

**STATG019 – Selected Topics in Statistics 2018**

# **Lecture 2**

**Theory and Methodology for  
Estimating the Generalization Error**

Dr Franz J. Király

# The statistical supervised learning setting

**Given data** from generative, unknown RV  $(X, Y)$  t.v.in  $\mathcal{X} \times \mathcal{Y}$  (not independent of each other in general)

features
labels
domains

where we observe

$$(X_1, Y_1), \dots, (X_N, Y_N) \underset{\text{i.i.d.}}{\sim} (X, Y)$$

## Estimate/learn

a prediction functional  $f$  t.v.in  $[\mathcal{X} \rightarrow \mathcal{Y}]$   
 (via an estimator using only the data  $(X_i, Y_i)$ )

**Such that** the expected generalization error

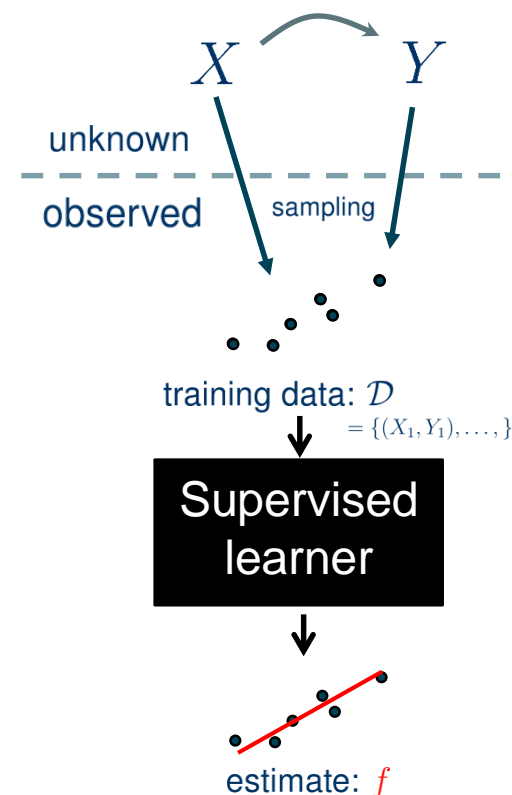
$$\varepsilon(f) := \mathbb{E} [L(f(X), Y)] \text{ is small}$$

where  $L : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  is a chosen convex loss

i.e.,  $[\hat{y} \mapsto L(\hat{y}, y)]$  is convex for all  $y \in \mathcal{Y}$

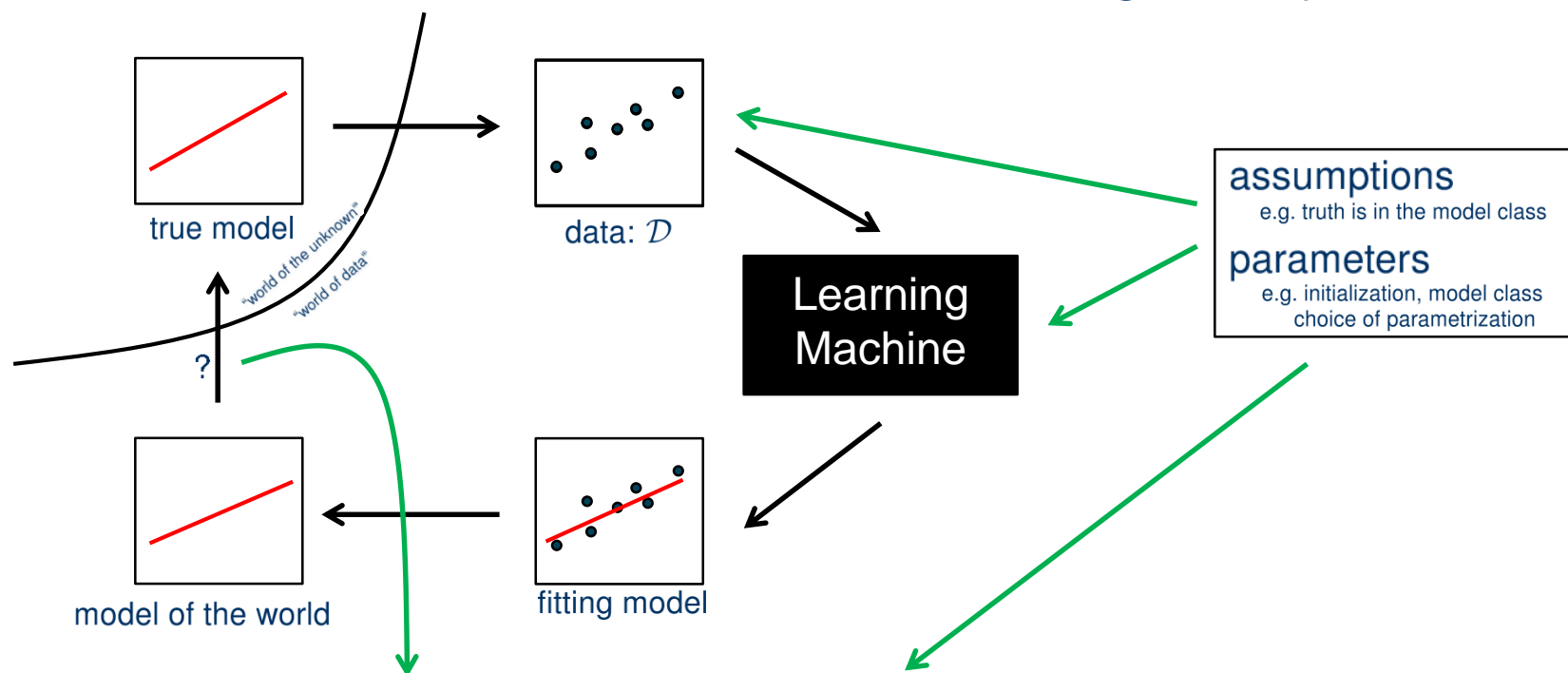
$$\text{e.g., } L : (\hat{y}, y) \mapsto (\hat{y} - y)^2$$

**Validity:** (i) under assumptions, prove that  $f$  has low error (model specific)  
 (ii) external estimate of  $f$ 's error with guarantees (model agnostic)



## (i) model-specific validity arguments

Classical statistics and statistical learning theory

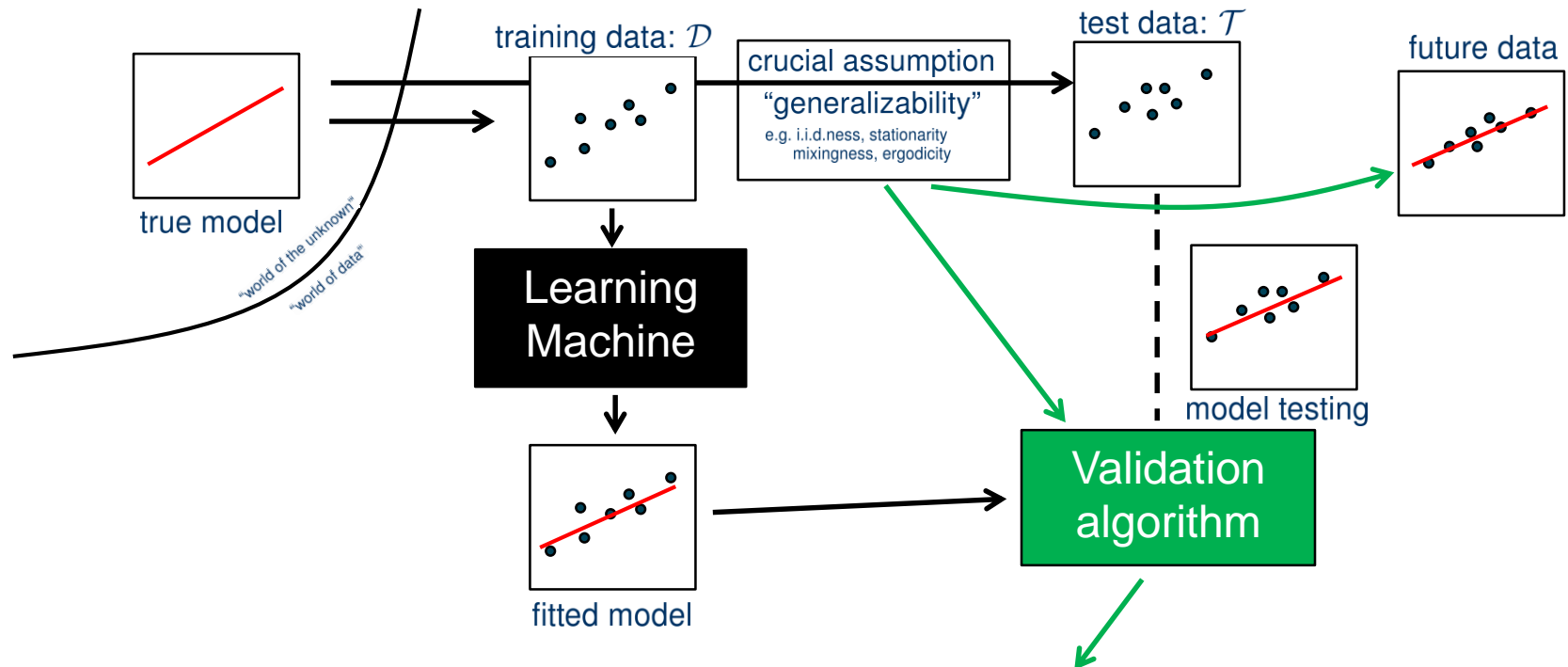


Guarantees implied by *assumptions on data/ and model properties*  
in the form of mathematical/statistical theorems

**Validity argument incomplete unless these are checked!**

## (ii) model-agnostic validity arguments

„model validation“, „model checking“, „model testing“



Guarantees implied by *empirical results and properties of the task*  
in the form results plus theorems about the validation algorithm

**Guarantees always hold and allow comparison**

*but may be weaker than model-specific ones if theory available*

# Table of contents for lecture 2

## Mathematical guarantees for model-agnostic validation

### Generic guarantees for generalization & external validity

- Finite sample properties of mean and variance

- Central limit theorems for mean and variance

- Estimating the generalization loss of a prediction strategy

### Supervised learning as function estimation

- Bias-variance trade-off for supervised learning

- Estimators of the generalization loss and the loss's variance

- Variance reduction through re-sampling

### Comparison of models and pairwise testing

- Prediction functionals versus prediction strategies

- Re-sampling estimates of black-box variance

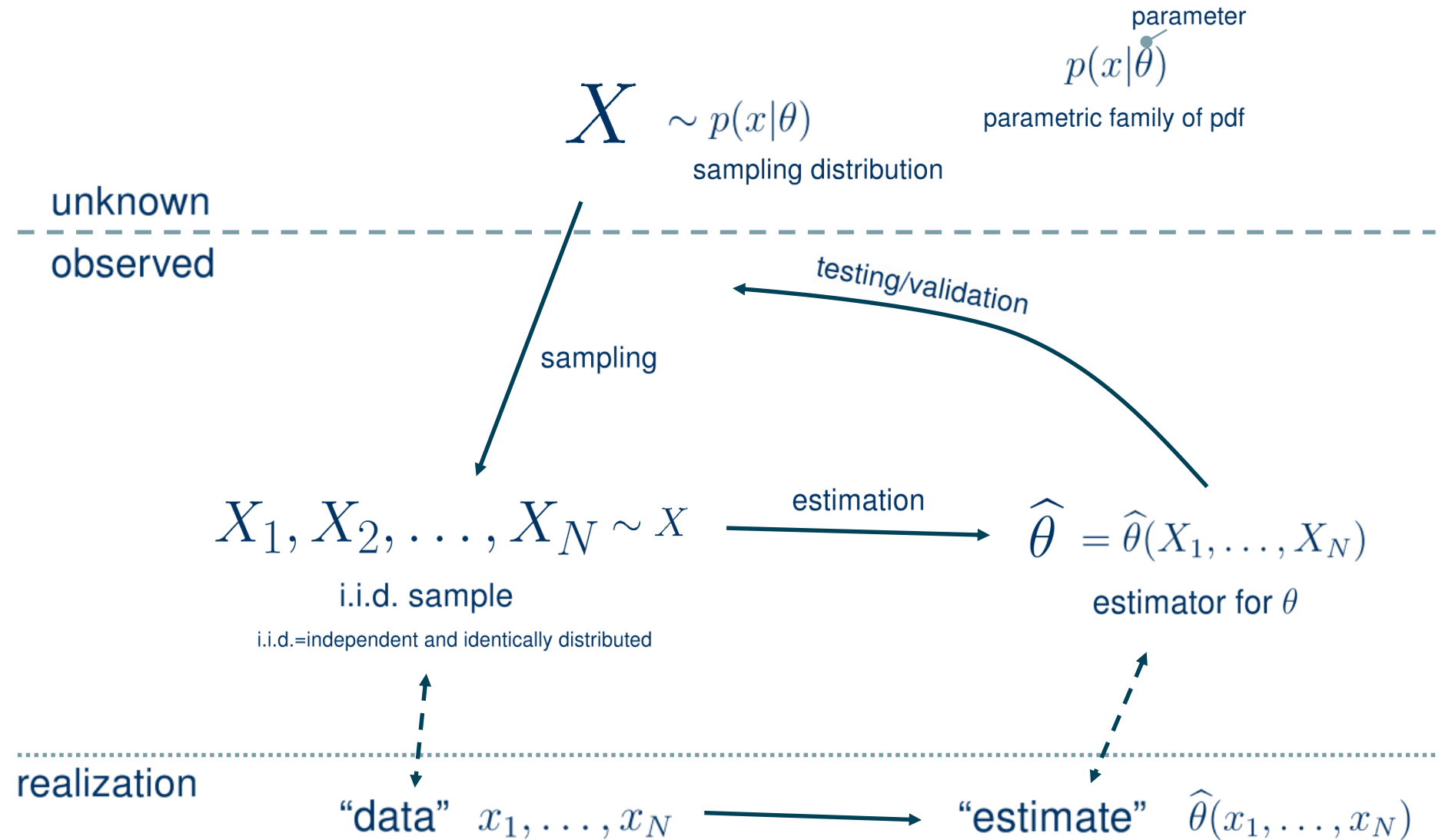
- Hypothesis testing for pairwise and portmanteau model comparison

- Interactions of testing and re-sampling

### Very briefly: overview of generic model-specific guarantees

# Generalization guarantees

# Reminder: Parametric Estimation



# The parametric estimation setting

Data is generated as:

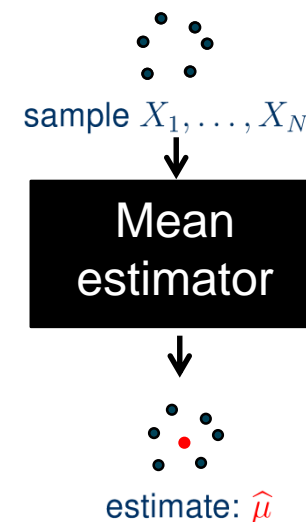
true parameter  $\theta \in \mathbb{R}^m$  (e.g.  $\theta = (\mu, \Sigma)$  for Gaussians)

$$X_1, \dots, X_N \underset{\text{i.i.d.}}{\sim} X \sim p(\cdot | \theta)$$

(parametric) estimator  $\hat{\theta} = \hat{\theta}(X_1, \dots, X_N)$

e.g. mean estimator  $\hat{\mu} = \frac{1}{N} \sum_{i=1}^N X_i$

covariance estimator  $\hat{\Sigma} = \frac{1}{N} \sum_{i=1}^N (X_i - \hat{\mu})(X_i - \hat{\mu})^\top$



**Accurate terminology:**  $\hat{\theta}(X_1, \dots, X_N)$  is „the estimate“ (a RV t.v.in  $\mathbb{R}^m$ )

$\hat{\theta} : (x_1, \dots, x_N) \mapsto \hat{\theta}(x_1, \dots, x_N)$  is „the estimator“ (an algorithm)  
function in  $[\mathcal{X}^n \rightarrow \mathbb{R}^m]$

often confounded, but distinction will become important in later section



# Guarantees: LLN, CLT and CI

hold in general, nonparametric setting  $X_1, \dots, X_N \underset{\text{i.i.d.}}{\sim} X$  t.v.in  $\mathbb{R}$  (assume finite moments)

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N X_i \text{ estimates } \mathbb{E}[X] \quad \hat{\Sigma} = \frac{1}{N-1} \sum_{i=1}^N (X_i - \hat{\mu})^2 \text{ estimates } \text{Var}[X]$$

$$\frac{\hat{\Sigma}}{N} \text{ estimates } \text{Var}[\hat{\mu}] \quad \text{„variance of the sample mean“}$$

**Theorems:** (i)  $\mathbb{E}[\hat{\mu}] = \mathbb{E}[X]$   $\text{Var}[\hat{\mu}] = \frac{1}{N} \text{Var}[X]$  (weak LLN implied by this & Chebyshev)

(statements  
you should  
have seen!)

(ii)  $\sqrt{N} (\hat{\mu} - \mathbb{E}[X]) \xrightarrow{d} \mathcal{N}(0, \text{Var}(X))$  as  $N \rightarrow \infty$  (CLT as usually stated)

(iii)  $\mathbb{E}[\hat{\Sigma}] = \text{Var}[X]$   $\mathbb{E}[\hat{\Sigma}/N] = \text{Var}[\hat{\mu}]$  (unbiasedness of estimates)

(statements  
you probably  
haven't seen  
but are simple  
computations)

(iv)  $\text{Var}[\hat{\Sigma}] = \frac{1}{N} \cdot [M_4 - \frac{N-3}{N-1} \text{Var}(X)^2]$  where  $M_4 := \mathbb{E}[(X - \mathbb{E}(X))^4]$  (LLN for sample variance)

(v)  $\sqrt{N} (\hat{\Sigma} - \text{Var}[X]) \xrightarrow{d} \mathcal{N}(0, M_4 - \text{Var}(X)^2)$  as  $N \rightarrow \infty$   
 $= \text{Var}[(X - \mathbb{E}[X])^2]$  (CLT for sample variance)

(stylized)  
Consequence:  $\left[ \hat{\mu} + \Phi^{-1}(\alpha/2) \cdot \sqrt{\hat{\Sigma}/N}, \hat{\mu} - \Phi^{-1}(\alpha/2) \cdot \sqrt{\hat{\Sigma}/N} \right]$  is „good“ CI

i.e.,  $\alpha$ -CI for  $\mu$  with good coverage probability as  $N > 50$

Important exception:  
imbalanced & binary  $X$

# Estimation of the generalization error

**Setting:** i.i.d. *test data*  $(X_1, Y_1), \dots, (X_M, Y_M) \underset{\text{i.i.d.}}{\sim} (X, Y)$  t.v.in  $\mathcal{X} \times \mathcal{Y}$

**prediction functional**  $f : \mathcal{X} \rightarrow \mathcal{Y}$  e.g.,  $\mathcal{X} = \mathbb{R}^n$  and  $\mathcal{Y} = \mathbb{R}$

**loss function**  $L : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  e.g.,  $L : (\hat{y}, y) \mapsto (\hat{y} - y)^2$

**To estimate:**  $\varepsilon(f) = \mathbb{E}[L(f(X), Y)]$  *expected generalization loss*

**Theorems suggest following estimators:**

$$\hat{\varepsilon}(f) := \frac{1}{M} \sum_{i=1}^M L(f(X_i), Y_i) \quad \text{Observation: } L_i := L(f(X_i), Y_i) \text{ are i.i.d.}$$

„empirical loss“ since pairs  $(X_i, Y_i)$  are i.i.d.

$$\hat{v}(f) := \frac{1}{M-1} \sum_{i=1}^M (L_i - \hat{\varepsilon})^2 \quad \text{„sample variance of the empirical losses“}$$

$$\text{Confidence interval: } \left[ \hat{\varepsilon}(f) + \Phi^{-1}(\alpha/2) \cdot \sqrt{\hat{v}(f)/M}, \hat{\varepsilon}(f) - \Phi^{-1}(\alpha/2) \cdot \sqrt{\hat{v}(f)/M} \right]$$

The end ... ?

**Big problem:** this is only valid if  $f$  is fixed (non-random), e.g., already trained/fitted!

Otherwise  $L_i$  are dependent through a random  $f$ . (*no guarantees for strategies!*)

*But statements & guarantees about the trained prediction functionals are correct!*

# **Bias and Variance in Parametric Estimation and Supervised Learning**

# Supervised learning as function estimation

(more restrictive formulation of the task due to stronger structural assumptions)

**Given data** from generative, unknown RV  $(X, Y)$  t.v.in  $\mathcal{X} \times \mathcal{Y}$  (not independent of each other in general)

features      labels      domains

where we observe

$$(X_1, Y_1), \dots, (X_N, Y_N) \underset{\text{i.i.d.}}{\sim} (X, Y)$$

## Parametric supervised assumption:

There is a „true“ labelling process  $f = f_\theta$   
(usually, one assumes an „additive error model“, that is:)

$$Y_i = f(X_i) + \varepsilon_i \quad \varepsilon_i \text{ is error with } \mathbb{E}[\varepsilon_i] = 0$$

(errors assumed independent)

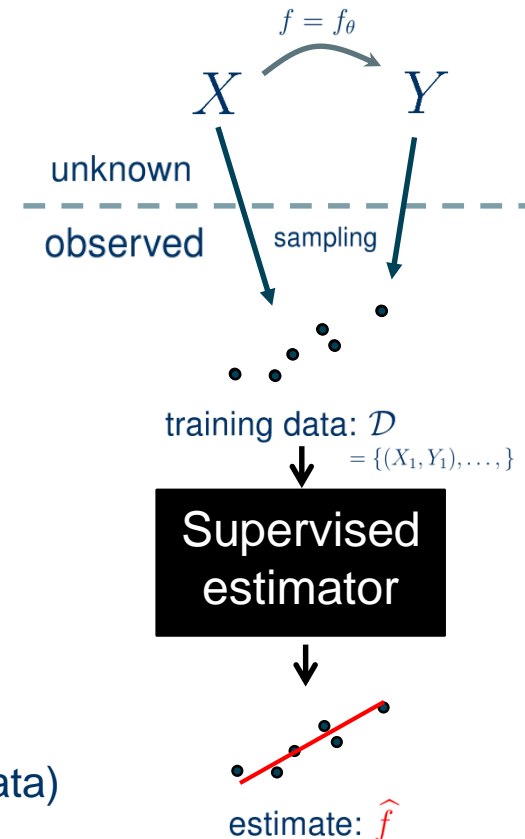
Example (Linear Regression):

$$f(x) = \beta^\top x + \alpha, \quad \theta = (\alpha, \beta)$$

## Supervised estimation task:

Learn a good approximation  $\hat{f} = f_{\hat{\theta}}$  (using the training data)

such that  $\varepsilon(\hat{f}) = \mathbb{E} \left[ L(\hat{f}(X), Y) \right]$  is small



# How good is an estimator?

## Bias and variance

parametric estimation:

true parameter  $\theta \in \mathbb{R}^m$  (e.g.  $\theta = (\mu, \Sigma)$  for Gaussians)

estimator  $\hat{\theta} = \hat{\theta}(X_1, \dots, X_N)$

bias of  $\hat{\theta}$ :

$$\text{Bias}(\hat{\theta}) = \mathbb{E}[\hat{\theta} - \theta]$$

measures expected deviation of the mean

variance of  $\hat{\theta}$ :

$$\text{Var}(\hat{\theta}) = \mathbb{E} \left[ (\hat{\theta} - \mathbb{E}[\hat{\theta}])^2 \right]$$

measures scatter around estimator mean

MSE of  $\hat{\theta}$ :

$$\text{MSE}(\hat{\theta}) = \mathbb{E} \left[ (\hat{\theta} - \theta)^2 \right]$$

measures total estimation error

Convention: for  $x \in \mathbb{R}^m$ , denote  $x^2 = x^\top x$

**Cave:** all these quantities may depend on  $\theta$

(yes, this includes the variance. Think why.)

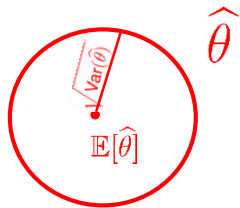
# Bias and variance

$$\text{Bias}(\hat{\theta}) = \mathbb{E}[\hat{\theta} - \theta]$$

measures expected deviation of the mean

$$\text{Var}(\hat{\theta}) = \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2]$$

measures scatter around estimator mean



low bias



high bias



low variance



high variance

# Bias-variance trade-off for estimators

$$\text{Bias}(\hat{\theta}) = \mathbb{E}[\hat{\theta} - \theta]$$

on measures expected deviation

$$\text{Var}(\hat{\theta}) = \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2]$$

measures scatter around estimator mean

$$\text{MSE}(\hat{\theta}) = \mathbb{E}[(\hat{\theta} - \theta)^2]$$

measures total estimation error

**Proposition:**  $\text{MSE}(\hat{\theta}) = \text{Bias}(\hat{\theta})^2 + \text{Var}(\hat{\theta})$

**proof:**  $\text{MSE}(\hat{\theta}) = \mathbb{E}[(\hat{\theta} - \theta)^2] = \mathbb{E}[\hat{\theta}^2] - 2\theta\mathbb{E}[\hat{\theta}] + \theta^2$

$$= \mathbb{E}[\hat{\theta}^2] - 2(\mathbb{E}[\hat{\theta}])^2 + (\mathbb{E}[\hat{\theta}])^2 + (\mathbb{E}[\hat{\theta}])^2 - 2\theta\mathbb{E}[\hat{\theta}] + \theta^2$$

$$= \mathbb{E}[\hat{\theta}^2] - 2\mathbb{E}[\hat{\theta}(\mathbb{E}[\hat{\theta}])] + \mathbb{E}(\mathbb{E}[\hat{\theta}])^2 + (\mathbb{E}[\hat{\theta}] - \theta)^2$$

$$= \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2] + (\mathbb{E}[\hat{\theta}] - \theta)^2$$

$$= \text{Var}(\hat{\theta}) + \text{Bias}(\hat{\theta})^2$$

stays also valid for taking  $xx^\top$  instead of  $x^\top x$

# Example: mean of Gaussian

$$\text{Bias}(\hat{\theta}) = \mathbb{E}[\hat{\theta} - \theta] \quad \text{Var}(\hat{\theta}) = \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2] \quad \text{MSE}(\hat{\theta}) = \mathbb{E}[(\hat{\theta} - \theta)^2]$$

$$\text{MSE}(\hat{\theta}) = \text{Bias}(\hat{\theta})^2 + \text{Var}(\hat{\theta})$$

## Example: estimation of mean

$X_1, \dots, X_N$  i.i.d Gaussian  $\sim \mathcal{N}(\mu, \sigma^2)$



$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N X_i$$

“natural” estimator

$$\hat{\mu}_s = 42$$

“universal” estimator



$$\text{Bias}(\hat{\mu}) = 0 \quad \text{MSE}(\hat{\mu}) = \text{Var}(\hat{\mu}) = \frac{\sigma^2}{N}$$

$$\text{Bias}(\hat{\mu}_s) = 42 - \theta \quad \text{Var}(\hat{\mu}_s) = 0 \quad \text{MSE}(\hat{\mu}_s) = (42 - \theta)^2$$



# Supervised learning as function estimation

**Given data** from generative, unknown RV  $(X, Y)$  t.v.in  $\mathcal{X} \times \mathcal{Y}$  (not independent of each other in general)  
 where we observe  
 features labels domains

$$(X_1, Y_1), \dots, (X_N, Y_N) \underset{\text{i.i.d.}}{\sim} (X, Y)$$

## Parametric supervised assumption:

There is a „true“ labelling process  $f = f_\theta$   
 (usually, one assumes an „additive error model“, that is:)

$$Y_i = f(X_i) + \epsilon_i \quad \epsilon_i \text{ is error with } \mathbb{E}[\epsilon_i] = 0$$

(errors assumed independent)

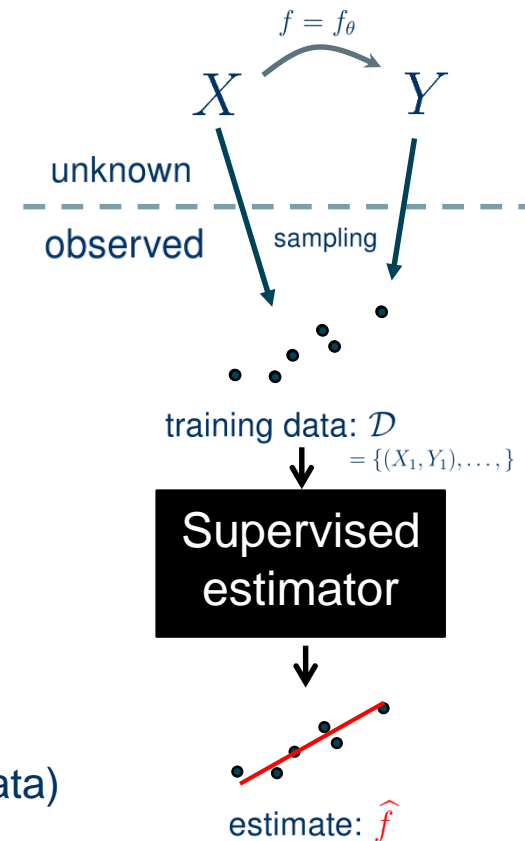
Example (Linear Regression):

$$f(x) = \beta^\top x + \alpha, \quad \theta = (\alpha, \beta)$$

## Supervised estimation task:

Learn a good approximation  $\hat{f} = f_{\hat{\theta}}$  (using the training data)

such that  $\varepsilon(\hat{f}) = \mathbb{E} \left[ L(\hat{f}(X), Y) \right]$  is small



# Bias and variance in supervised regression

assume  $\mathcal{Y} = \mathbb{R}$

test point  $X_*$ , test label  $Y_*$  where  $(X_*, Y_*) \sim (X, Y)$

parameter  $\theta$  “is” a generative function  $f = f_\theta$ :

$$Y_* = f(X_*) + \epsilon_*$$

$\hat{f} = f_{\hat{\theta}}$  learnt prediction rule (possibly depending on seen training data)

$$\begin{aligned} \text{bias of } \hat{f} \text{ at } X_*: \quad \text{Bias}(\hat{f}|X_*) &= \mathbb{E}_{Y|X_*}[\hat{f}(X_*) - f(X_*)|X_*] \\ &= \mathbb{E}_{Y|X_*}[\hat{f}(X_*) - Y_*|X_*] \end{aligned}$$

$$\begin{aligned} \text{variance of } \hat{f} \text{ at } X_*: \quad \text{Var}(\hat{f}|X_*) &= \text{Var}_{Y|X_*} \left[ \hat{f}(X_*)|X_* \right] = \mathbb{E}_{Y|X_*} \left[ (\hat{f}(X_*)^2|X_*) \right] \\ &\quad - \left( \mathbb{E}_{Y|X_*} [\hat{f}(X_*)|X_*] \right)^2 \end{aligned}$$

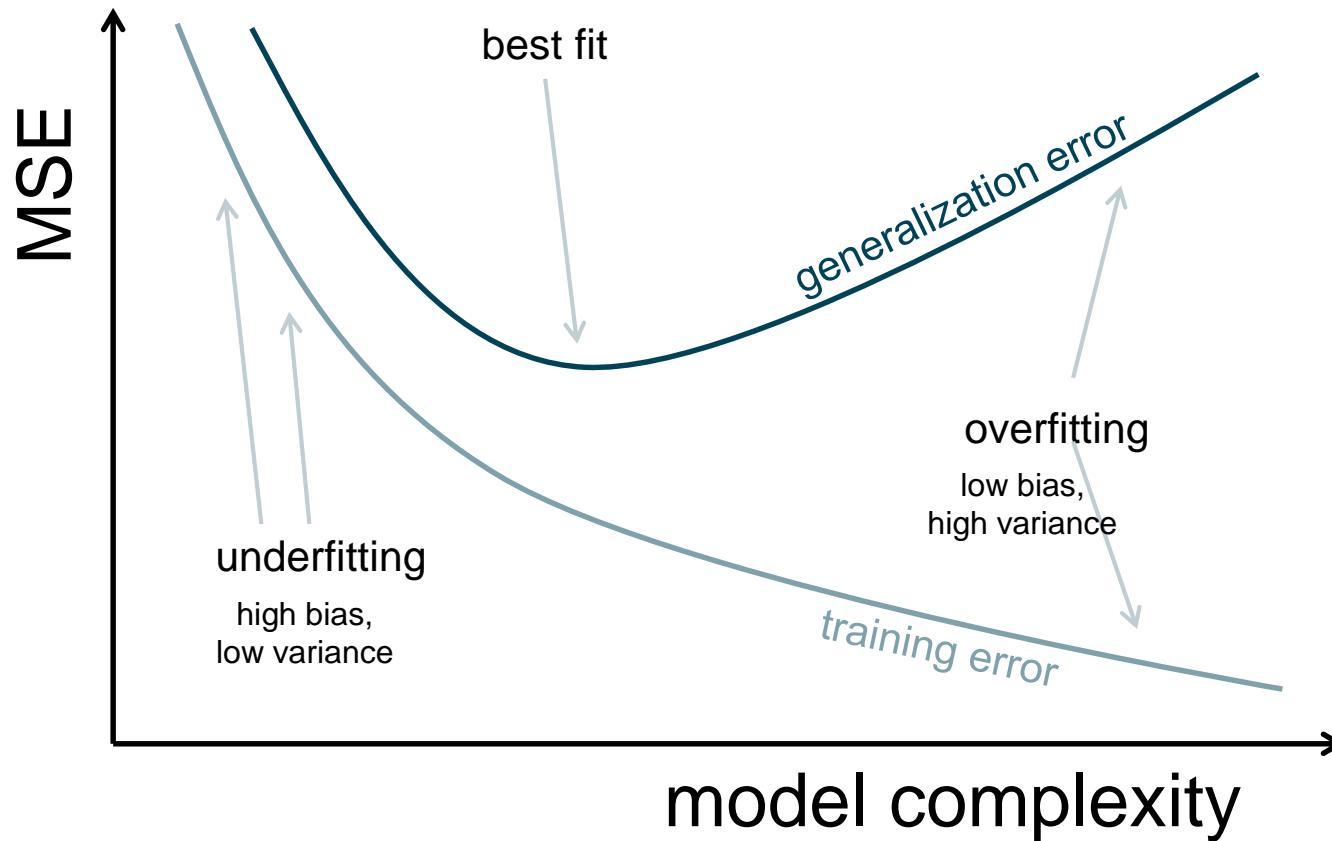
$$\begin{aligned} \text{MSE of } \hat{f} \text{ at } X_*: \quad \text{MSE}(\hat{f}|X_*) &= \mathbb{E}_{Y|X_*} \left[ (\hat{f}(X_*) - Y_*)^2|X_* \right] \\ &= \mathbb{E}_{Y|X_*} \left[ L(\hat{f}(X_*), Y_*)|X_* \right] \text{ for } L: (\hat{y}, y) \mapsto (\hat{y} - y)^2 \end{aligned}$$

$$\textbf{Proposition: } \text{MSE}(\hat{f}|X_*) = \text{Var}(\epsilon_*) + \text{Bias}(\hat{f}|X_*)^2 + \text{Var}(\hat{f}|X_*)$$

expected out-of-sample- MSE    „irreducible error“    Bias and variance of prediction  
from measurement noise

(proof in analogy to earlier bias-variance)

# The Bias-variance-trade-off in prediction



# Bias-variance-trade-off

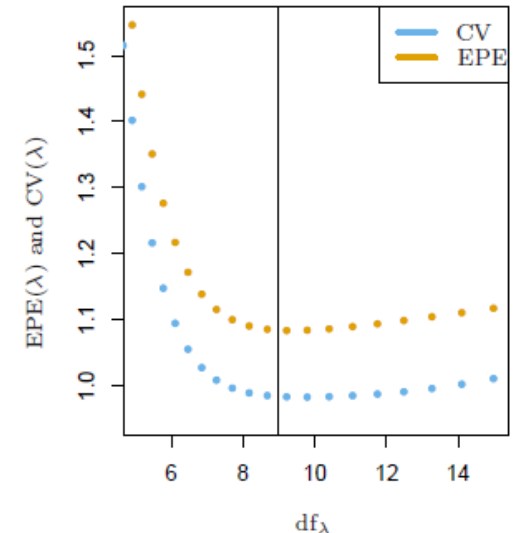
An experiment from  
*The Elements of Statistical Learning* (Section 5.5)

Interpolation using regularized splines

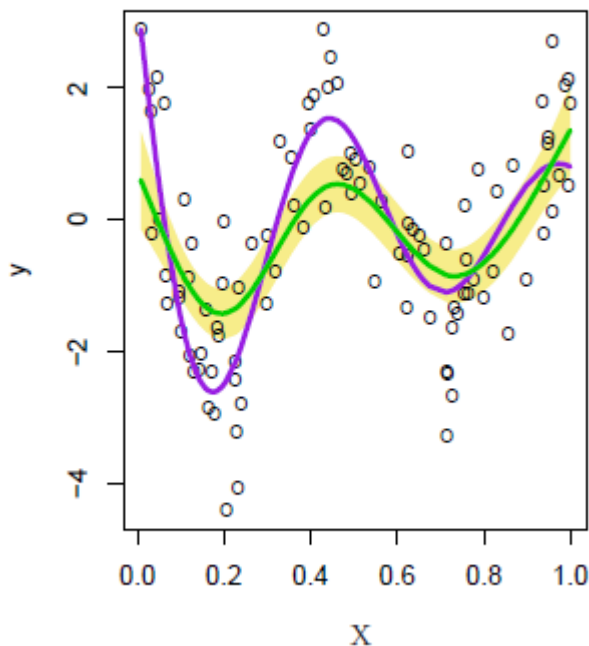
Strength of regularization set by parameter  $df_\lambda$

Prediction variance grows with decreasing  $df_\lambda$

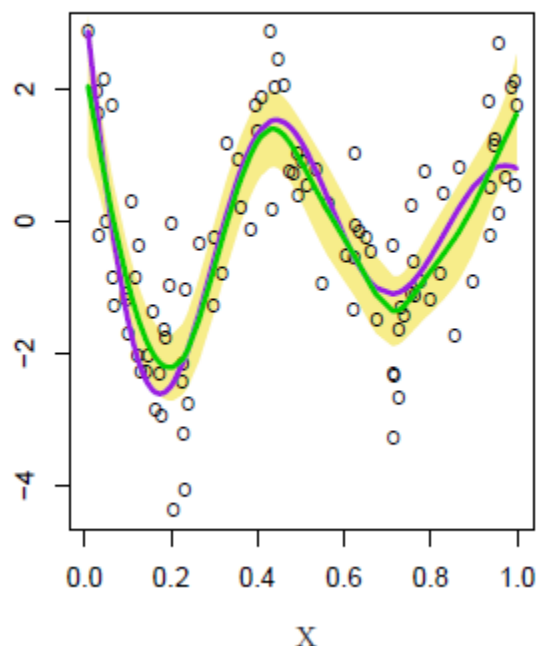
Prediction bias decreases with increasing  $df_\lambda$



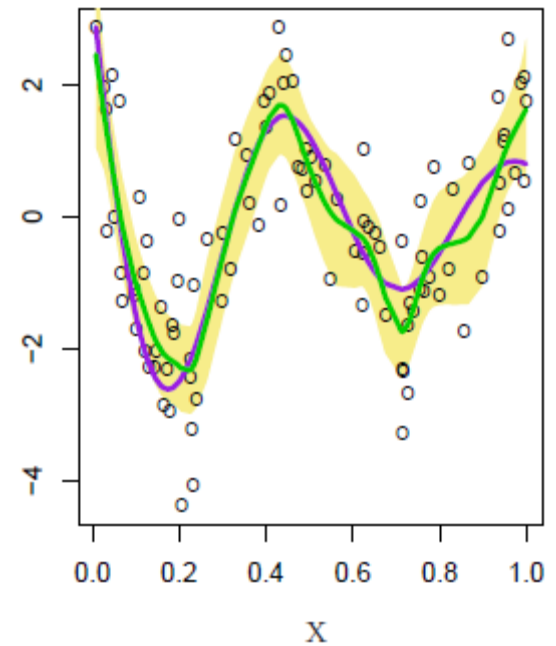
$df_\lambda = 5$



$df_\lambda = 9$



$df_\lambda = 15$



# Over-optimism of the training error

$$\text{bias of } \hat{f} \text{ at } X_*: \quad \text{Bias}(\hat{f}|X_*) = \mathbb{E}_{Y|X_*}[\hat{f}(X_*) - f(X_*)|X_*]$$

$$\text{variance of } \hat{f} \text{ at } X_*: \quad \text{Var}(\hat{f}|X_*) = \text{Var}_{Y|X_*} [\hat{f}(X_*)|X_*]$$

$$\text{MSE of } \hat{f} \text{ at } X_*: \quad \text{MSE}(\hat{f}|X_*) = \mathbb{E}_{Y|X_*} [(\hat{f}(X_*) - Y_*)^2|X_*]$$

**Proposition:**  $\text{MSE}(\hat{f}|X_*) = \text{Var}(\epsilon_*) + \text{Bias}(\hat{f}|X_*)^2 + \text{Var}(\hat{f}|X_*)$

is true only if  $\epsilon_*$  is independent of  $\hat{f}$

i.e., if  $(X_*, Y_*)$  has not already been seen in the training set

(this was assumed on the earlier slide)

If it has been *seen*, i.e., we test on the training set, then:

$$\begin{aligned} \text{MSE}(\hat{f}|X_*) &= \text{Var}(\epsilon_*) + \text{Bias}(\hat{f}|X_*)^2 + \text{Var}(\hat{f}|X_*) \\ &\quad - 2 \cdot \text{Cov}(\hat{f}(X_*)|X_*, Y_*|X_*) \end{aligned}$$

(expected training MSE)

Sensible learning machines' prediction usually co-varies with observations

hence **expected training MSE < expected test MSE**

*similar statements hold in more general settings in terms of noise and error statistic*

# Explicit form for Ordinary Least Squares

bias of  $\hat{f}$  at  $X_*$ :  $\text{Bias}(\hat{f}|X_*) = \mathbb{E}_{Y|X_*}[\hat{f}(X_*) - f(X_*)|X_*]$

variance of  $\hat{f}$  at  $X_*$ :  $\text{Var}(\hat{f}|X_*) = \text{Var}_{Y|X_*} \left[ \hat{f}(X_*)|X_* \right]$

MSE of  $\hat{f}$  at  $X_*$ :  $\text{MSE}(\hat{f}|X_*) = \mathbb{E}_{Y|X_*} \left[ (\hat{f}(X_*) - Y_*)^2|X_* \right]$

**Linear Regression:**  $f(x) = \beta^\top x + \alpha$ ,  $Y_i = f(X_i) + \varepsilon_i$ ,  $\hat{f}(x) = \hat{\beta}^\top x + \hat{\alpha}$

where objects with hats may (but don't need to) depend on training data

**OLS:**  $\mathbb{E}[\hat{\beta}] = \beta$ ,  $\mathbb{E}[\hat{\alpha}] = \alpha$  hence  $\text{Bias}(\hat{f}|X_*) = \langle \mathbb{E}[\hat{\beta}] - \beta, X_* \rangle + \mathbb{E}[\hat{\alpha}] - \alpha = 0$

for variance, use decomposition  $\hat{f}(x) = \bar{Y} + \hat{\beta}^\top (x - \bar{X})$

so  $\text{Var}(\hat{f}|X_*) = \text{Var}(\bar{Y}|X_*) + \text{Var}(\hat{\beta}^\top (X_* - \bar{X})|X_*)$  „error in estimating regression line“

$$= \sigma^2 \left( \frac{1}{N} + (X_* - \bar{X})^\top C_{xx}^{-1} \right) (X_* - \bar{X}) = \frac{\sigma^2}{N} \left( 1 + (X_* - \bar{X})^\top \Sigma^{-1} (X_* - \bar{X}) \right)$$

Full MSE decomposition:  $\text{MSE}(\hat{f}|X_*) = \sigma^2 + 0 + \frac{\sigma^2}{N} \left( 1 + (X_* - \bar{X})^\top \Sigma^{-1} (X_* - \bar{X}) \right)$

Overfitting optimism for already seen test point  $X_*$ : (= coefficient of  $\varepsilon_*$  in  $\hat{f}(X_*)$ )

$$\text{Cov}(\hat{f}(X_*)|X_*, Y_*|X_*) = \text{Cov}(\hat{f}(X_*)|X_*, \varepsilon_*|X_*) = \frac{\sigma^2}{N} \left( 1 + (X_* - \bar{X}) \cdot \Sigma^{-1} \cdot (e_i - \frac{1}{N}) \right)$$

# Unconditional bias-variance trade-off

test point  $X_*$ , test label  $Y_*$  where  $(X_*, Y_*) \sim (X, Y)$

*unconditional quantities: taking expectations*

$$\text{total bias of } \hat{f}: \quad \text{Bias}^2(\hat{f}) := \mathbb{E} \left[ \text{Bias}^2(\hat{f}|X_*) \right] \neq \mathbb{E} \left[ \text{Bias}(\hat{f})|X_* \right]^2$$

$$\text{total variance of } \hat{f}: \quad \text{Var}(\hat{f}|X_*) := \mathbb{E} \left[ (\hat{f}(X_*) - \mathbb{E}[\hat{f}(X_*)])^2 \right]$$

$$\text{total MSE of } \hat{f}: \quad \text{MSE}(\hat{f}) := \mathbb{E} \left[ (\hat{f}(X_*) - Y_*)^2 \right] = \epsilon(\hat{f})$$

**Proposition:**  $\text{MSE}(\hat{f}) = \text{Var}(\epsilon_*) + \text{Bias}^2(\hat{f}) + \text{Var}(\hat{f})$

**Proof:** take expectations in conditional trade-off

*Note: the „obvious“ unconditional generalization*

$$\text{MSE}(\hat{f}) = \text{Var}(\epsilon_*) + \text{Bias}(\hat{f})^2 + \text{Var}(\hat{f})$$

where  $\text{Bias}(\hat{f}) := \mathbb{E} \left[ \text{Bias}(\hat{f}|X_*) \right]$  is **wrong** in general!

# Prediction functionals vs prediction strategies

test point  $X_*$ , test label  $Y_*$  where  $(X_*, Y_*) \sim (X, Y)$

*unconditional quantities: taking expectations*



# Re-sampling strategies

# Re-sampling

Given data vector  $\mathcal{D} = (X_1, \dots, X_N)$   $X_1, \dots, X_N \sim X$  t.v.in  $\mathcal{X}$

estimator  $\hat{\theta} : \mathcal{X}^* \rightarrow \mathbb{R}$   $\mathcal{X}^* =$  vectors in  $\mathcal{X}$  of arbitrary length

estimate  $\hat{\theta}(\mathcal{D})$

„Re-sampling estimator“ is constructed from re-samples

$\hat{\theta}(\mathcal{D}[\pi_1]), \dots, \hat{\theta}(\mathcal{D}[\pi_k])$  (random!)

where  $\pi_i$  t.v.in  $\{1, \dots, N\}^*$  are re-sampling index vectors

Important cases:

$\pi_i$  are i.i.d. random with/without replacement of fixed size  $m$

$\pi_i$  invariant under permutation e.g., missing block of size  $k$  non-overlapping

Important applications:

obtaining an improved version of  $\hat{\theta}$ , e.g., variance-reduced

obtaining non-parametric estimates of expectation and variance

# Estimation of the generalization error

**Setting:** i.i.d. *test data*  $(X_1, Y_1), \dots, (X_M, Y_M) \underset{\text{i.i.d.}}{\sim} (X, Y)$  t.v.in  $\mathcal{X} \times \mathcal{Y}$

**prediction functional**  $f : \mathcal{X} \rightarrow \mathcal{Y}$  e.g.,  $\mathcal{X} = \mathbb{R}^n$  and  $\mathcal{Y} = \mathbb{R}$

**loss function**  $L : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  e.g.,  $L : (\hat{y}, y) \mapsto (\hat{y} - y)^2$

**To estimate:**  $\varepsilon(f) = \mathbb{E} [L(\hat{f}(X), Y)]$  *expected generalization loss*

**Theorems suggest following estimators:**

$$\hat{\varepsilon}(f) := \frac{1}{M} \sum_{i=1}^M L(f(X_i), Y_i) \quad \text{Observation: } L_i := L(f(X_i), Y_i) \text{ are i.i.d.}$$

„empirical loss“ since pairs  $(X_i, Y_i)$  are i.i.d.

$$\hat{v}(f) := \frac{1}{M(M-1)} \sum_{i=1}^M (L_i - \hat{\varepsilon})^2 \quad \text{„standard error of the empirical loss“}$$

Confidence interval:  $\left[ \hat{\varepsilon}(f) + \Phi^{-1}(\alpha/2) \cdot \sqrt{\hat{v}(f)}, \hat{\varepsilon}(f) - \Phi^{-1}(\alpha/2) \cdot \sqrt{\hat{v}(f)} \right]$

The end ... ?

**Big problem:** this is only valid if  $f$  is constant, e.g., already trained/fitted!

Otherwise  $L_i$  are dependent through a random  $f$ . (*no guarantees for strategies!*)

*But statements & guarantees about the trained prediction functionals are correct!*

# Estimation of the generalization error

**Setting:** i.i.d. data  $(X_1, Y_1), \dots, (X_M, Y_M) \underset{\text{i.i.d.}}{\sim} (X, Y)$  t.v.in  $\mathcal{X} \times \mathcal{Y}$

**prediction strategy**  $f_{\mathcal{T}}$  t.v.in  $[\mathcal{X} \rightarrow \mathcal{Y}]$  “trained” on subset  $(X_i, Y_i), i \in \mathcal{T}$

**loss function**  $L : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  e.g.,  $L : (\hat{y}, y) \mapsto (\hat{y} - y)^2$

**To estimate:**  $\varepsilon(f) = \mathbb{E}[L(f_{\mathcal{T}}(X), Y)]$  *expected generalization loss*

**One-split estimation:** training indices  $\mathcal{T} \subseteq \{1, \dots, N\}$  with  $\#\mathcal{T} = M$   
test indices  $\mathcal{V} = \{1, \dots, N\} \setminus \mathcal{T}$

$$\hat{\varepsilon}(f|\mathcal{T}) := \frac{1}{\#\mathcal{V}} \sum_{j \in \mathcal{V}} L(f_{\mathcal{T}}(X_j), Y_j) \quad \begin{array}{l} \text{„out-of-sample estimated“} \\ \text{since } f_{\mathcal{T}} \text{ indep. of } \{(X_i, Y_i) : i \in \mathcal{T}\} \\ \text{„empirical loss“} \end{array}$$

**Corollary of CLT:**  $\mathbb{E}[\hat{\varepsilon}(f|\mathcal{T})] = \varepsilon(f)$  i.e.,  $\hat{\varepsilon}(f|\mathcal{T})$  unbiasedly estimates  $\varepsilon(f)$

$\hat{\varepsilon}(f|\mathcal{T}) \xrightarrow{P} \varepsilon(f|\mathcal{T}) = \mathbb{E}[L(f_{\mathcal{T}}(X), Y)|f]$  when  $N \rightarrow \infty$  (with  $\mathcal{T}$  fixed)  
i.e.,  $\hat{\varepsilon}(f|\mathcal{T})$  consistently estimates  $\varepsilon(f|\mathcal{T})$

**Problem:** one-split estimation depends on single training set and single run  
so variance of the one-split estimator may be high if algorithm is unstable

**Solution:** Re-sample averaging,  $\hat{\varepsilon}_{CV} = \frac{1}{K} \sum_{\kappa=1}^K \hat{\varepsilon}(f|\mathcal{T}_{\kappa})$  **Potential issue:**  
 $\hat{\varepsilon}(f|\mathcal{T}_{\kappa})$  are correlated

# Variance reduction by averaging

**Variance reduction lemma:** for any number of correlated random variables

$Z_1, \dots, Z_M$  with  $\text{Var}(Z_i) = \sigma^2$  and  $\text{Corr}(Z_i, Z_j) = \rho$  for  $i \neq j$

one has 
$$\text{Var} \left( \frac{1}{M} \sum_{i=1}^M Z_i \right) = \rho \cdot \sigma^2 + \frac{1-\rho}{M} \cdot \sigma^2 \quad \begin{matrix} \leq \sigma^2 & \text{if } \rho > 0 \\ \neq \sigma^2 & \text{(unless } Z_i \text{ are equal)} \end{matrix}$$

So if  $Z_i$  are correlated re-samples with the same bias

averaging reduces the variance while not changing the bias

*Hence by bias-variance trade-off, expected mean-squared-error is reduced*

**For cross-validation estimator**  $\hat{\varepsilon}_{CV} = \frac{1}{K} \sum_{\kappa=1}^K \hat{\varepsilon}(f|\mathcal{T}_{\kappa})$  *(any type of CV or re-sampling where folds are exchangeably sampled!)*

$$\text{Var}(\hat{\varepsilon}_{CV}) = \left( \rho + \frac{1-\rho}{K} \right) \cdot \text{Var}(\hat{\varepsilon}[\mathcal{T}]) \quad \text{where} \quad \rho = \text{Corr}(\hat{\varepsilon}(f|\mathcal{T}_i), \hat{\varepsilon}(f|\mathcal{T}_j))$$

„variance of CV“      „1 or smaller“      „variance of one-split“      „correlation of fold-wise estimates“      for  $i \neq j$

*Since one-split estimates are unbiased, CV-estimator also is.*

*Hence by bias-variance trade-off, resample averaging reduces expected error.*

# Estimating the CV performance estimates' variance

**Why?**  $\hat{\varepsilon}_{CV}$  is an estimate for the expected prediction error of algorithm  $f$

Variance is needed for *confidence intervals*; also useful in *comparing* strategies

**Naive approach:**  $\hat{v}[\mathcal{T}] := \frac{1}{\#\mathcal{V} \cdot (\#\mathcal{V} - 1)} \sum_{j \in \mathcal{V}} (L(f_{\mathcal{T}}(x_j), Y_j) - \hat{\varepsilon}(f|\mathcal{T}))^2$

is an (unbiased, consistent) estimator for  $v[\mathcal{T}] = \text{Var}[L(f(X|\mathcal{T}), Y)|\mathcal{T}]$

**Problem (?):** conditional on training set, so variance from that is not included  
may be fine if algorithm is „stable“ w.r.t. training set choice and repetitions

**Bad but frequently seen approach:** sample variance of  $\hat{\varepsilon}(f|\mathcal{T}_i)$ ,  $i = 1 \dots K$   
*LLN/CLT does not apply:  $K$  is usually too small, and sample is correlated!*

**Averaging approach:**  $\hat{v}_{CV} := \frac{1}{K} \sum_{\kappa=1}^K \hat{v}[\mathcal{T}_i]$  reduces variance from training set  
(by variance reduction lemma)

**Problem:** this is not the variance of  $\hat{\varepsilon}_{CV}$  which it usually overestimates!  
(but being conservative is fine for avoiding type I errors)

**Theorem:** There is no unbiased estimator for the variance of  $\hat{\varepsilon}_{CV}$  (of certain form)  
(Bengio, Grandvalet 2004) Key realization: non-identifiability of inter-fold co-variance

*There could be biased estimators with low MSE... but no substantial (?) progress since.*

# Black-box estimates of mean and variance

Given data vector  $\mathcal{D} = (X_1, \dots, X_N)$   $X_1, \dots, X_N \underset{\text{i.i.d.}}{\sim} X$  t.v.in  $\mathcal{X}$

estimator  $\hat{\theta} : \mathcal{X}^* \rightarrow \mathbb{R}$  estimate  $\hat{\theta}(\mathcal{D})$

**Jackknife estimation of**  $\mathbb{E} [\hat{\theta}(\mathcal{D})]$  and  $\text{Var} [\hat{\theta}(\mathcal{D})]$

$\hat{\theta}_{/i} := N \cdot \hat{\theta}(\mathcal{D}) - (N-1) \cdot \hat{\theta}(\mathcal{D}_{/i})$  where  $\mathcal{D}_{/i} = (X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_N)$   
 „ordinary jackknife pseudo-samples“ „dataset minus i-th data point“

**Bootstrap estimation of**  $\text{Var} [\hat{\theta}(\mathcal{D})]$

$B_i := \mathcal{D}[\pi_i]$  where  $\pi_1, \dots, \pi_B$  are i.i.d. t.v.in  $\mathbb{N}^N$

$\pi_i = (n_1, \dots, n_N), n_i \underset{\text{i.i.d.}}{\sim} \text{Unif}\{1, \dots, N\}$

Sample mean and variance are „good“ estimates of  $\mathbb{E} [\hat{\theta}(\mathcal{D})]$  and  $\text{Var} [\hat{\theta}(\mathcal{D})]$

under certain regularity assumptions on  $\hat{\theta}$  (these are often unclear in literature!)

**Applicability:** Re-sample estimates of full training/test variance

Confidence intervals for loss statistics which are not means

**Cave:** Regularity conditions may be strong! Inappropriate for medians/quantiles

# Testing & Comparison



# Principles of model comparison testing

First interesting and simpler case: prediction *functionals*

**Setting: i.i.d. test data**  $(X_1, Y_1), \dots, (X_M, Y_M) \stackrel{\text{i.i.d.}}{\sim} (X, Y)$  t.v.in  $\mathcal{X} \times \mathcal{Y}$

**prediction functionals**  $f_1, \dots, f_S : \mathcal{X} \rightarrow \mathcal{Y}$  e.g.,  $\mathcal{X} = \mathbb{R}^n$  and  $\mathcal{Y} = \mathbb{R}$

**loss function**  $L : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  e.g.,  $L : (\hat{y}, y) \mapsto (\hat{y} - y)^2$

**Testing approach:** obtain loss residuals

$$L_{ij} := L(f_j(X_i), Y_i), \quad i = 1 \dots M, \quad j = 1 \dots S$$

the  $S$  samples  $\mathcal{S}_j := (L_{1j}, \dots, L_{Mj})$  are tupled/paired

any paired and tupled portmanteau test for location comparison is applicable

Recommendation: non-parametric Wilcoxon (rank-sign) test

„comparing confidence intervals“ is also fine (unpaired z-test or t-test, conservative)

**Prediction strategies?** *Unclear what to do.*

Nadeau, Bengio (2004) discuss a few options based on the t-test

**Re-sampling compatibility?** *2x median p-value is conservative aggregation.*

Generally for exchangeable re-samples:  $\gamma^{-1}p_{(\gamma k)}$  is conservative aggregation of  $(p_1, \dots, p_k)$

# WORQ - widely open research questions!

(aka MSc/PhD topics for the ambitious and theoretically inclined)

## Theory of variance estimators for black-box functionals

Bootstrap & re-sampling seem to be the only semi-solid strategies

Also, just semi-solid, without source that has exact assumptions (?)

How best to separate test set variance and training set variance?

## Theory of variance estimators for re-sampled statistics

Bengio/Grandvalet: no unbiased estimator *of a special form*

This does not preclude a good estimate of another, simple form

... such as re-sampling the re-sample statistic (??)

## Do all this for the „complicated tasks“

Time series, on-line learning, anomaly detection, reinforcement learning

Structured and heterogeneous prediction tasks

Probabilistic and Bayesian modelling (credibility intervals?)

## Hypothesis testing & portmanteau comparison

Best way for strategies unclear – how to incorporate training set variance?

# **Misunderstandings and Statistical Learning Theory**

## **The No Free Lunch Theorems** (Wolpert and Macready, 1997 onwards)

*... for all data there is a model, for all models there is data...*

### **Frequent mis-interpretation:**

All of statistics and machine learning is arbitrary anyway.

### **More correct interpretation:**

No meaningful definition of „learning“ is possible without assumptions on how training and test data relate. For example, they should be similar (e.g. distributionally).

## **Shao – Linear Model Selection by Cross-Validation (1993)**

*... some types of CV fail to identify the „correct“ model...*

### **Frequent mis-interpretation:**

Cross-validation should not be done. Or: Bayesian statistics is the only way.

### **More correct interpretation:**

Model identification is more difficult than accurate prediction, sometimes considerably so.

## **Bengio, Grandvalet – No Unbiased Estimator for the Variance of K-Fold Cross-Validation (2004)**

### **Frequent mis-interpretation:**

Confidence intervals for error metrics cannot be computed, quantitative comparisons between different methods are futile.

### **More correct interpretation:**

Predictions on different folds are correlated, one needs to be careful in aggregating them.

***Recall: all models are wrong (George Box), but some are useful.***

***This is similarly true for meta-methods and model checking.***

# Overview: Model-specific Learning Theory

multiple „flavours“ of model-specific guarantees based on „model class complexity“  
 bound generalization loss  $R(f) = \varepsilon(f)$  by empirical *training* loss  $R_N(f)$  plus “complexity term”

<i>approach/field</i>	<i>Scope and assumptions</i>	<i>some notable statements</i>
<b>Statistical Learning Theory</b> Vapnik, Chervonenkis	Training and test data follow the <i>same</i> distribution Learning machine is in a fixed set of functions $\mathcal{F}$ Questions: asymptotic behaviour of machine Relation between training and test error	$R(f) \leq R_N(f) + \frac{3}{\sqrt{N}} \cdot \sqrt{\log S_{\mathcal{F}}(2n) + \log \frac{2}{\varepsilon}}$ for any learner $f$ with probability $\geq 1 - \varepsilon$ $R(f), R_N(f)$ expected and empirical loss/risk $S_{\mathcal{F}}(n)$ number of classification rules on $n$ points general extensions via Rademacher/covering theory
<b>Bayesian/ Parametric</b>	Training and test data follow the <i>same</i> distribution Learning machine in parametric function class $\mathcal{F}$ Predictions are distributional „posterior“	$\text{AIC} = -2\mathcal{L}(f) + 2d \quad \text{BIC} = -2\mathcal{L}(f) + d \log N$ „Akaike information criterion“ „Bayesian IC“ for many simple model classes: Selection by AIC is asymptotically equivalent to LOOCV Selection by BIC is as. eq. to certain leave-out-CV
<b>PAC-Bayesian</b>	Training and test data follow the <i>same</i> distribution Learning machine in stochastic function class $\mathcal{F}$ Inspired by SLT and Bayesian paradigm	$R(f) \leq R_N(f) + \frac{3}{\sqrt{N}} \cdot \sqrt{\text{KL}(f \pi) + \log \frac{2}{\varepsilon}}$ for any learner $f$ with probability $\geq 1 - \varepsilon$ $\text{KL}(f \pi)$ Kullback-Leibler divergence $\pi$ prior belief $\pi$ in reference class $\mathcal{F}$
<b>Minimum Description Length</b>	Training and test data follow the <i>same</i> distribution Parametric function class $\mathcal{F}$ similar to Bayesian Based on information theoretical argumentation maximizing Bayes posterior is posited best	similar to PAC-Bayesian (quantities are interpreted information theoretically)