

MLRF Lecture 06

J. Chazalon, LRDE/EPITA, 2021

More theory on ML

Lecture 05 part 03

What is our goal?

*Given **samples** (described by features) and **true labels**,
find a **good** function
which will correctly **predict labels**
given **new data samples***

Problems:

- Which family for our function?
- What is “good”?
- How to train / find such function?

Let us step back a little bit.

What are the sources of error?

Figures from Bradski & Kaehler,
Learning OpenCV, O'Reilly 2008

Noise

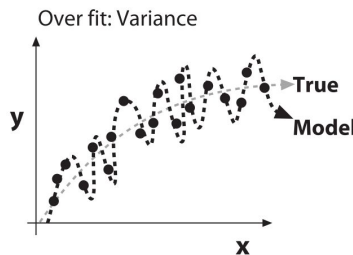
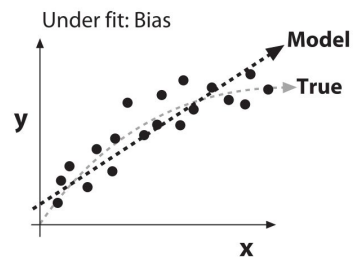
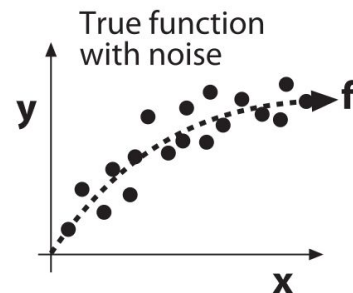
- Your data is not perfect. (or “*Every model is wrong.*”)
- Even if there exist an optimal underlying model, the observations are corrupted by noise (e.g. multiple y for a given x).
- **Even the optimal solution could be wrong.**

Bias

- You need to simplify to generalize.
- Your classifier needs to drop some information about the training set to have generalization power.
- **The set of solutions explored does not contain the optimal solution.**

Variance

- You have many ways to explain your training dataset.
- It is hard to find an optimal solution among those many possibilities.
- **If we draw another training set from the same distribution, we would obtain another solution.**

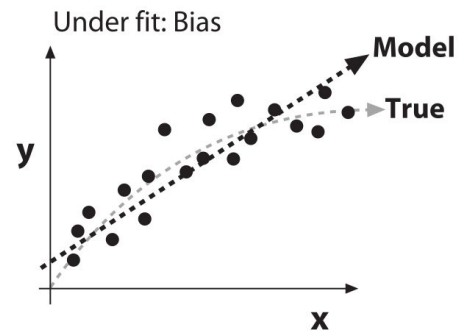


Two big issues

Figures from Bradski & Kaehler,
Learning OpenCV, O'Reilly 2008

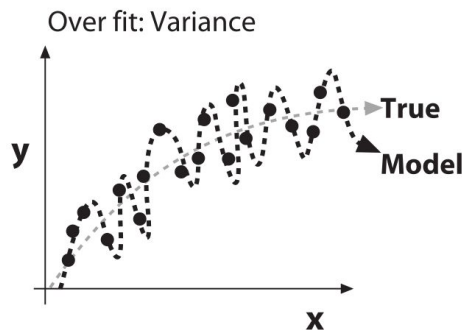
Under-fitting

- Caused by **bias**
- Your model assumptions are too strong for the data, so the model won't fit well.



Over-fitting

- Caused by **variance**
- Your algorithm has memorized the data *including* the noise, so it can't generalize.



(plus **performance limit**: due to noise and model capacity)

The theory

Bias (statistical definition)

From [Wikipedia](#)

Let T be a statistic used to estimate a parameter θ .

If $E[T] = \theta + \text{bias}(\theta)$ then $\text{bias}(\theta)$ is called the bias of the statistic T , where $E[T]$ represents the expected value of the statistics T .

If $\text{bias}(\theta) = 0$, then $E[T] = \theta$. So, T is an unbiased estimator of the true parameter, say θ .

Expected Risk

Let \mathbf{D}_n be a training set of examples \mathbf{z}_i , drawn independently from an unknown distribution $p(\mathbf{z})$

We need a set of functions \mathbf{F} . Example: linear functions $\mathbf{f}(\mathbf{x}) = \mathbf{a} \cdot \mathbf{x} + b$

We need a loss function $L(\mathbf{z}, \mathbf{f})$. Example: $L((\mathbf{x}, \mathbf{y}), \mathbf{f}) = (\mathbf{f}(\mathbf{x}) - \mathbf{y})^2$

The **Expected Risk**, i.e. the expected generalization error, is:

$$R(f) = E_Z[L(z, f)] = \int_Z L(z, f)p(z)dz$$

But we do not know $p(\mathbf{z})$, and we cannot test all \mathbf{z} !

Empirical Risk

Because we **cannot measure** the real **Expected Risk**, we have to **estimate it** using the **Empirical Risk**:

$$\hat{R}(f, D_n) = \frac{1}{n} \sum_{i=1}^n L(z_i, f)$$

Integral to sum...

(D_n is our dataset)

And our training procedure then relies on **Empirical Risk Minimization** (ERM):

$$f^*(D_n) = \arg \min_{f \in \mathcal{F}} \hat{R}(f, D_n)$$

Pick the best on the dataset.

And the training error is given by:

$$\hat{R}(f^*(D_n), D_n)$$

Empirical risk of the best found, evaluated on our dataset.

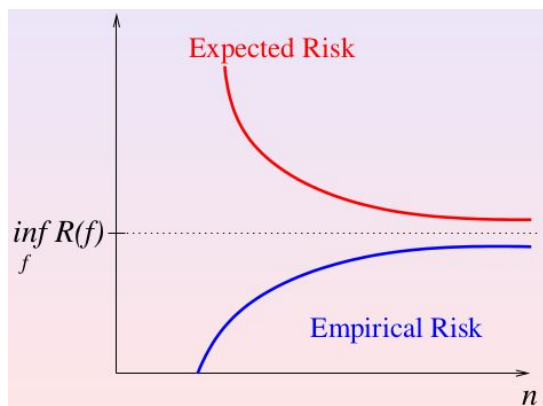
Does this make sense?

The empirical risk is an unbiased estimate of the risk, i.e. the more test samples we have, the more accurate our estimate is, **under iid assumption**.

$$\begin{aligned} E[\hat{R}(f, D)] &= E \left[\frac{1}{n} \sum_{i=1}^n L(f, Z_i) \right] \\ &= \frac{1}{n} \sum_{i=1}^n E[L(f, Z_i)], && Z_i \text{ s are independent} \\ &= \frac{1}{n} \sum_{i=1}^n E[L(f, Z)], && Z_i \text{ s are identically distributed} \\ &= \frac{1}{n} n E[L(f, Z)] \\ &= R(f) \end{aligned}$$

Estimate the Expected Risk with the Empirical Risk

For a given *capacity*, using more samples to train and evaluate your predictor **should** make your Empirical Risk converge toward the best possible Expected Risk, if the ERM is consistent for \mathcal{F} , given your training set \mathcal{D}_n .



Fixed h

But the training risk is biased

The **training** error is a biased estimate of the risk, i.e. the solution $f^\star(D_n)$ found by minimizing the training error is better on D_n than on any other set D'_n drawn from $p(z)$.

$$E \left[R(f^\star(D_n)) - \hat{R}(f^\star(D_n), D_n) \right] \geq 0$$

Expected, on training set

Empirical, on training set

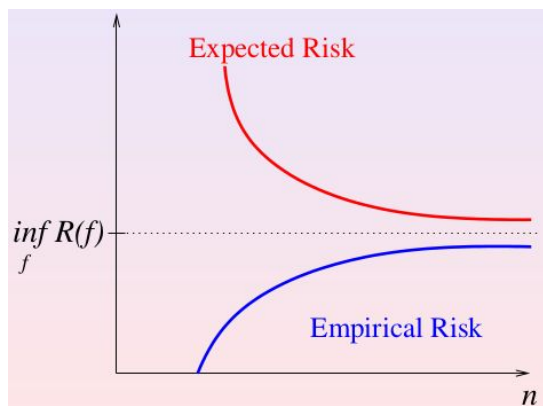
However, under certain assumptions, the difference between the expected and the empirical risks can be bounded. This is an important result from the work of Vapnik [Vapnik (2000). [The nature of statistical learning theory](#). Springer.].

Note that the empirical risk on the test set is an unbiased estimate of the risk.

$$E \left[\hat{R}(f^\star(D_{\text{train}}), D_{\text{test}}) \right] = R(f^\star(D_{\text{train}}))$$

Estimate the Expected Risk with the Empirical Risk

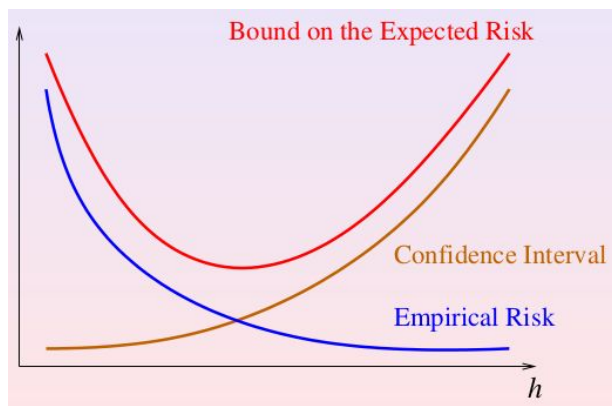
For a given *capacity*, using more samples to train and evaluate your predictor **should** make your Empirical Risk converge toward the best possible Expected Risk, if the ERM is consistent for F , given your training set D_n .



Fixed h

The difference between Expected Risk and Empirical Risk is **bounded** but depends on the **capacity** of F (set of possible functions).


There is an **optimal** capacity for a given number of training samples n .



Fixed n

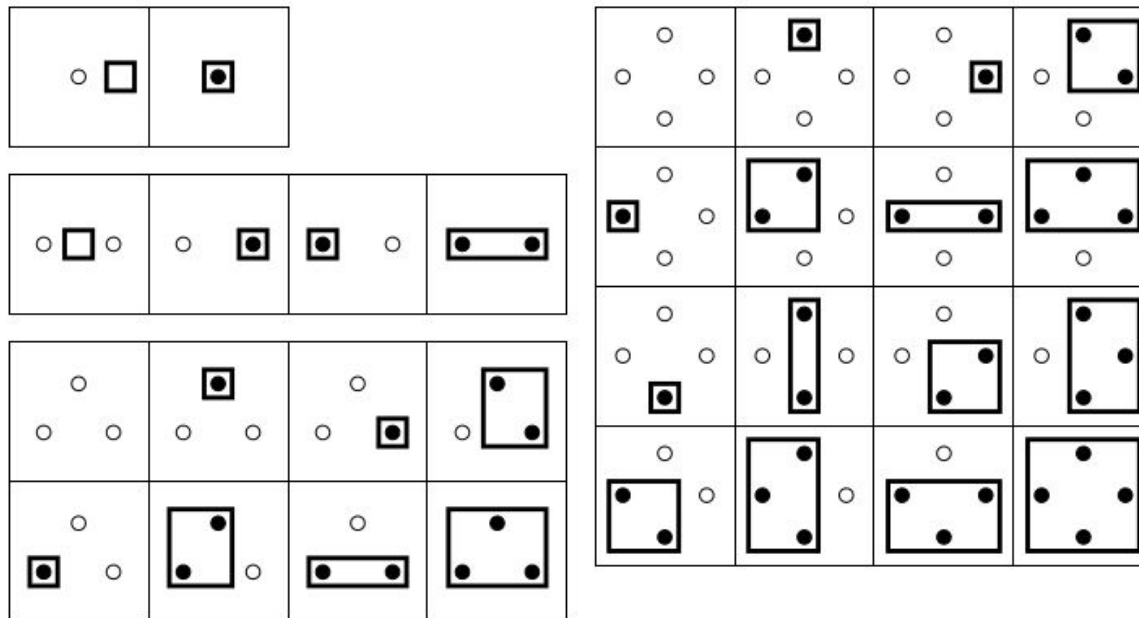
Capacity

The capacity $h(\mathbf{F})$ is a measure of its size, or complexity (or VC dimension).

For classification, the capacity of \mathbf{F} is defined by Vapnik & Chervonenkis as the largest n  such that there exist a set of examples \mathbf{D}_n such that one can always find an $\mathbf{f} \in \mathbf{F}$ which gives the correct answer for all examples in \mathbf{D}_n , for any possible labeling.

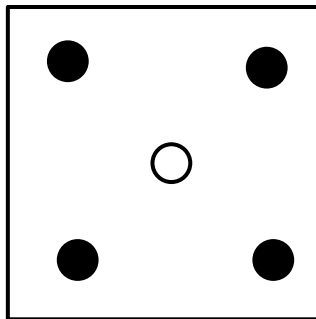
Capacity

Consider for F the characteristic functions of rectangles. We can find families of 1, 2, 3 or 4 points which can be labelled arbitrarily:



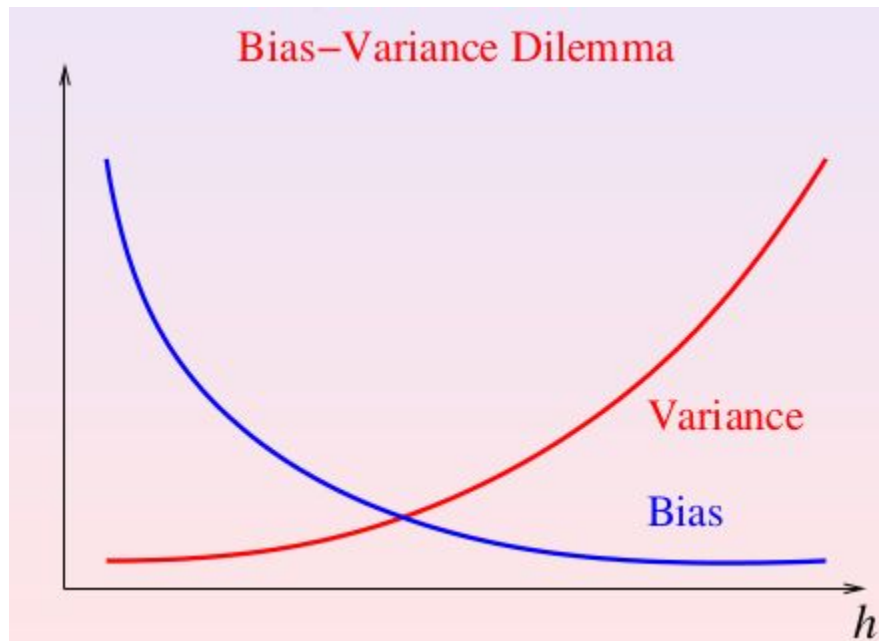
Capacity

However, given a family of 5 points, if the four external points are labelled 1 and the center point labelled 0, than no function from \mathcal{F} can predict that labelling.
Hence here $D = 4$.

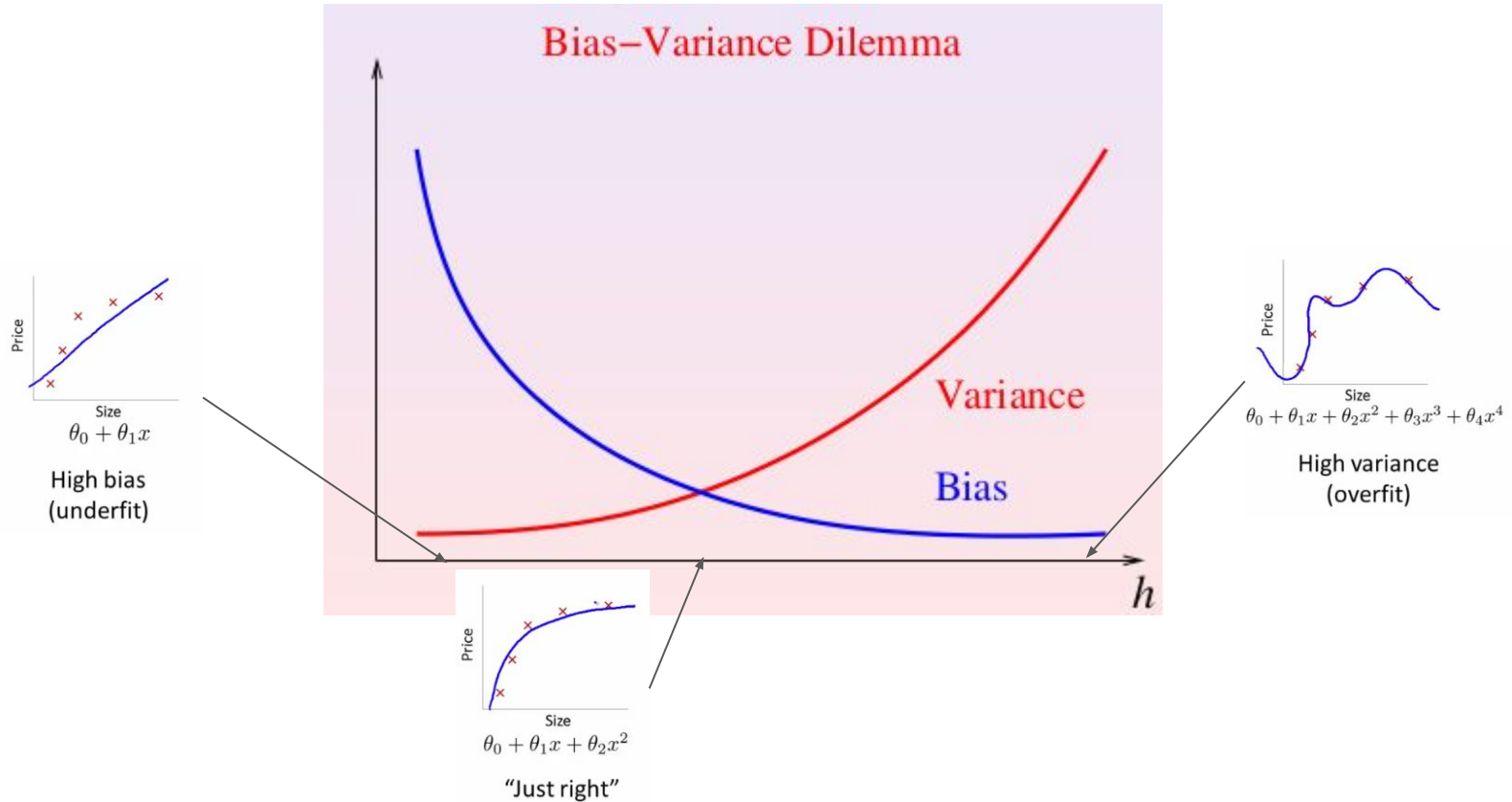


The Bias-Variance Dilemma

Intrinsic dilemma: when the capacity $h(\mathcal{F})$ grows, the bias goes down, but the variance goes up!



Look for an optimal balance



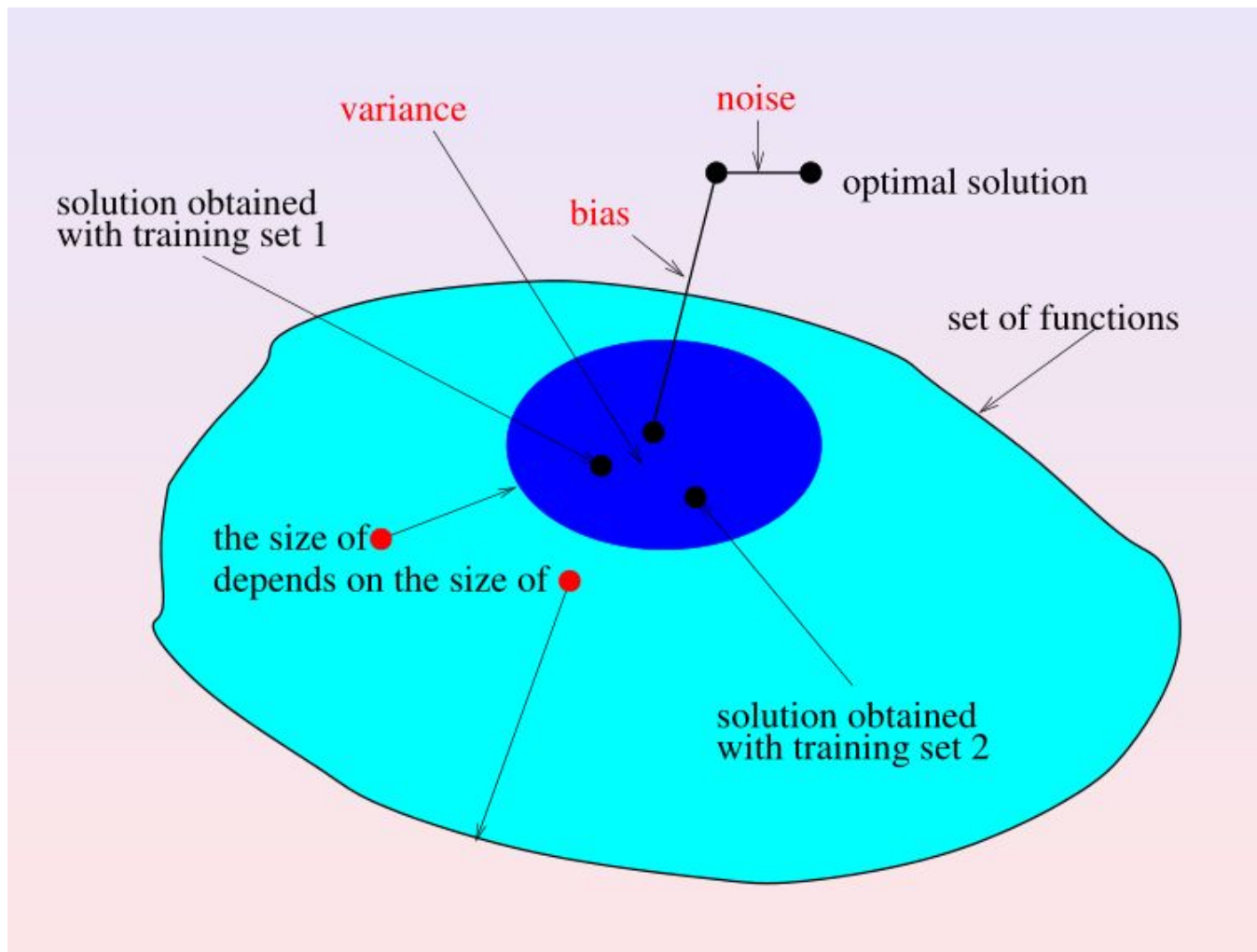
Decomposing the bias-variance-error for MSE

For a regression problem with a mean square loss, we have the following decomposition. Let $\mathbf{Y} = \mathbf{f}(\mathbf{X}) + \epsilon$, with $\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma_\epsilon^2)$ and $f_D(\mathbf{X})$ an estimator of $\mathbf{f}(\mathbf{X})$, learned over the training set \mathbf{D} . The error at a particular point $\mathbf{X} = \mathbf{x}_0$ is:

$$\begin{aligned} \text{Err}(x_0) &= E[(Y - f_D(x_0))^2 | X = x_0] \\ &= E[(Y - E[f_D(x_0)] + E[f_D(x_0)] - f_D(x_0))^2 | X = x_0] \\ &= E[(Y - E[f_D(x_0)])^2 | X = x_0] + E[(E[f_D(x_0)] - f_D(x_0))^2] \\ &\quad + 2 \cdot E[(Y - E[f_D(x_0)]) \cdot (E[f_D(x_0)] - f_D(x_0)) | X = x_0] \end{aligned}$$

Given that
 $E[f_D(x_0)] - f_D(x_0) = 0$

$$\begin{aligned} \text{Err}(x_0) &= E[(Y - E[f_D(x_0)])^2 | X = x_0] + E[(E[f_D(x_0)] - f_D(x_0))^2] \\ &= E[(f(x_0) + \epsilon - E[f_D(x_0)])^2] + \text{Var}(f_D(x_0)) \\ &= E[(f(x_0) - E[f_D(x_0)])^2] + E[\epsilon^2] + \text{Var}[f_D(x_0)] \\ &= [\text{Bias}[f_D(x_0)]]^2 + \sigma_{\text{epsilon}}^2 + \text{Var}[f_D(x_0)] \end{aligned}$$



In practice

In practice: Empirical Risk and Expected Risk

Measure train and test error.

Use hold-out sets, cross-validations, etc. to get a test error.

Train error: Empirical Risk.

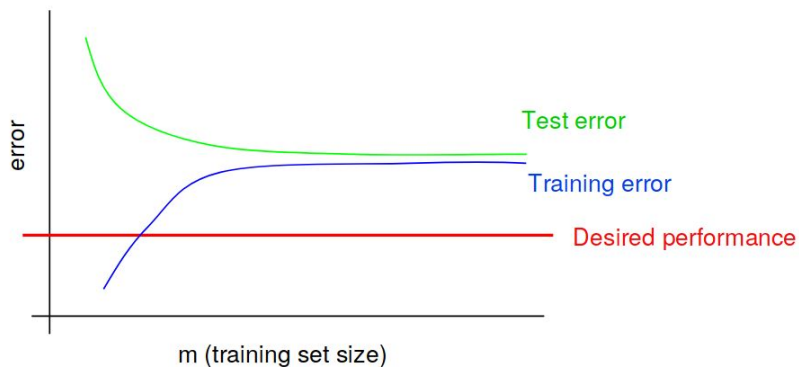
⇒ *Can my model learn something (by heart)?*

Test error: Coarse estimate of the Expected Risk.

⇒ *Can my model generalize to unseen data?*

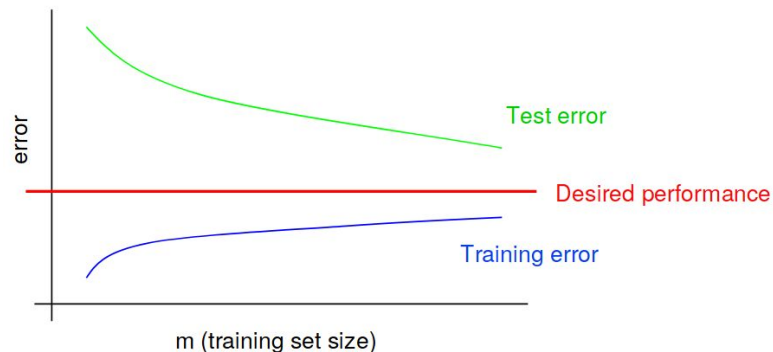
Detect under-fitting and over-fitting

High bias: This learning curve shows high error on both the training and test sets, so the algorithm is suffering from high bias.



Even training error is unacceptably high.
Small gap between training and test error.

High variance: This learning curve shows a large gap between training and test set errors, so the algorithm is suffering from high variance.



Test error still decreasing as m increases.
Suggests larger training set will help.
Large gap between training and test error.

Some solutions / hints

From Bradski & Kaehler, Learning OpenCV,
O'Reilly 2008 ↓

Problem	Possible Solutions
Bias	<ul style="list-style-type: none"> • More features can help make a better fit. • Use a more powerful algorithm.
Variance	<ul style="list-style-type: none"> • More training data can help smooth the model. • Fewer features can reduce overfitting. • Use a less powerful algorithm.
Good test/train, bad real world	<ul style="list-style-type: none"> • Collect a more realistic set of data.
Model can't learn test or train	<ul style="list-style-type: none"> • Redesign features to better capture invariance in the data. • Collect new, more relevant data. • Use a more powerful algorithm.

From C. Aggarwal, Data Mining: The Textbook,
Springer 2015 →

Technique	Source/Level of Bias	Source/Level of Variance
Simple Models	Oversimplification increases bias in decision boundary.	Low variance. Simple models do not overfit.
Complex Models	Generally lower than simple models. Complex boundary can be modeled.	High variance. Complex assumptions will be overly sensitive to data variation.
Shallow Decision Trees	High bias. Shallow tree will ignore many relevant split predicates.	Low variance. The top split levels do not depend on minor data variations.
Deep Decision Trees	Lower bias than shallow decision tree. Deep levels model complex boundary.	High variance because of overfitting at lower levels.
Rules	Bias increases with fewer antecedents per rule.	Variance increases with more antecedents per rule.
Naive Bayes	High bias from simplified model (e.g., Bernoulli) and naive assumption.	Variance in estimation of model parameters. More parameters increase variance.
Linear Models	High bias. Correct boundary may not be linear.	Low variance. Linear separator can be modeled robustly.
Kernel SVM	Bias lower than linear SVM. Choice of kernel function.	Variance higher than linear SVM.
k-NN Model	Simplified distance function such as Euclidean causes bias. Increases with k .	Complex distance function such as local discriminant causes variance. Decreases with k .
Regularization	Increases bias	Reduces variance

How to get started?

1. Get enough data in the right format from your customer (*hard*)
2. Check and split data (*boring but mandatory*)
3. Agree on a loss function and minimum performance goal (*moderate*)
4. Try to overfit a predictor on some samples (train set loss), increase complexity only if needed (*capacity check*)
5. Fit on more data (*more = better*)
6. Check for overfitting (val set loss) and add regularization if needed
7. Evaluate performance thoroughly (test set loss) (*reports, identify failure cases, etc.*)
8. Do some hyper-parameter optimization, try other models...\$
9. ...