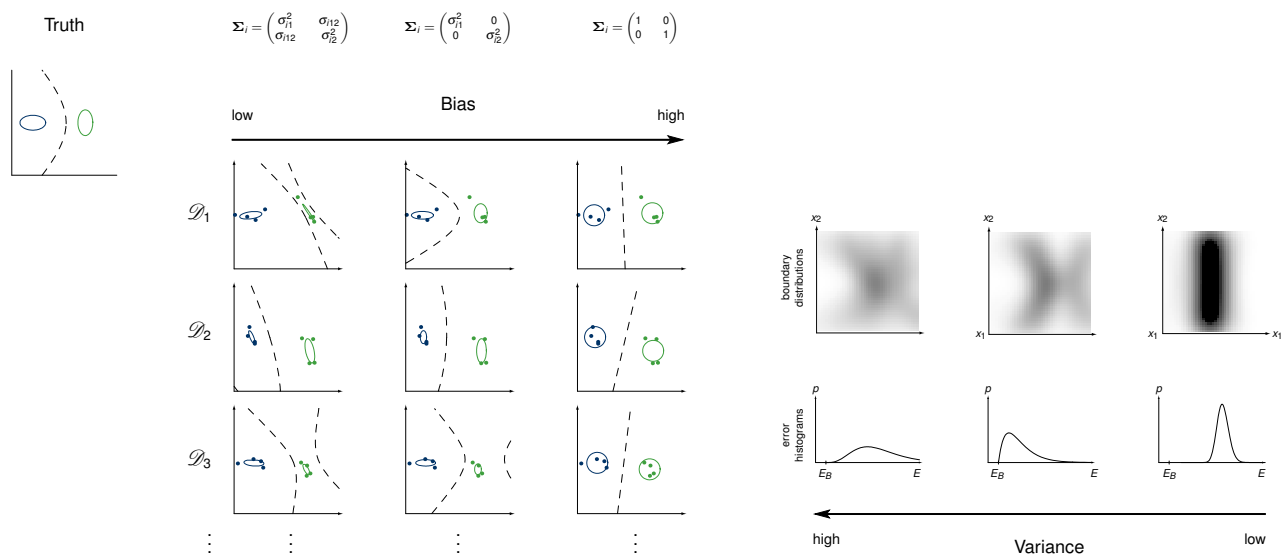
## Bias and Variance for Classification (cont.)

### Conclusions:

- In **regression** the bias-variance relation is additive in  $(\text{bias})^2$  and variance.
- For **classification** the relation is multiplicative and nonlinear.
- In classification the sign of the boundary bias affects the role of the variance in the error.
- Therefore, low variance is generally important for accurate classification.

Variance generally dominates bias in classification!

## Bias and Variance for Classification (cont.)



Adapted from Duda, Hart

# Next Time in Pattern Recognition



## Resampling for Estimating Statistics

### Problem:

- Determine the bias and variance for some learning algorithm applied to a new problem with unknown distributions.

From what we have seen so far, bias and variance change with varying samples.

Resampling techniques can be used to yield more informative estimates of a general statistics.

## Resampling for Estimating Statistics (cont.)

Formally:

- Suppose we want to estimate a parameter  $\theta$  that depends on a random sample set  $X = (x_1, \dots, x_n)$ .
- Assume we have an estimator  $\phi_n(X)$  of  $\theta$  but do not know its distribution.
- Resampling methods try to estimate the bias and variance of  $\phi_n(X)$  using subsamples from  $X$ .

## Jackknife

Let  $PS_i(X)$  be the  $i$ -th pseudo-value of  $\phi_n(X)$ :

$$\begin{aligned} PS_i(X) &= n\phi_n(X) - (n-1)\phi_{n-1}(X_{(i)}) \\ &= \phi_n(X) - \underbrace{(n-1)(\phi_{n-1}(X_{(i)}) - \phi_n(X))}_{\text{bias}_{\text{jack}}} \end{aligned}$$

where  $X_{(i)} = (x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$  is the set without the  $i$ -th element.

Notes:

- $PS_i(X)$  can be interpreted as a bias-corrected version of  $\phi_n(X)$ :
- The bias trend is assumed to be in the estimators from  $\phi_{n-1}(X_{(i)})$  to  $\phi_n(X)$ .

## Jackknife (cont.)

### Jackknife Principle:

- The pseudovalues  $PS_i(X)$  are treated as independent random variables with mean  $\theta$ .
- Using the central limit theorem, the ML estimators for the mean  $\mu_{PS}$  and variance  $\sigma_{PS}^2$  of the pseudovalues are:

$$\mu_{PS} = \frac{1}{n} \sum_{i=1}^n PS_i(X)$$

$$\sigma_{PS}^2 = \frac{1}{n-1} \sum_{i=1}^n (PS_i(X) - \mu_{PS})^2$$

## Jackknife (cont.)

### Example

Estimator for the sample mean:  $\phi_n(X) = \frac{1}{n} \sum_{i=1}^n x_i = \bar{X}$

Pseudovalues of  $\phi_n(X)$ :

$$PS_i(X) = n\bar{X} - (n-1)\bar{X}_{(i)} = x_i$$

Jackknife estimates:

$$\mu_{PS} = \frac{1}{n} \sum_{i=1}^n PS_i(X) = \bar{X}$$

$$\sigma_{PS}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{X})^2$$

## Jackknife (cont.)

### Example

Estimator for sample variance:  $\phi_n(X) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{X})^2$

Pseudovalues of  $\phi_n(X)$ :

$$PS_i(X) = \frac{n}{n-1} (x_i - \bar{X})^2$$

Which implies that:

$$\mu_{PS} = \frac{1}{n} \sum_{i=1}^n PS_i(X) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{X})^2$$

Interestingly:

- $E\{\phi_n(X)\} = \frac{n-1}{n} \sigma^2$  whereas  $E\{\mu_{PS}\} = \sigma^2$
- $\mu_{PS}$  is a bias-corrected version of  $\phi_n(X)$

## Bootstrap

Literary Sidenote:

The term bootstrap comes from the story: *The adventures of Baron Münchhausen*.

- A *bootstrap* data set is created by randomly selecting  $n$  points from the sample set with replacement.
- In *bootstrap estimation* this selection process is independently repeated  $B$  times.
- The  $B$  bootstrap data sets are treated as independent sets.

## Bootstrap (cont.)

The bootstrap estimate of a statistic  $\theta$  and its variance are the mean of the  $B$  estimates  $\hat{\theta}^B$  and its variance:

$$\mu_{BS} = \frac{1}{B} \sum_{i=1}^B \hat{\theta}_i^B$$

$$\sigma_{BS}^2 = \frac{1}{B-1} \sum_{i=1}^B \left( \hat{\theta}_i^B - \mu_{BS} \right)^2$$

The bias is the difference between the bootstrap estimate and the estimator  $\phi_n(X)$ :

$$\text{bias}_{BS} = \mu_{BS} - \phi_n(X)$$

## Bootstrap (cont.)

Properties of the bootstrap estimate:

- Bootstrapping does not change the prior of the data (choose with replacement).
- The larger the number  $B$ , the more will the bootstrap estimate tend towards the true statistic  $\theta$ .
- In contrast, the jackknife estimator requires exactly  $n$  repetitions:
  - less than  $n$  repetitions yield poorer estimates
  - more than  $n$  repetitions merely duplicate information already provided

## Estimating and Comparing Classifiers

Two reasons why we want to know the generalization rate of a classifier on a given problem:

1. to see if the classifier performs well enough to be useful
2. to compare its performance with a competing design

## Cross-Validation

- In **cross-validation**, the training samples are split into two disjoint parts:
  - The first set is the training set used for the traditional training.
  - The second set is the test set used to estimate the classification error.
  - In a second step, both sets are swapped.
  - By that, the classification error can be estimated on the complete data set.
  - Yet training and test set are always disjoint.
- An  **$m$ -fold cross-validation** splits the data into  $m$  disjoint sets of size  $n/m$ :
  - 1 set is used as test set.
  - The other  $m - 1$  sets are used for training.
  - Each set is used once for testing.
- In the **extreme case of  $m = n$** , we have a jackknife estimate of the classification accuracy.

## Cross-Validation (cont.)

The classifier is trained until a minimum validation error is reached (good generalization vs. overfitting):

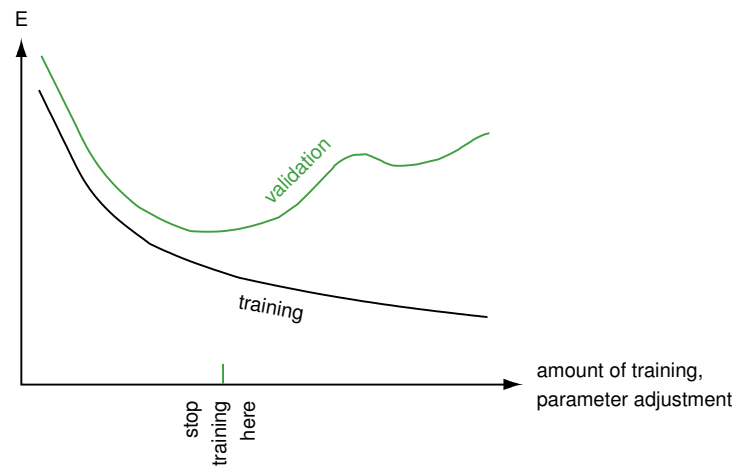


Fig.: The validation error plotted against the amount of training data (adapted from Duda, Hart).

## Lessons Learned

- There is no such thing as a free lunch!
- Bias-variance trade-off
- Jackknife
- Bootstrap
- Cross-Validation



# Next Time in Pattern Recognition



## Further Readings

Examples and various content have been taken from:

- Richard O. Duda, Peter E. Hart, David G. Stork: Pattern Classification, 2nd Edition, John Wiley & Sons, New York, 2000.
- S. Sawyer: Resampling Data: Using a Statistical Jackknife, Washington University, 2005.

Further reading:

- T. Hastie, R. Tibshirani, J. Friedman: The Elements of Statistical Learning, 2nd Edition, Springer, 2009.

## Comprehensive Questions

- What is the meaning of the terms bias and variance?
- What is the difference in bias-variance trade-off between regression and classification?
- How do you estimate the bias and variance of a method?
- What is cross-validation and how can it be used to train a classifier?