

MTAT.03.227 MACHINE LEARNING

**No Free Lunch Theorems and
Statistical Learning Theory**

Sven Laur
University of Tartu

Quick recap of the basic concepts

Let $L(a, b)$ be a loss function that characterises the cost of incorrect guesses. Then the aim of machine learning algorithms is to minimise *risk*

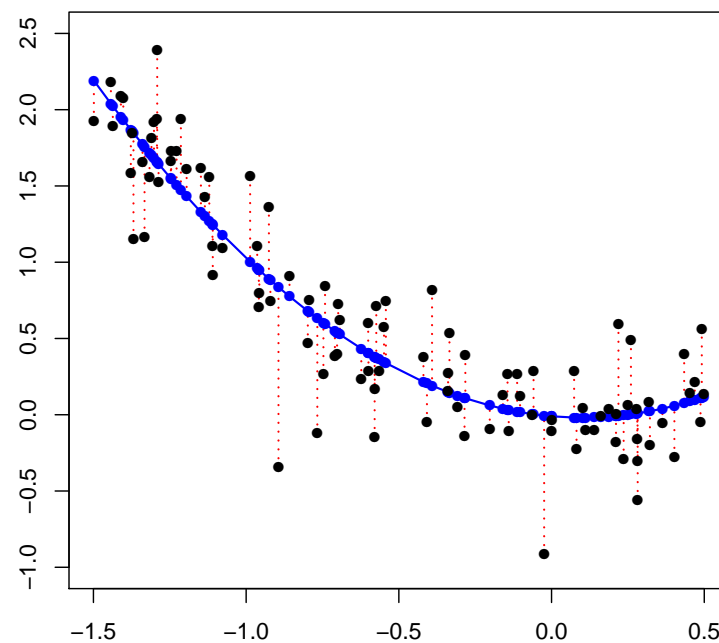
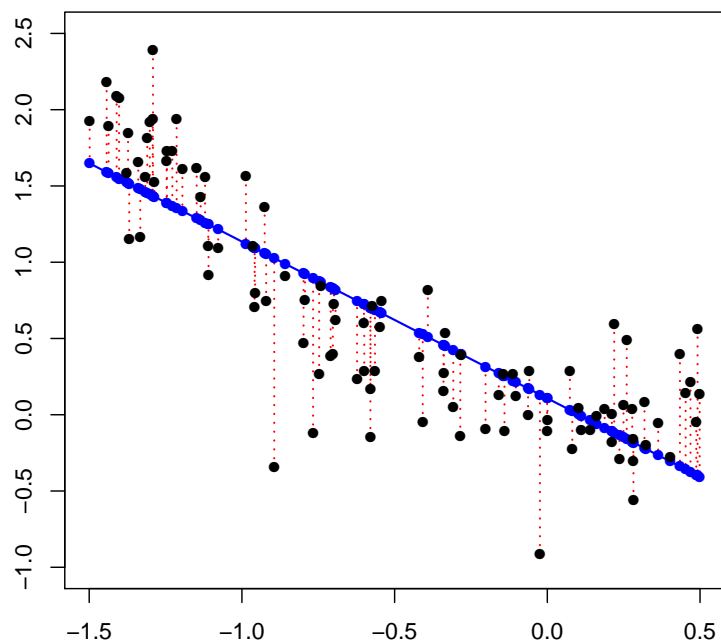
$$R(f) = \mathbf{E}_{\mathbf{x}, y \leftarrow \mathcal{D}} (L(f(\mathbf{x}), y))$$

over an existing but unknown distribution \mathcal{D} . Usually, we seek a solution from a fixed function class \mathcal{H} and try to minimise empirical risk

$$R_m(f) = \frac{1}{m} \cdot \sum_{i=1}^m L(f(\mathbf{x}_i), y_i)$$

where samples are assumed to be independent and identically distributed.

Mean square error as an example



Mean square error for a linear and quadratic model that minimise the empirical risk. Dotted lines visualise individual residues.

Bias-Variance Dilemma

Model bias and consistency

As the set of potential solutions is limited by \mathcal{H} , we might never discover the true function. A *model bias* for \mathcal{H} and true distribution \mathcal{D} is

$$\text{Bias}(\mathcal{H}|\mathcal{D}) = \min_{f \in \mathcal{H}} \mathbf{E}_{\mathbf{x}, y \leftarrow \mathcal{D}} (L(f(\mathbf{x}), y))$$

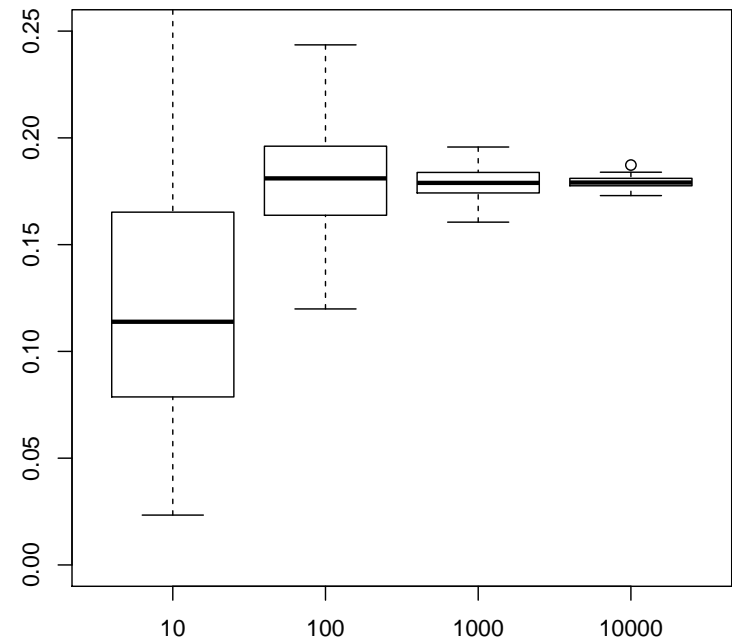
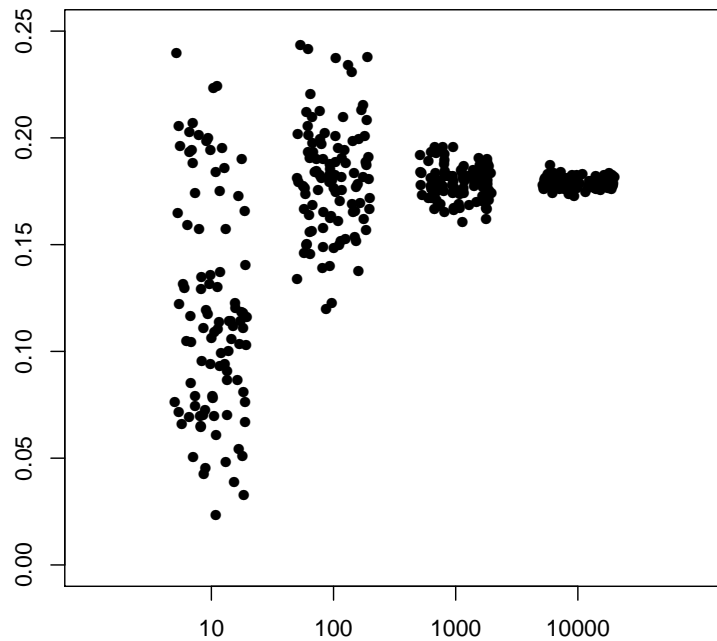
A machine learning algorithm \mathcal{A} is *asymptotically consistent* for a class of data distributions \mathcal{D} if for any data distribution \mathcal{D} from the class \mathcal{D}

$$\mathbf{E}_{\mathcal{D}}(R(f_m)) \xrightarrow{m} \text{Bias}(\mathcal{H}|\mathcal{D})$$

where f_m is the function returned by \mathcal{A} given m samples and \mathcal{H} is the function class from which \mathcal{A} chooses its output.

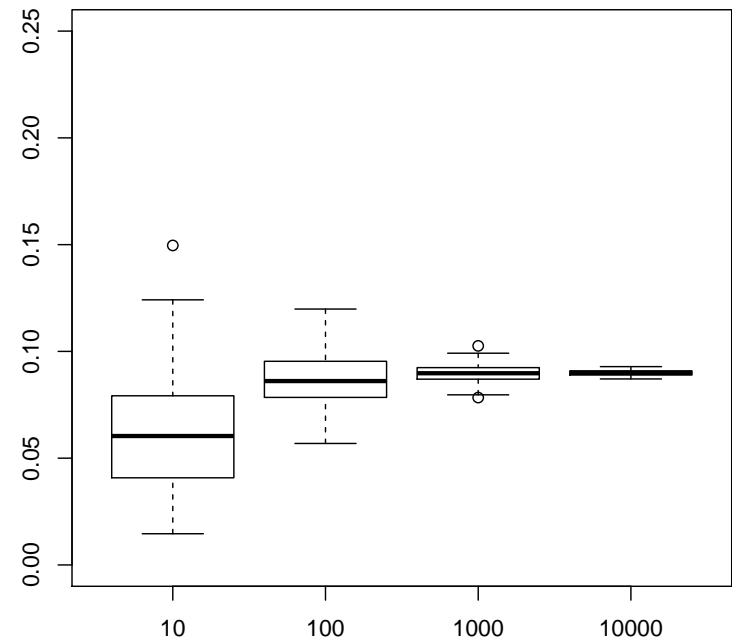
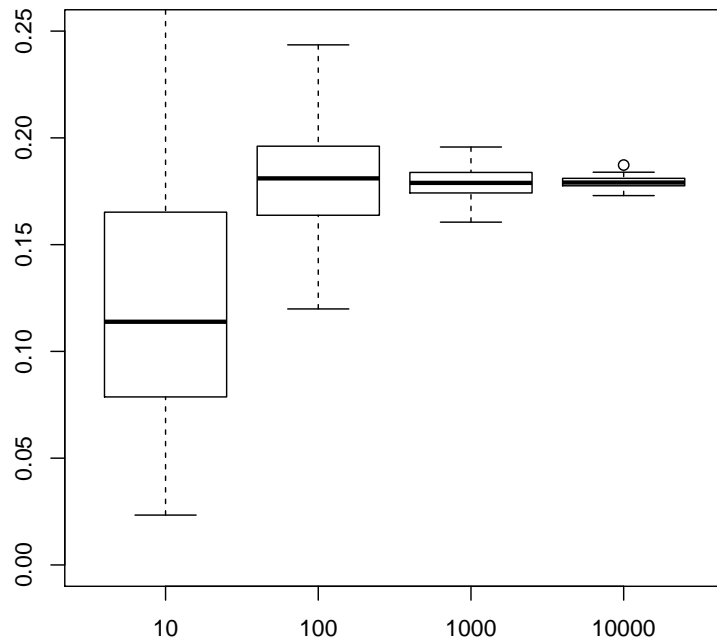
Asymptotic consistency of a linear model

Linear models are known to be asymptotically consistent. Figures depict the empirical risk for the data distribution visualised before.



Different model classes have a different bias

Empirical risk for linear and quadratic models for the data distribution generated by a quadratic model with noise. The error does not go to zero.



Inherent noise in the data

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be the true dependence between the input \mathbf{x} and the output y corrupted by the noise. Let \mathcal{D} be the corresponding data distribution. Then we can consider the risk of the true model

$$R_{\text{true}} = \mathbf{E}_{\mathcal{D}}(L(f(\mathbf{x}), y))$$

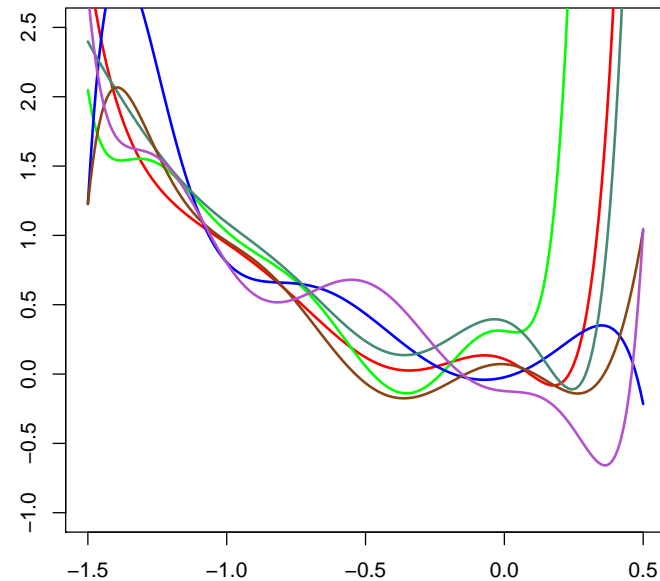
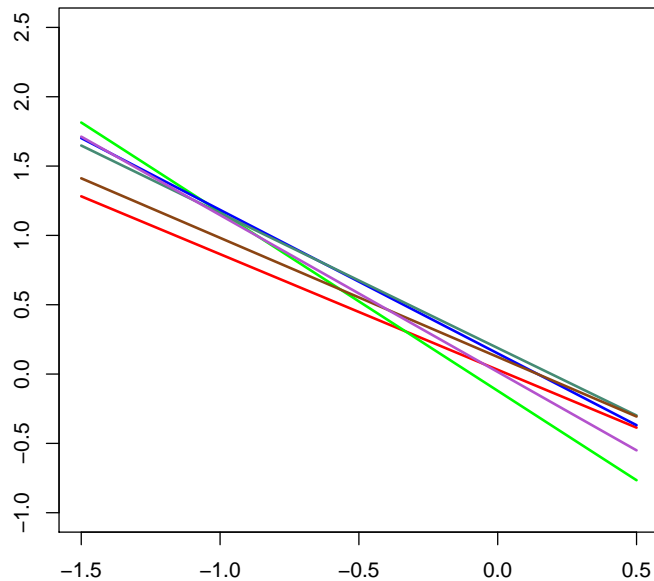
For a quadratic loss function and symmetric error distribution, the minimal achievable risk over all functions and the risk of the true model coincide.

We will refer to this inherent residual risk as *inherent noise*.

If additive error is distributed according to the normal distribution then the inherent noise is equal to the variance of the normal distribution.

Smallest bias is not always the best

Predictions of the model depend on the training set. A method is good if the inferred model is not radically different for different training sets.



Linear model on the left and 8th order polynomial on the right. Models are trained on 20 data points drawn from the quadratic model with the noise.

Bias-variance-noise decomposition

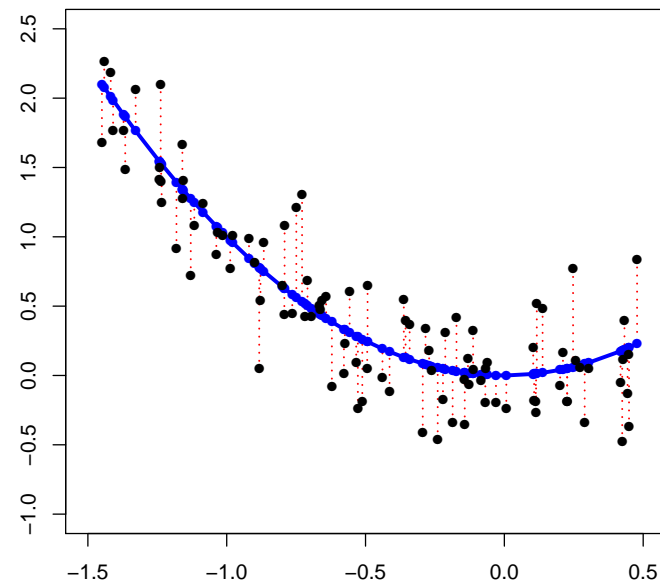
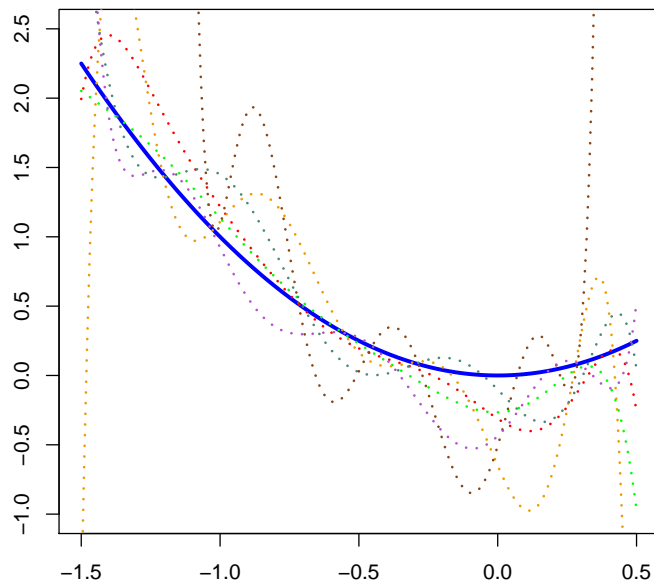
Let \mathcal{S} denote training sample and let $\mathcal{A}_{\mathcal{S}} : \mathbb{R}^n \rightarrow \mathbb{R}$ be the predictor function inferred by the machine learning algorithm \mathcal{A} based on the sample \mathcal{S} .

Assume that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be the true dependence between the input \mathbf{x} and output y that is corrupted by additive noise that is independent from $f(\mathbf{x})$ and has zero mean. Now if we consider quadratic loss function, then the risk can be further decomposed:

$$\begin{aligned} \mathbf{E}_{\mathcal{D}}((\mathcal{A}_{\mathcal{S}}(\mathbf{x}) - y)^2) &= \mathbf{E}_{\mathcal{S}} \left[\mathbf{E}_{\mathbf{x}} \left[\mathbf{E}_y((\mathcal{A}_{\mathcal{S}}(\mathbf{x}) - f(\mathbf{x}) + f(\mathbf{x}) - y)^2) \right] \right] \\ &= \underbrace{\mathbf{E}_{\mathcal{S}, \mathbf{x}}((\mathcal{A}_{\mathcal{S}}(\mathbf{x}) - f(\mathbf{x}))^2)}_{\text{variability caused by } \mathcal{S}} + \underbrace{\mathbf{E}_{(\mathbf{x}, y)}((f(\mathbf{x}) - y)^2)}_{\text{inherent noise}} \end{aligned}$$

Corresponding illustration

The left pane shows the fluctuations of inferred predictors. The right pane shows fluctuations of data points due to the noise.



Bias-variance-noise decomposition

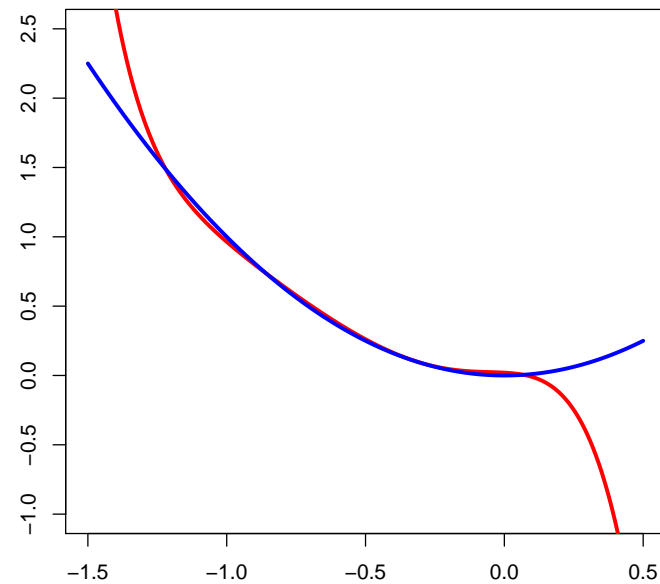
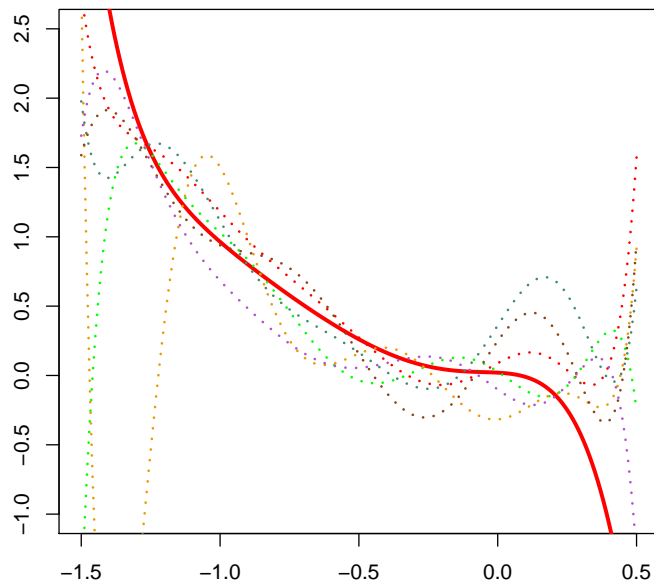
The component estimating the risk caused by the variability of training data can be further decomposed. Let \mathcal{H} be the set of potential prediction functions of the machine learning algorithm \mathcal{A} . Let $f \in \mathcal{H}$ be the optimal function among them. The average of the first term simplifies

$$\begin{aligned}\mathbf{E}_{\mathcal{S}, \mathbf{x}} \left((\mathcal{A}_{\mathcal{S}}(\mathbf{x}) - f(\mathbf{x}))^2 \right) &= \mathbf{E}_{\mathcal{S}, \mathbf{x}} \left((\mathcal{A}_{\mathcal{S}}(\mathbf{x}) - \overline{\mathcal{A}_{\mathcal{S}}}(\mathbf{x}) + \overline{\mathcal{A}_{\mathcal{S}}}(\mathbf{x}) - f(\mathbf{x}))^2 \right) \\ &= \underbrace{\mathbf{E}_{\mathcal{S}, \mathbf{x}} \left((\mathcal{A}_{\mathcal{S}}(\mathbf{x}) - \overline{\mathcal{A}_{\mathcal{S}}}(\mathbf{x}))^2 \right)}_{\text{variance}} + \underbrace{\mathbf{E}_{\mathbf{x}} \left((\overline{\mathcal{A}_{\mathcal{S}}}(\mathbf{x}) - f(\mathbf{x}))^2 \right)}_{\text{bias}}\end{aligned}$$

where $\overline{\mathcal{A}_{\mathcal{S}}}(\mathbf{x})$ is averaged prediction of $\mathcal{A}_{\mathcal{S}}(\mathbf{x})$ over all training sets \mathcal{S} .

Corresponding illustration

The left pane shows the fluctuations of inferred predictors compared to the averaged predictor. The right pane shows the difference between the actual model and the averaged prediction.



Bias-variance trade-off

- ▷ Simple models have small variance but high bias.
- ▷ Complex models have small bias but high variance.
- ▷ We need to balance both components with regularisation.

Minimal training error can be viewed as estimate on model bias

$$E_{tr} = \min_{f \in \mathcal{H}} R_m(f) \approx \min_{f \in \mathcal{H}} R(f) = \text{Bias}(\mathcal{H}|\mathcal{D}) .$$

Hence we somehow need to estimate the variance term in terms of model coefficients count or something similar.

Statistical Learning Theory

General treatment of generalisation error

Sampling bounds for training error

Let us consider finite set of functions $\mathcal{H} = \{f_1, \dots, f_k\}$. The iid assumption allows us to find confidence intervals for $R(f_j)$. They are of type

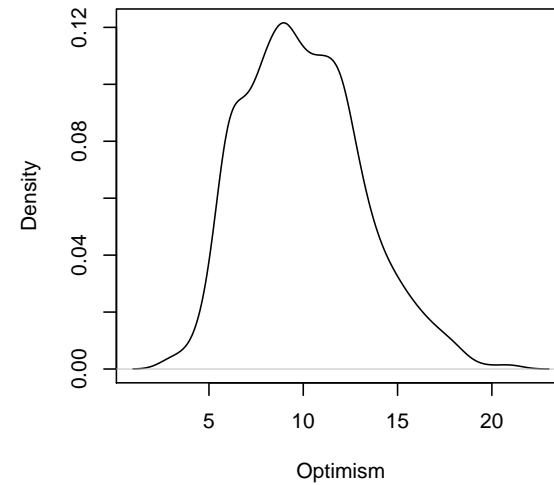
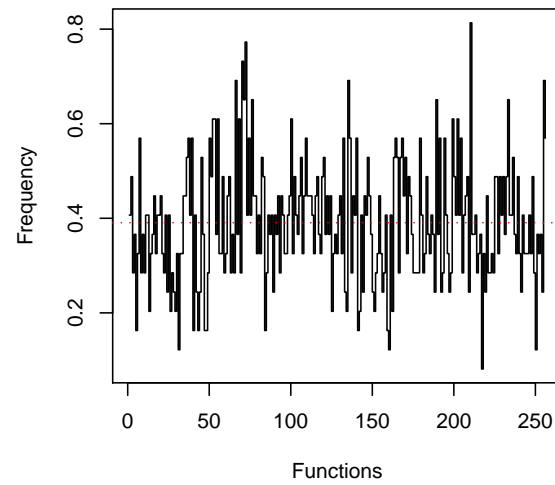
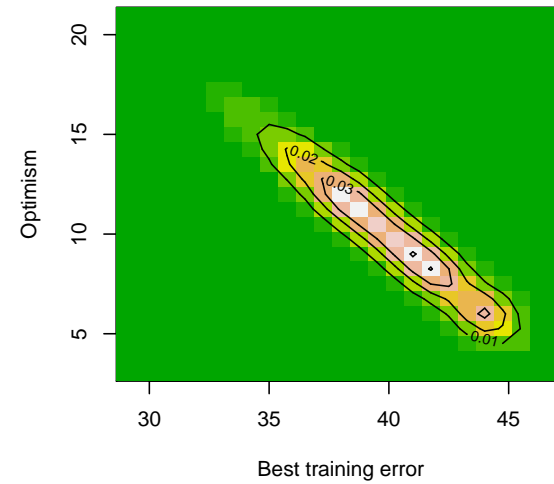
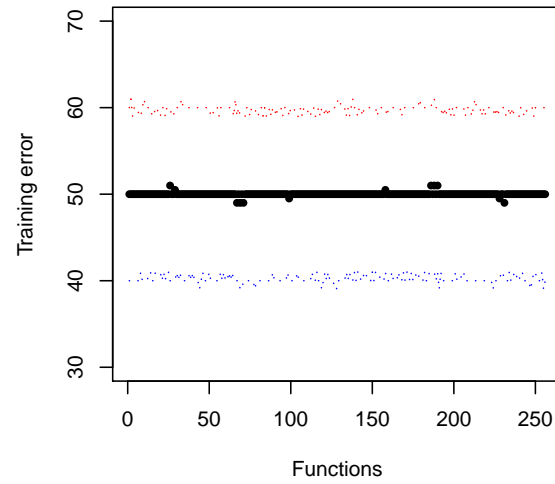
$$\Pr [|R_m(f_j) - R(f_j)| \geq \varepsilon] \leq c \cdot \exp(-\beta m^2 \varepsilon^2) =: \delta$$

for some constants $c, \beta > 0$. They follow from Chernoff, Hoeffding or McDiarmid inequalities. Now applying the union bound we get

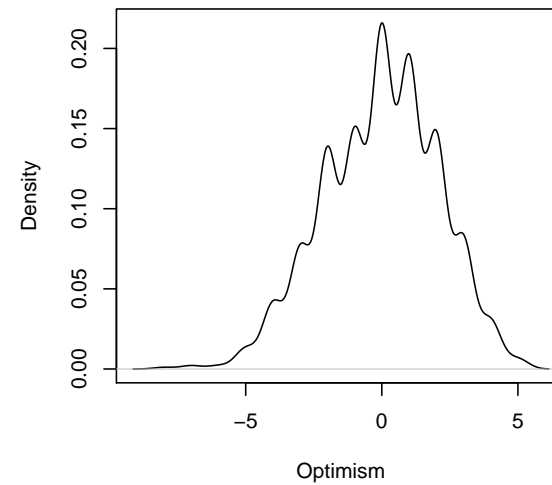
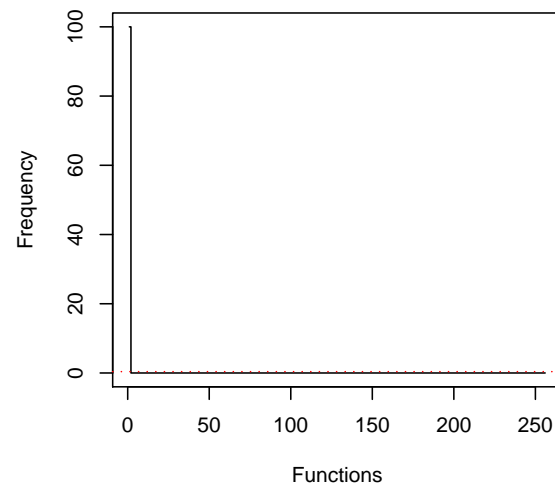
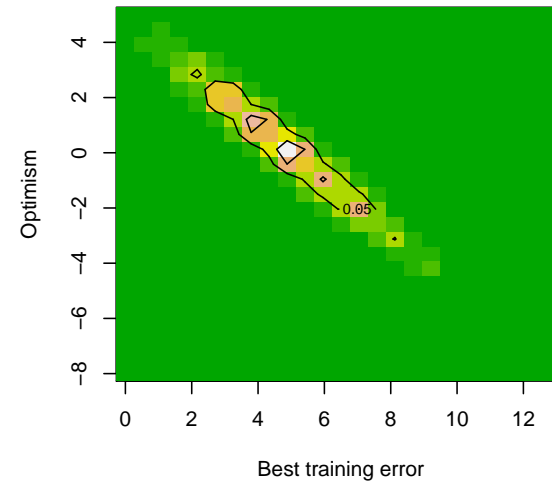
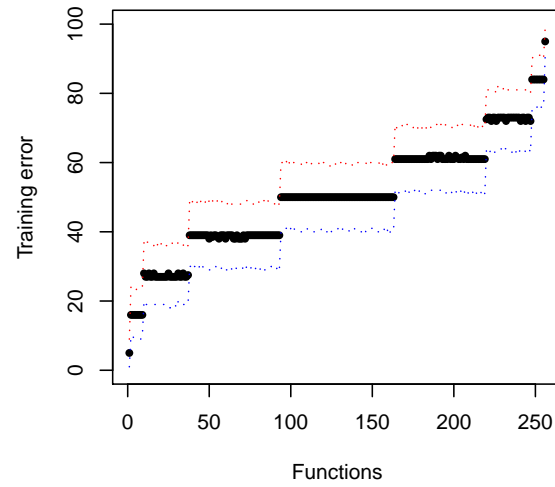
$$\Pr [\exists j : |R_m(f_j) - R(f_j)| \geq \varepsilon] \leq |\mathcal{H}| \cdot \delta .$$

This inequality bounds optimism $\Delta = R_m(f_*) - R(f_*)$ for the proposed solution f_* that minimises training error. Since δ decreases as a function of m the method is asymptotically consistent.

Tightness of inequality. Bad case



Tightness of inequality. Good case



Effective dimension of function families

- ▷ A good function family has few functions that are near optimal.
- ▷ Very similar functions *imply* many near-optimal solutions.

Vapnik-Chervonenkis dimension is one way to characterise flexibility of \mathcal{H}

The *VC dimension* of a function class \mathcal{H} is the largest number of samples d for which one can find $\mathbf{x}_1, \dots, \mathbf{x}_d$ such that any labelling y_1, \dots, y_d can be realised.

If the number of samples $m > d$ then only up to

$$G_{\mathcal{H}}(m) \leq \left(\frac{em}{d}\right)^d$$

labellings can be implemented. For $m > 2d$, \mathcal{H} is really sparse net.

Uniform bound on the optimism

Suppose you have two independent m element sets \mathcal{S}_1 and \mathcal{S}_2 and you optimise classifier on one and test it on other how to estimate optimism.

- ▷ Although the number of functions in \mathcal{H} is potentially infinite, the number of realisable labellings in the $2m$ element data set $\mathcal{S} = \mathcal{S}_1 \cup \mathcal{S}_2$ is

$$G_{\mathcal{H}}(2m) \leq \left(\frac{2em}{d} \right)^d$$

where d is the VC dimension of \mathcal{H} .

- ▷ Each split $\mathcal{S} = \mathcal{S}_1 \cup \mathcal{S}_2$ and labelling ℓ defines empirical errors

$$R_m^1(\mathcal{S}_1, \mathcal{S}_2, \ell) \quad \text{and} \quad R_m^2(\mathcal{S}_1, \mathcal{S}_2, \ell)$$

and in training we choose ℓ that minimises $R_m^1(\ell)$.

Bad splits for each labelling

- ▷ Lets consider worst case $R_m^1(\ell) = 0$ and $R_m^2(\ell) = 1$.
- ▷ For a fixed labelling ℓ , the number of splits where all errors are in the second term is at most 1.
- ▷ As all splits are equiprobable, the probability of such a bad split is

$$1 / \binom{2m}{m} = O(2^{-m/2})$$

- ▷ Thus the probability of large optimism is decreasing exponentially

$$\Pr [\exists f \in \mathcal{H} : R_m^1(f) = 0 \wedge R_m^2(f) = 1] \leq \left(\frac{em}{d}\right)^d O(2^{-\frac{m}{2}}) = O(2^{-\frac{m}{2}})$$

More advanced SLT bound

More careful analysis gives you the bound

$$\Pr \left[\forall f \in \mathcal{H} : R(f) \leq R_m(f) + c(\delta) \sqrt{\frac{\text{VCdim}(\mathcal{H})}{n}} \right] \geq 1 - \delta$$

Alternatives to VC-dimension

There are many alternatives to VC-dimension:

▷ Rademacher complexity

$$\mathcal{R}_m(\mathcal{H}|\mathcal{S}) = \mathbf{E} \left[\sup_{f \in \mathcal{H}} \frac{2}{m} \cdot \sum_{i=1}^m \sigma_i f(\mathbf{x}_i) \right] \quad \text{for } \sigma_i \in \{-1, 1\}$$
$$\mathcal{R}_m(\mathcal{H}) = \mathbf{E}_{\mathcal{S} \leftarrow \mathcal{D}} (\mathcal{R}_m(\mathcal{H}|\mathcal{S}))$$

▷ Covering numbers

They are important since sometimes we cannot compute VC-dimension while the other complexity measures do exist and are easy to compute.

NFL Theorems

Why do we need NFL theorems?

No Free Lunch theorems show the following important facts

- ◇ **There are no universal solutions.** All machine learning methods are *equally good* if we do not place *strong* assumptions on the input data.
- ◇ **Benchmarks can lie:** For each machine learning algorithm there exists a *sample or sample class* where it outperforms some other method.
- ◇ **All clever bounds have fine print conditions:** Sampling bounds and SLT results are *correct* but do not tell *exactly* what we need and hope.

We state NFL theorems in maximally simple way to show their banal nature.

Setup for the NFL theorem

Consider a finite learning problem.

- ▷ Let $\mathcal{S} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$ be the training set.
- ▷ Let $\mathcal{T} = \{(\mathbf{x}_{m+1}, y_{m+1}), \dots, (\mathbf{x}_n, y_n)\}$ be the test data.
- ▷ Assume that the inputs of the training and test data are *disjoint*.
- ▷ Let $\mathcal{A}_{\mathcal{S}}(\mathbf{x})$ prediction of machine learning method on \mathbf{x} .

How large is the test error (*out-of-training sample error*)?

$$R_{ot}(\mathcal{A}) = \frac{1}{n - m} \cdot \sum_{i=m+1}^n [\mathcal{A}_{\mathcal{S}}(\mathbf{x}_i) \neq y_i]$$

The simplest version of NFL theorem

Theorem 1. If we assume that labels y_i are computed as $f(\mathbf{x}_i)$ for a function chosen uniformly from all possible functions then

$$\mathbf{E}_f (R_{ot}(\mathcal{A}|f)) = \frac{1}{2} .$$

Proof

▷ Let $f^c(\mathbf{x}_{m+i}) = \neg f(\mathbf{x}_{m+i})$ denote the complement of f over \mathcal{T} . Then

$$R_{ot}(\mathcal{A}|f) + R_{ot}(\mathcal{A}|f^c) = 1$$

▷ As f and f° have equal weight $\Pr[f|\mathcal{S}] = \Pr[f^\circ|\mathcal{S}]$, we get

$$\mathbf{E}_f(R_{ot}(\mathcal{A}|f)) = \frac{1}{2} \cdot \sum_{f:X \rightarrow Y} \Pr[f|\mathcal{S}] (R_{to}(\mathcal{A}|f) + R_{ot}(\mathcal{A}|f^c)) = \frac{1}{2}$$

Obvious conclusions

The claim can be extended to different loss functions L and function sets \mathcal{F} as long as the symmetry condition on the test set is preserved:

- ▷ We can define complement operator such that

$$\forall f \in \mathcal{F} : \exists f^\circ \in \mathcal{F} \wedge \forall f \in \mathcal{F} : (f^\circ)^\circ = f$$

- ▷ Corresponding losses are constant for any prediction

$$L(y, f(\mathbf{x}_{m+i})) + L(y, f^\circ(\mathbf{x}_{m+i})) = g(\mathbf{x}_{m+i})$$

- ▷ Probabilities in complement pairs are equal

$$\Pr [f|\mathcal{S}] = \Pr [f^\circ|\mathcal{S}]$$

NFL theorem for problem classes

- ▷ In real life tasks, some problems are harder than others.
- ▷ Let \mathcal{F} denote probability distribution over all functions
- ▷ Which problem instances \mathcal{F} the algorithm \mathcal{A} can handle?

Theorem II. Consider a uniform distribution over all possible problem classes. Then

$$\mathbf{E}_{\mathcal{F}} \left(\mathbf{E}_{f \leftarrow \mathcal{F}} (R_{ot}(\mathcal{A})) \right) = \frac{1}{2} .$$

As a result, no algorithm can be better than the other—there are some problem classes for which the other algorithm is better.

Concentration equation for NFL

By combining Theorem II and Law of large numbers we obtain the following surprising result about benchmarking.

Theorem III. Let $\mathcal{F}_1, \dots, \mathcal{F}_k$ be the set of problem classes chosen uniformly over the design space. Then

$$\frac{1}{k} \cdot \sum_{j=1}^k \mathbf{E}_{f \leftarrow \mathcal{F}_j} (R_{ot}(\mathcal{A})) \approx \frac{1}{2}$$

for any learning algorithm \mathcal{A} and thus for any two learning algorithms \mathcal{A} and \mathcal{B} there exist exist a particular example problem

$$\mathbf{E}_{f \leftarrow \mathcal{F}_j} (R_{ot}(\mathcal{A})) \lesssim \mathbf{E}_{f \leftarrow \mathcal{F}_j} (R_{ot}(\mathcal{B})) \ .$$

Wait SLT results contradict NFL theorems!

NFL theorem says that you cannot learn arbitrary function. On the same time, SLT bounds assure that $R_m(f) + \Delta$ is rather accurate estimate.

What happens then?

- ▷ We can assume that n large enough so that $R(f) \approx R_{ot}(f)$.
- ▷ For most functions $R_m(f) + \Delta \approx 0.5$ and we do not care.
- ▷ For few functions $R_m(f) + \Delta \approx 0.0$
 - ◇ For half of these functions $R_{ot}(f) \approx 0.0$ and we are happy.
 - ◇ For half of these functions $R_{ot}(f) \approx 1.0$ and we are happy.

Small print notice. For fixed error estimate SLT bounds say nothing.

- ▷ Rejection of uninteresting results creates bias!
- ▷ We must additionally assume that the learning task is easy.