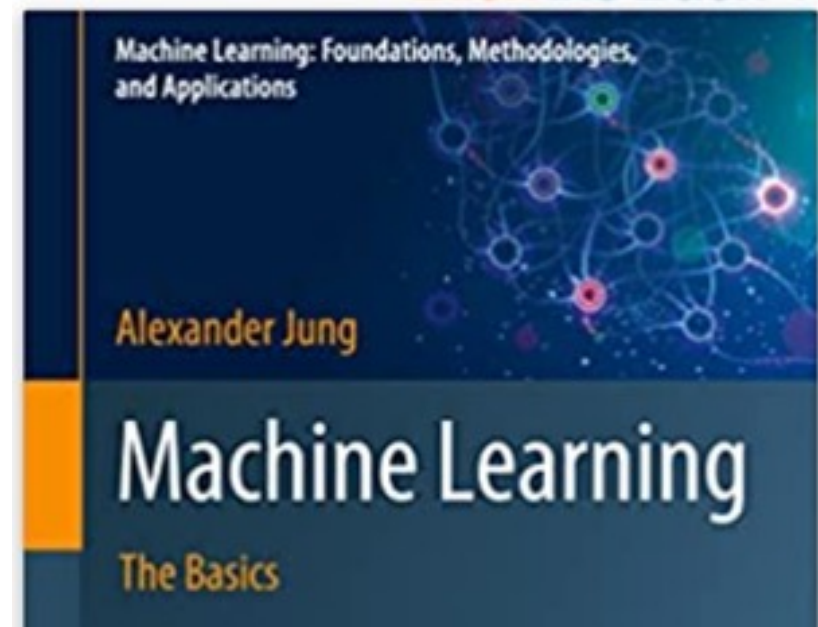# Model Validation and Selection

Alex(ander) Jung
Assistant Professor for Machine Learning
Department of Computer Science
Aalto University

# Reading.

Ch. 6 of https://mlbook.cs.aalto.fi

**Machine Learning: Foundations, Methodologies, and Applications**

Alexander Jung

**Machine Learning**

**The Basics**

scikit learn

Install  User Guide  API  Examples  Community  More ▾

Prev        Up        Next

scikit-learn 1.1.1
Other versions

Please cite us if you use the

## 3. Model selection and evaluation

### 3.1. Cross-validation: evaluating estimator performance

- 3.1.1. Computing cross-validated metrics

https://scikit-learn.org/stable/model_selection.html

# Model Validation

How do we know a ML method is any good ?

A. Jung HCML Summer School'22

# Model Selection

How to choose between different alternative methods?

# Learning Goals

- know train err is bad quality measure for ML method

- val.err. is more useful as quality measure for a ML model

- basic idea of k-fold CV

- hyper-parameter tuning = model selection

- Python implementations of k-fold CV / gridsearch

# "Model"
# =
# Hypothesis Space

# What are three main components of machine learning?

# 1. Data

A. Jung HCML Summer School'22

# Data

- set of "data points" (atomic unit of information)

- data point has <span style="color:red">features and labels</span>

- <span style="color:red">features</span> are properties that <span style="color:red">can measured easily</span>

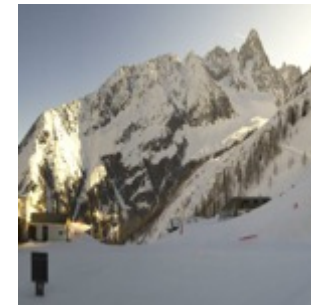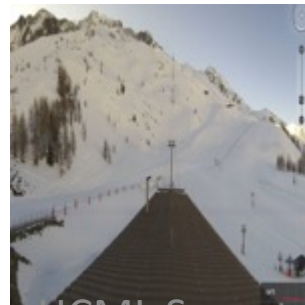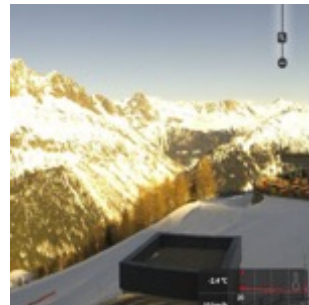- <span style="color:red">labels</span> =higher-level facts or quantities of interest
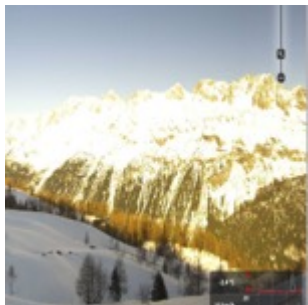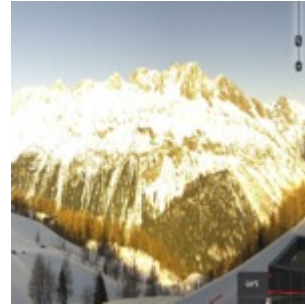
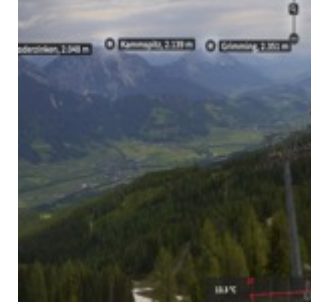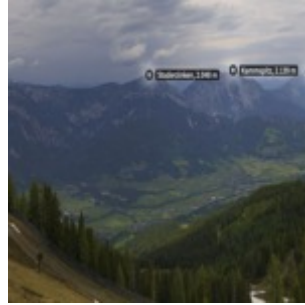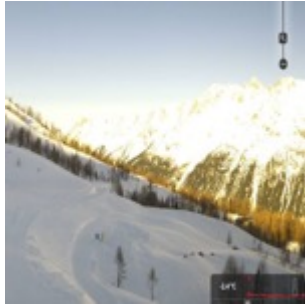# Data Point = "Some Ski Day"

feature x : morning temperature

label y : maximum daytime temperature

A. Jung HCML Summer School 22

# Data = Bunch of Data Points

# Sample Size

"sample size" m

=

number of (labeled) data points

# Sample Size m = 4



label y

feature x

# 2. Hypothesis Space

# How Many Hypotheses Are There?
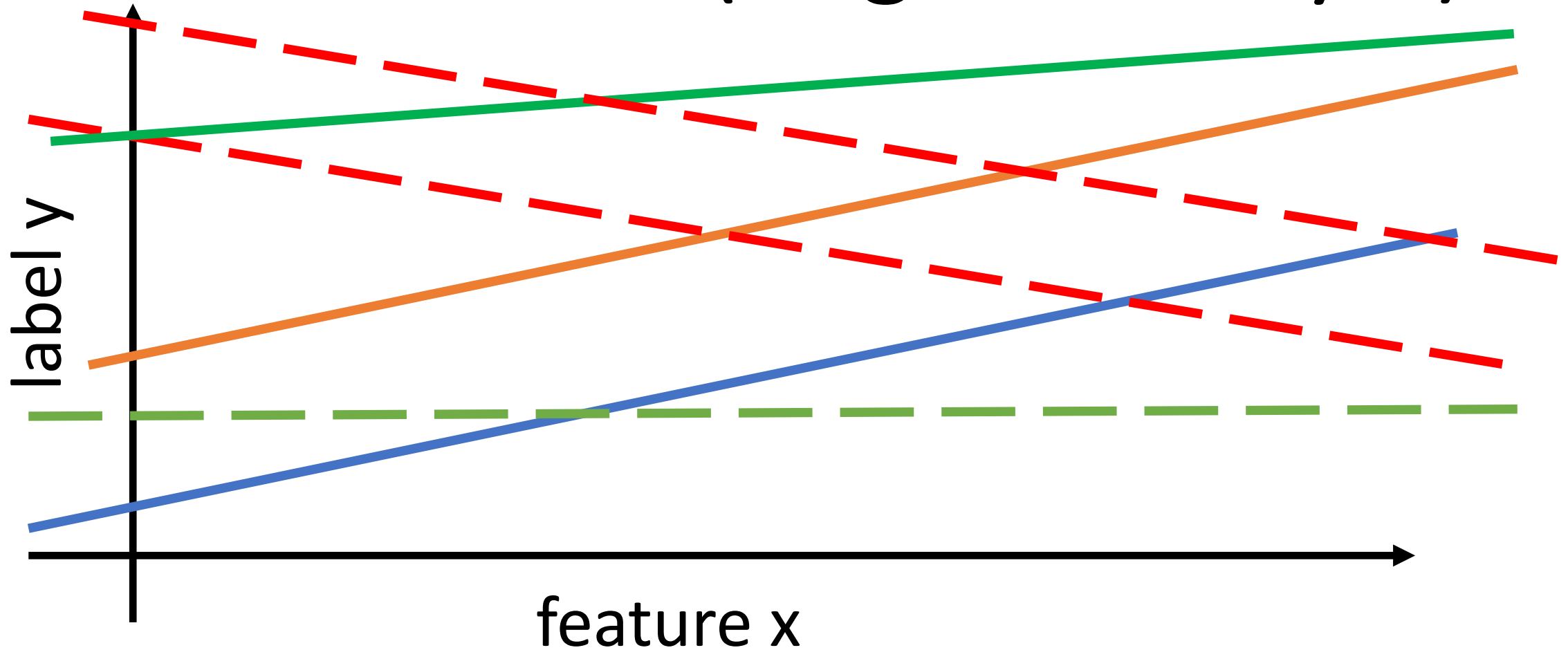
predict

h(x)

predicted max temp.

$$\hat{y} = -2.3$$

feature x = -10

# Model 1:
# Linear Predictors (Degree 1 Polyn.)



label y
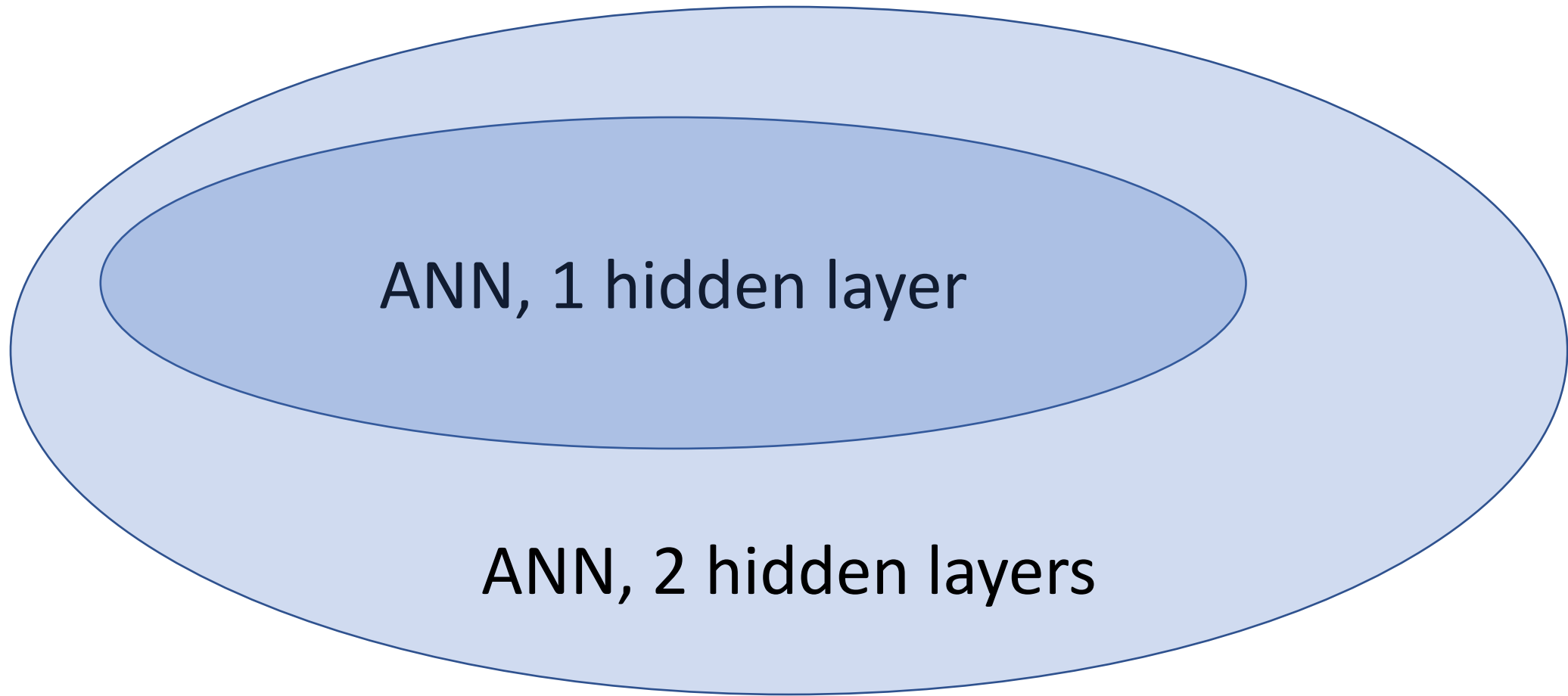
feature x

# Model 2:
# Degree 3 Polyn. Predictors



label y

feature x

# Nested Models – I



Model 1: linear predictors

Model 2: degree 3 polyn.

# Nested Models - II



ANN, 1 hidden layer

ANN, 2 hidden layers

# Nested Models - III



effective hyp. space @
1 GD step

2 GD steps

3 GD steps

# Math Notation

$$\mathcal{H}^{(n)} = \left\{ h(x) = \sum_{l=0}^{n-1} w_l x^l \; with \; some \; w_l \right\}$$

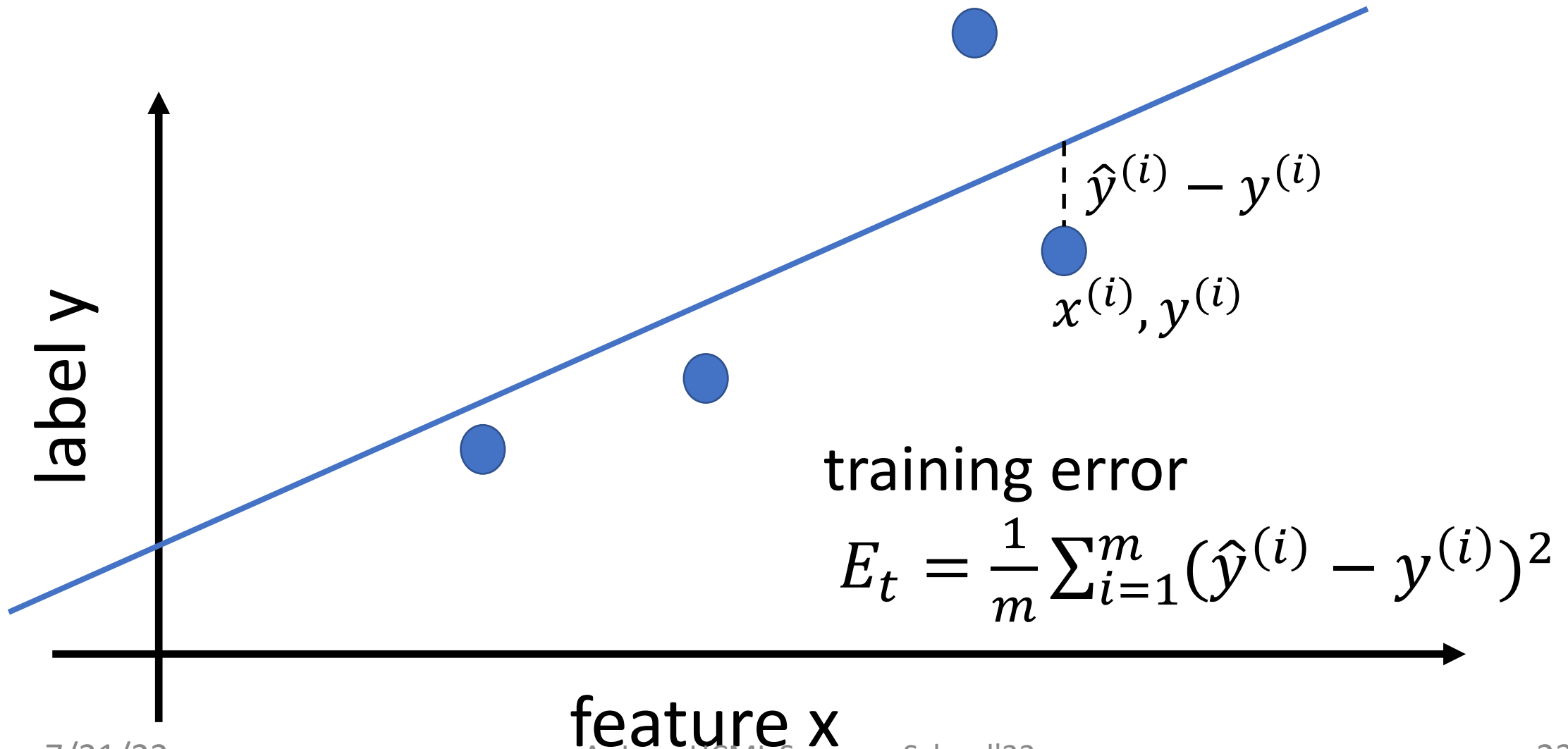$\mathcal{H}^{(2)}$ … linear hypotheses

$\mathcal{H}^{(4)}$ … degree 3 polyn.

$$\mathcal{H}^{(1)} \subseteq \mathcal{H}^{(2)} \subseteq \mathcal{H}^{(3)} \subseteq \mathcal{H}^{(4)} \subseteq \ldots$$

# 3. Loss Function

# Learn Linear Predictor



$\hat{y}^{(i)} - y^{(i)}$

$x^{(i)}, y^{(i)}$

training error

$$E_t = \frac{1}{m} \sum_{i=1}^{m} (\hat{y}^{(i)} - y^{(i)})^2$$

label y

feature x

# Learn Degree 3 Polyn.

$$\hat{y}^{(i)} - y^{(i)}$$

label y

$$x^{(i)}, y^{(i)}$$

training error

$$E_t = \frac{1}{m} \sum_{i=1}^{m} (\hat{y}^{(i)} - y^{(i)})^2$$

feature x

# Training Errors



model 1
linear predictors

model 2:
degree 3 polyn.

# Overfitting



label y

$\hat{y}^{(i)} - y^{(i)}$

$x^{(i)}, y^{(i)}$

★ "new" datapoint

training error

$E_t = \frac{1}{m}\sum_{i=1}^{m}(\hat{y}^{(i)} - y^{(i)})^2$

feature x

# Small Training Error Does Not Imply Good Performance on New Data Points!

# Small Training Error Merely Indicates That Optimization/Training Algorithm Works

A. Jung HCML Summer School'22

# A Case in Point

we can perfectly fit (almost) any m data points using polynomials of degree n-1 as soon as

$$n \geq m$$

m=2, n=2

# Reminder: Probabilistic Model

- data points are realizations of RVs

- joint pdf p(x,y) of features and label

- training set is a RV

- learnt hypothesis h(.) is a RV

- prediction h(x) is a RV

# Why is Train. Err. Misleading?

- consider expected loss of hypothesis
- estimate expectation using sample average
- this only works if hypothesis does not depends on data points used in average
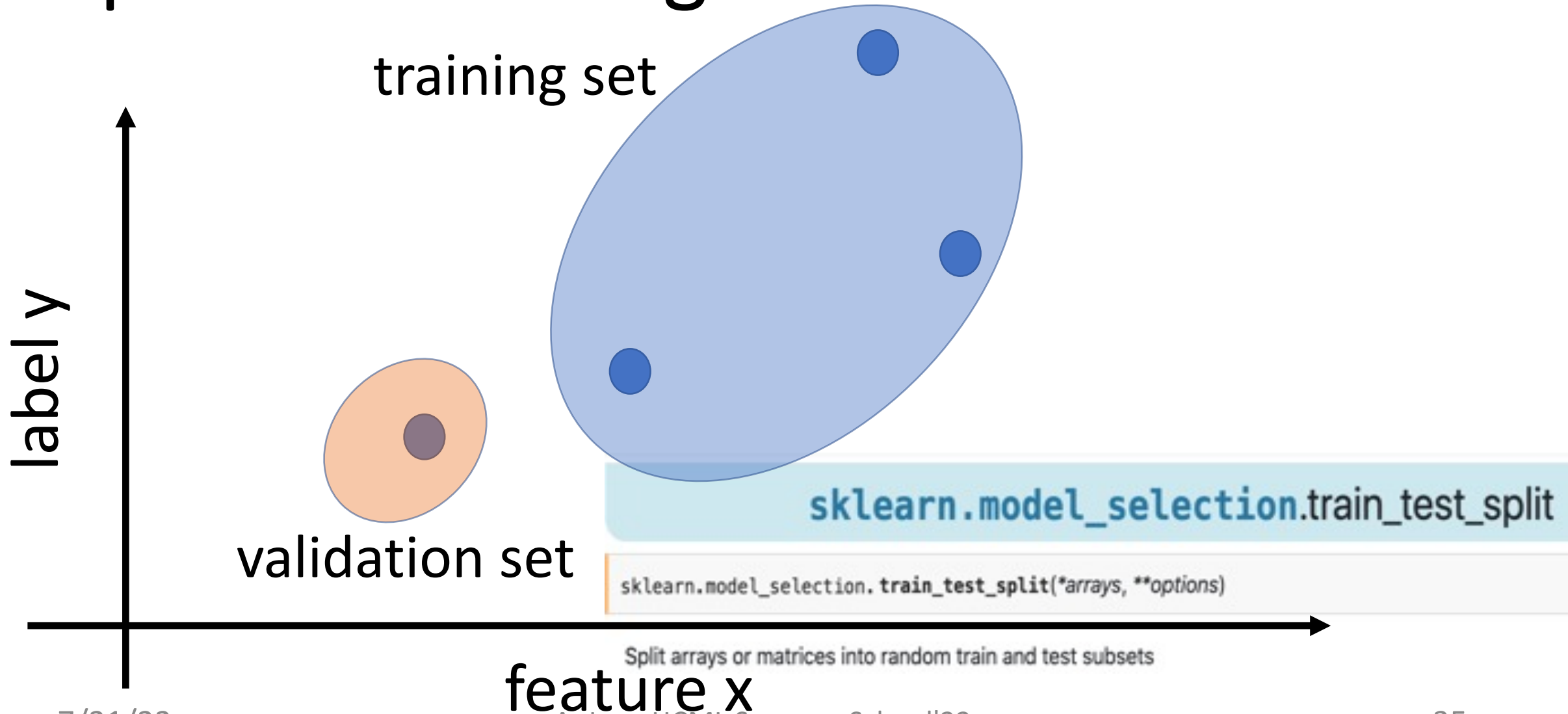- does not hold for training error
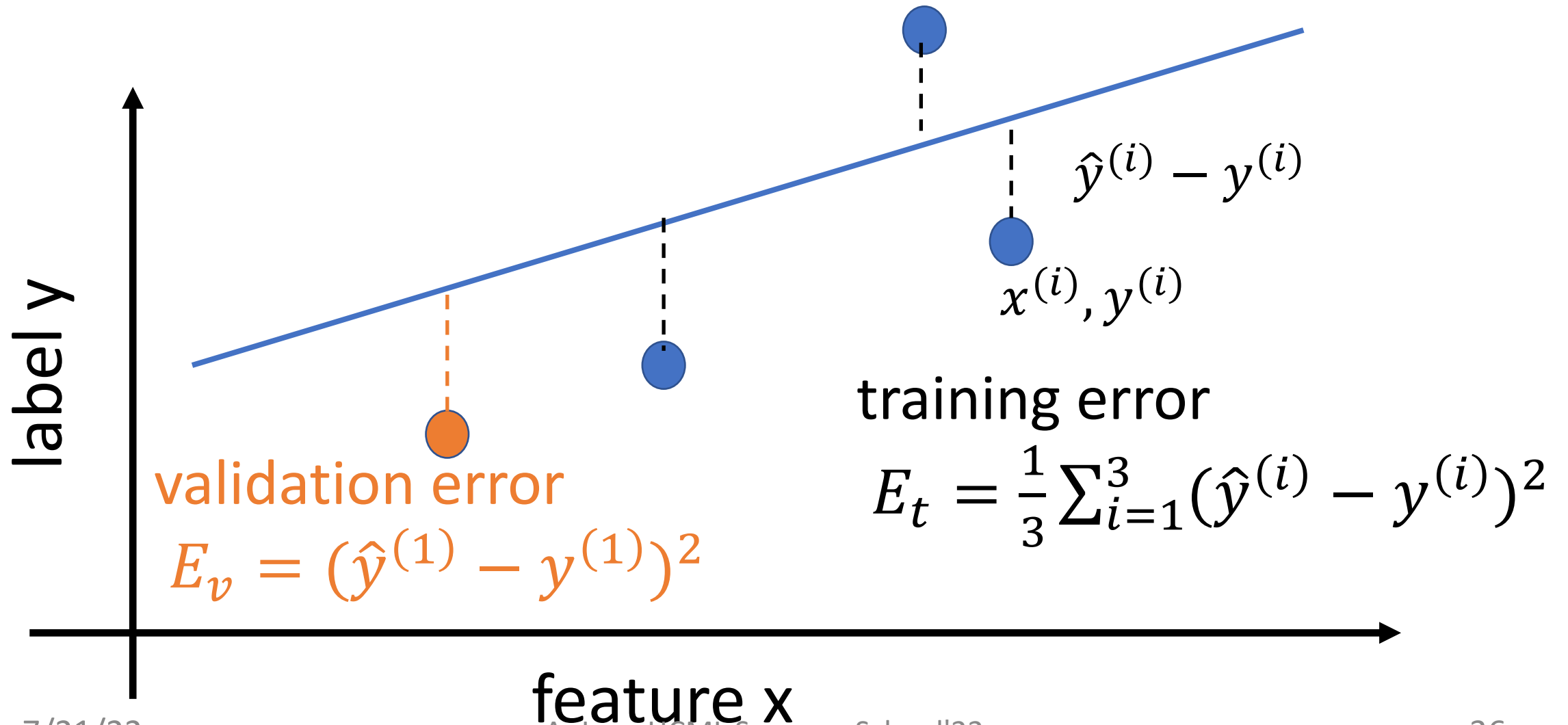
# Model Validation and Selection

A. Jung HCML Summer School'22

# Basic Idea of Validation

- divide data points into two subsets

- use <span style="color:red">training set</span> to learn predictor

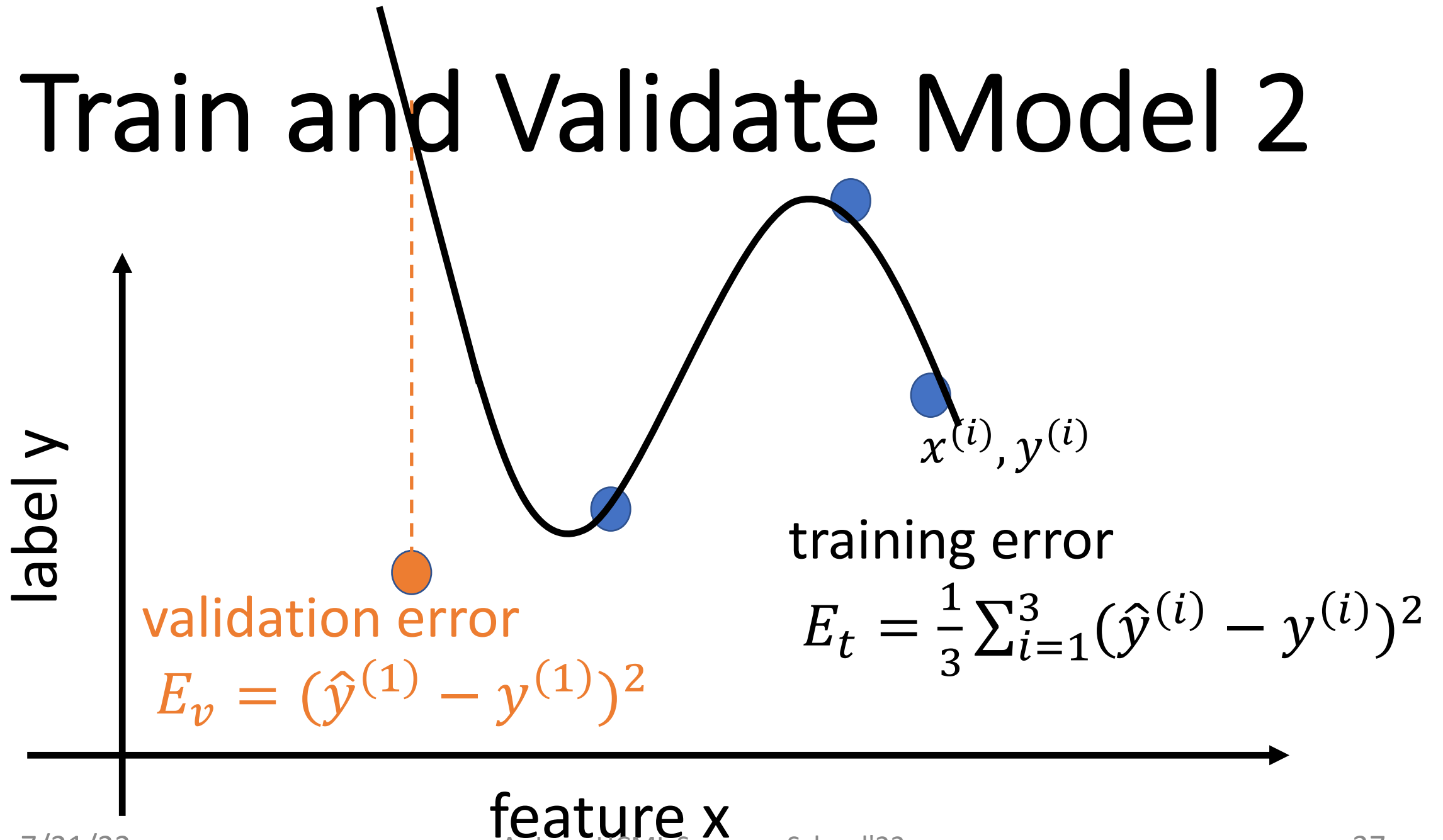- use <span style="color:red">validation set</span> to estimate loss

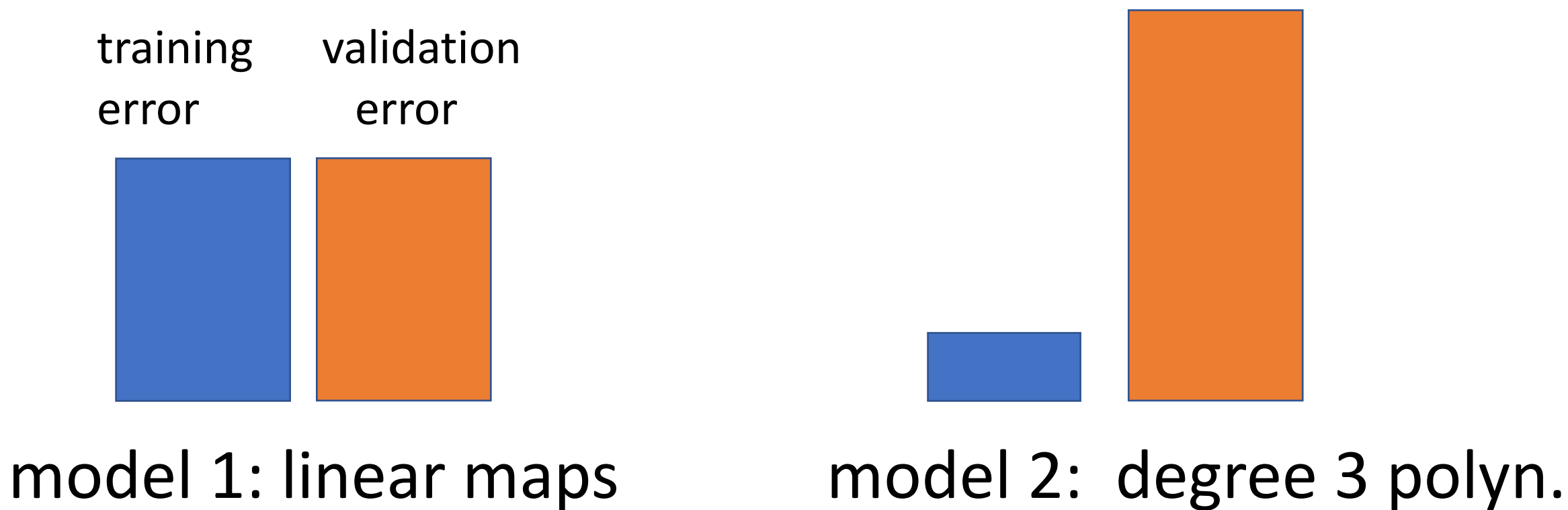# Split into Training and Validation Set



training set

label y

feature x

validation set

sklearn.model_selection.train_test_split

sklearn.model_selection. **train_test_split**(*arrays, **options*)

Split arrays or matrices into random train and test subsets

# Train and Validate Model 1



validation error
$$E_v = (\hat{y}^{(1)} - y^{(1)})^2$$

training error
$$E_t = \frac{1}{3}\sum_{i=1}^{3}(\hat{y}^{(i)} - y^{(i)})^2$$

$$\hat{y}^{(i)} - y^{(i)}$$

$$x^{(i)}, y^{(i)}$$

label y

feature x

# Train and Validate Model 2



label y

feature x

$x^{(i)}, y^{(i)}$

training error
$$E_t = \frac{1}{3}\sum_{i=1}^{3}(\hat{y}^{(i)} - y^{(i)})^2$$

validation error
$$E_v = (\hat{y}^{(1)} - y^{(1)})^2$$

# Basic Idea of Model Selection
## choose model via validation error

training
error

validation
error

model 1: linear maps

model 2: degree 3 polyn.

# Train/Val Error vs Model Complexity

$$\mathcal{H}^{(n)} = \left\{ h(x) = \sum_{l=0}^{n-1} w_l x^l \text{ with } \textcolor{red}{\text{weights } w_l} \right\}$$

model dimension/complexity n

# k-Fold Cross Validation

- might be unlucky with train/val split

- problematic for small datasets

- IDEA: randomly split several times

- "average out" unlucky splits

# K-Fold Cross Validation
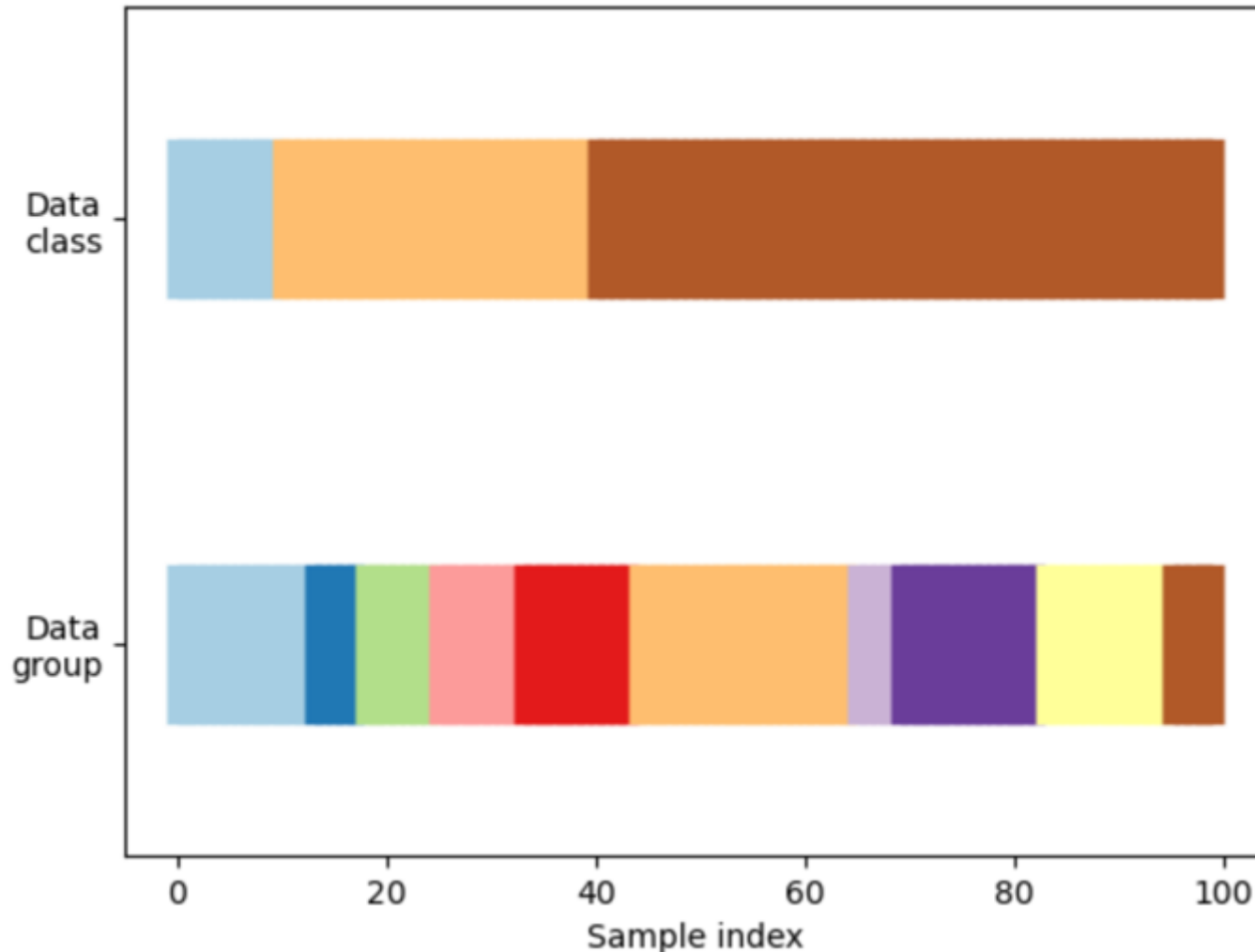


fold 1

fold 2

fold 3

# k-Fold Cross Validation

how to choose nr of folds (the "k" in k-fold CV) ?

- train fold should be sufficiently large (avoid overfitting)

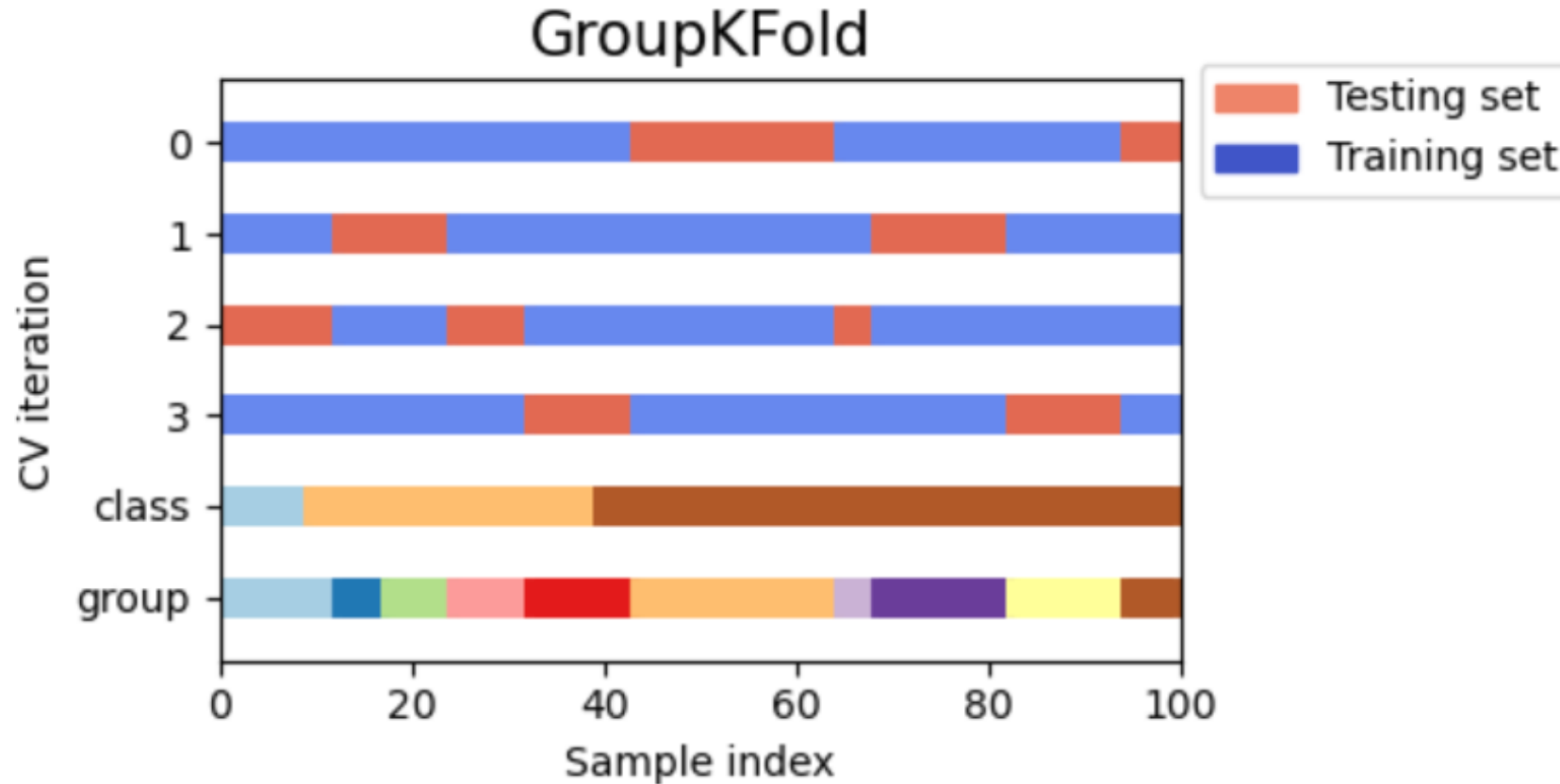- val  folds should sufficiently large (to get reliable estimate of generalization)

# CAUTION!

- k-fold CV requires a method to split into folds

- most basic method: evenly divide into k folds

- works if data is i.i.d. ("order of data points is arbitrary")

- fails if data points are grouped or ordered

# Imbalanced Classes and Group Structure



- e.g. data points with same label are contiguous blocks

- or data points are obtained at consecutive time instants ($\rightarrow$ correlations)
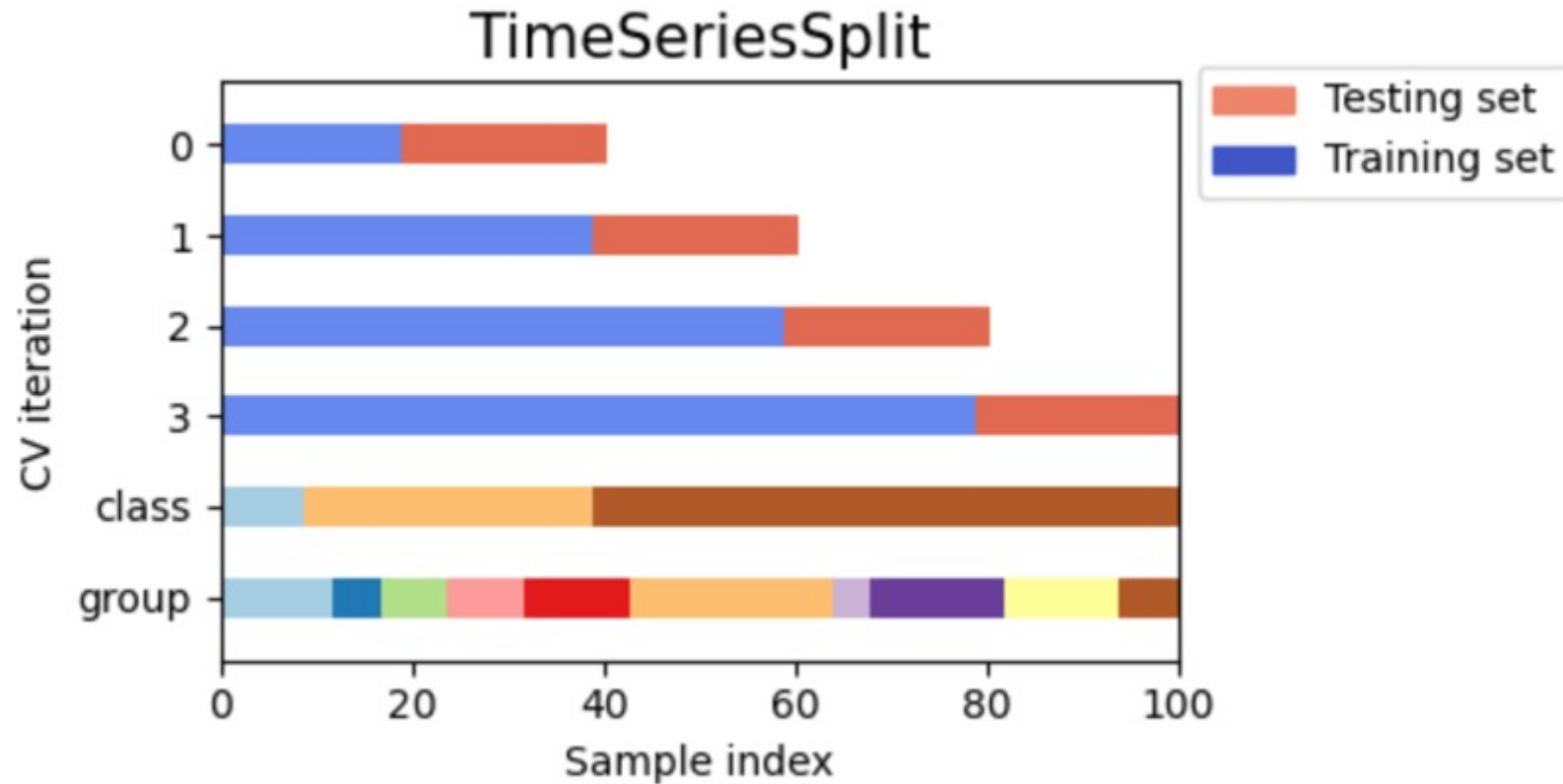
# Group-Preserving Splitting



https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GroupKFold.html
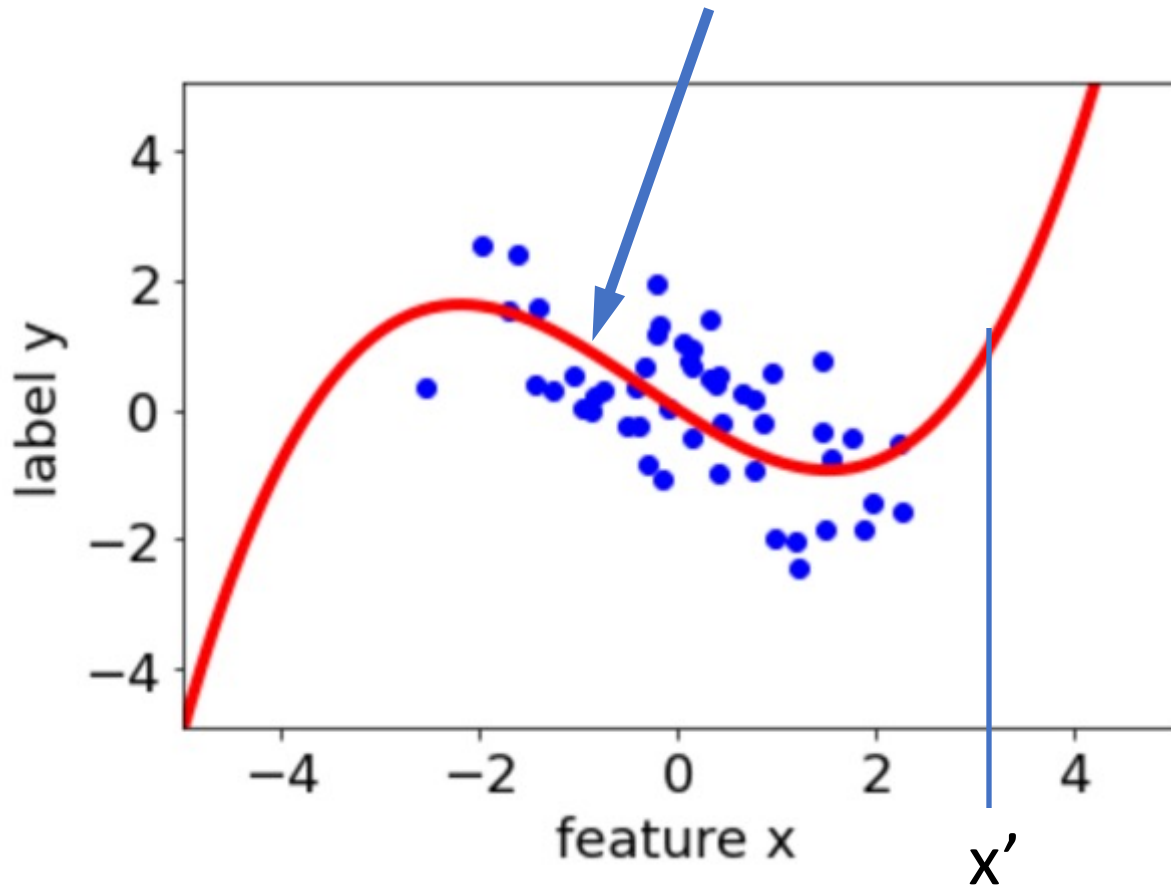
# Class-Ratio Preserving Splitting



StratifiedKFold

# Temporal Successive Splitting



source: https://scikit-learn.org/stable/

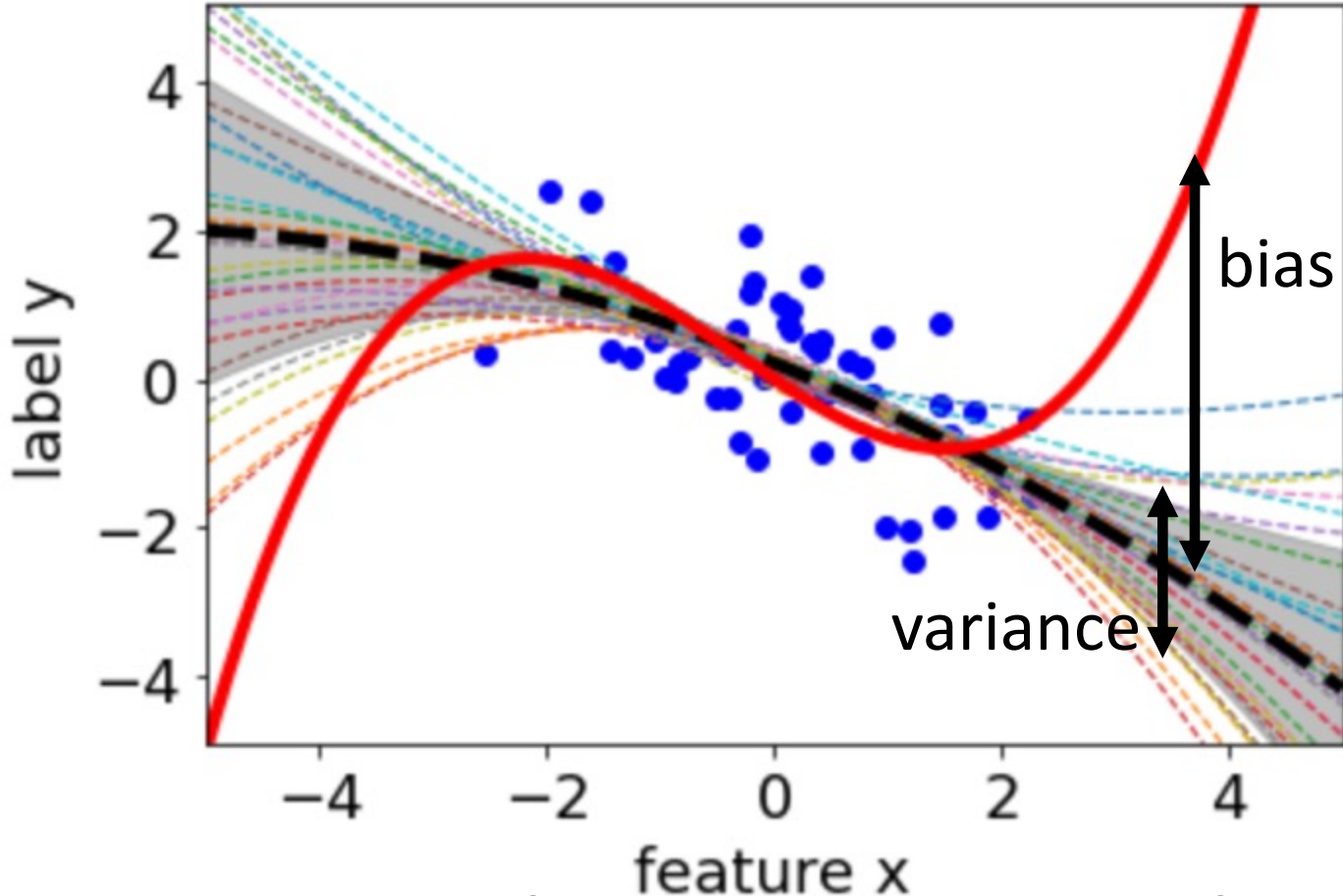# Bias and Variance Decomposition

# Toy Data

y = g(x) + "noise"



learn hypothesis h(.) using a randomly selected training set

compute prediction h(x') for a fixed feature value x'
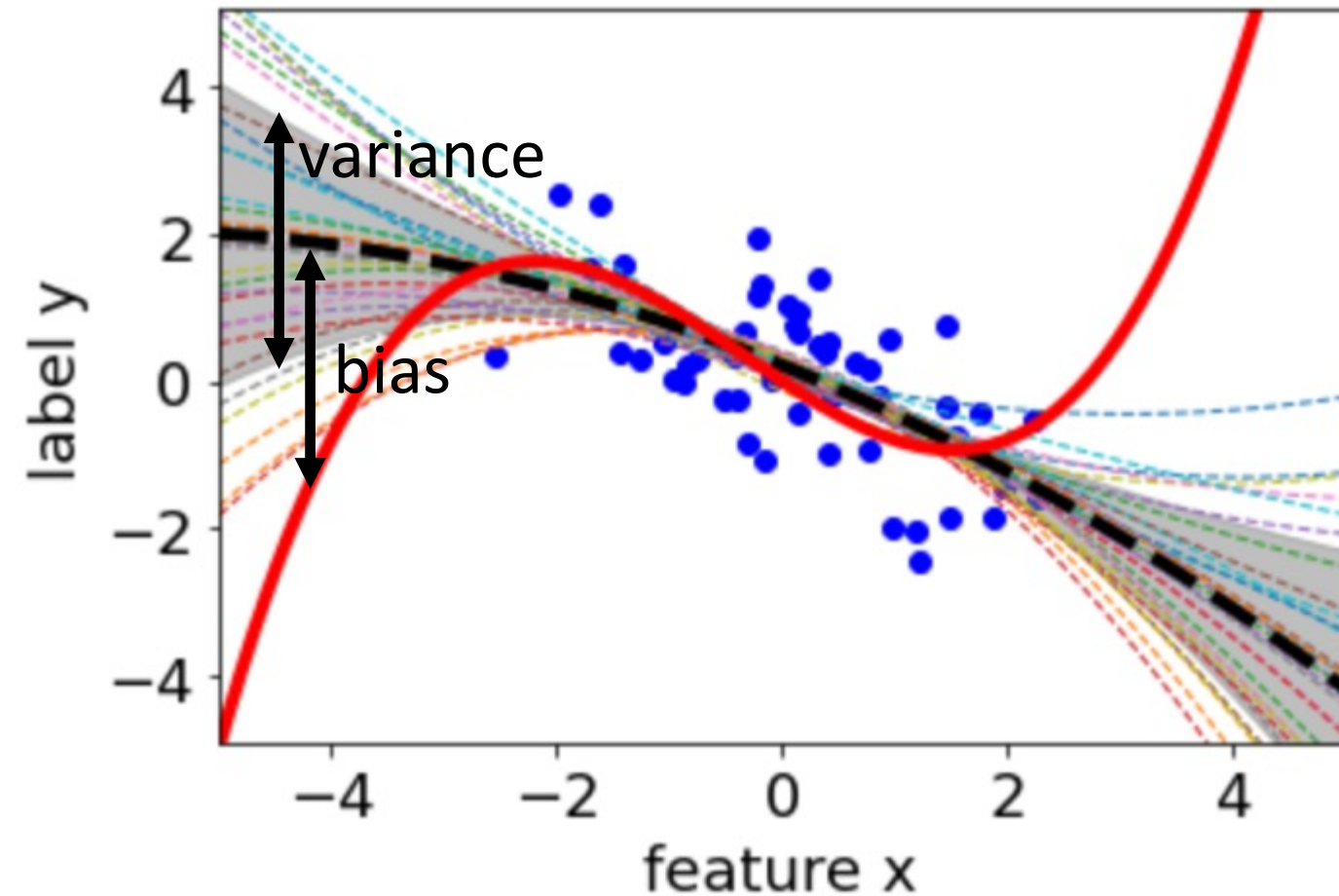
# Ensemble of Learnt Hypotheses



hypothesis learn on train set #1

hypothesis learn on train set #2

# Bias and Variance



$$\hat{y} = h(x')$$

RV since obtained
from a randomly
selected training set

$$\mathsf{E}\{(\hat{y}-y)^2\} = (\mathsf{E}\{\hat{y}\}-y)^2 + \mathsf{E}\{(\hat{y}-\mathsf{E}\{\hat{y}\})^2\}$$
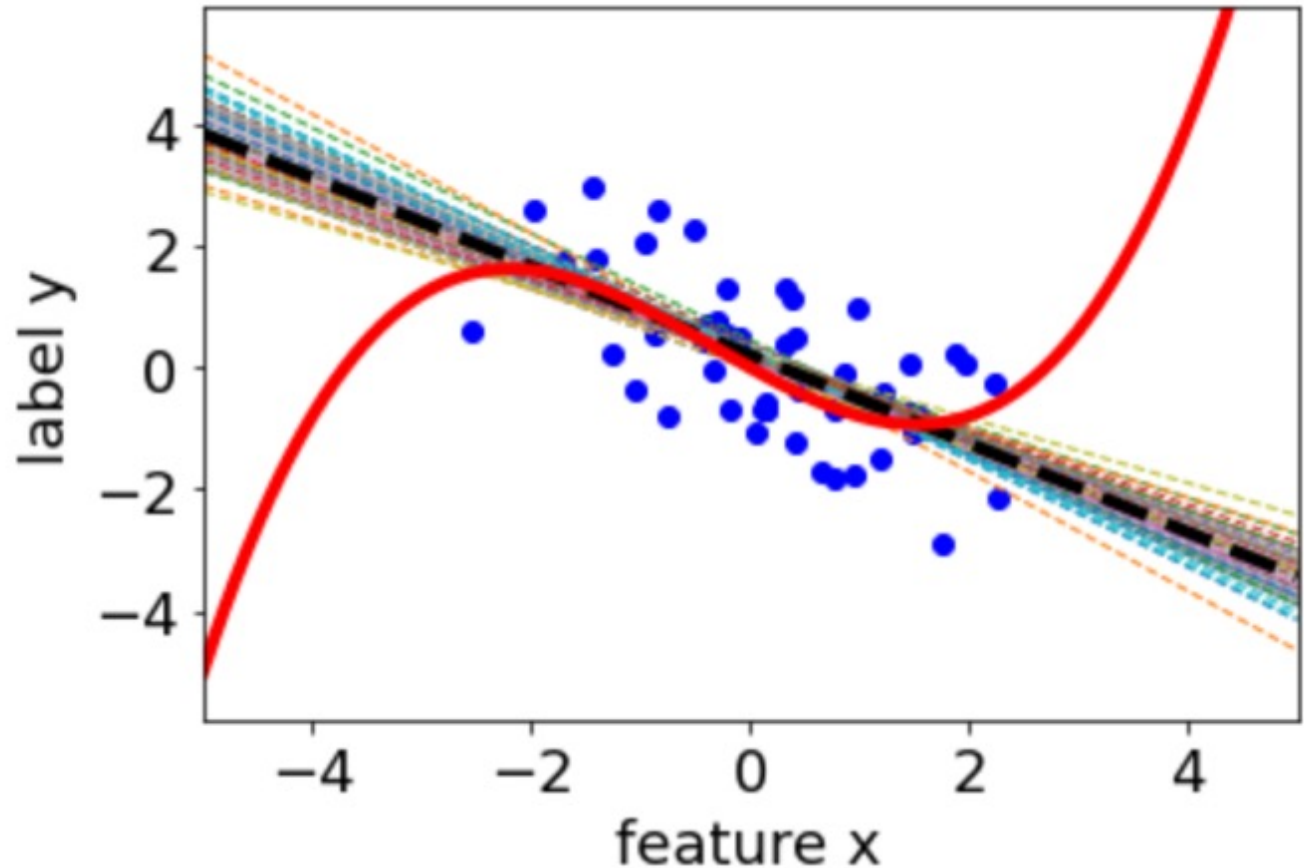
# Bias and Variance Tradeoff



"Prediction Error = Bias + Variance"

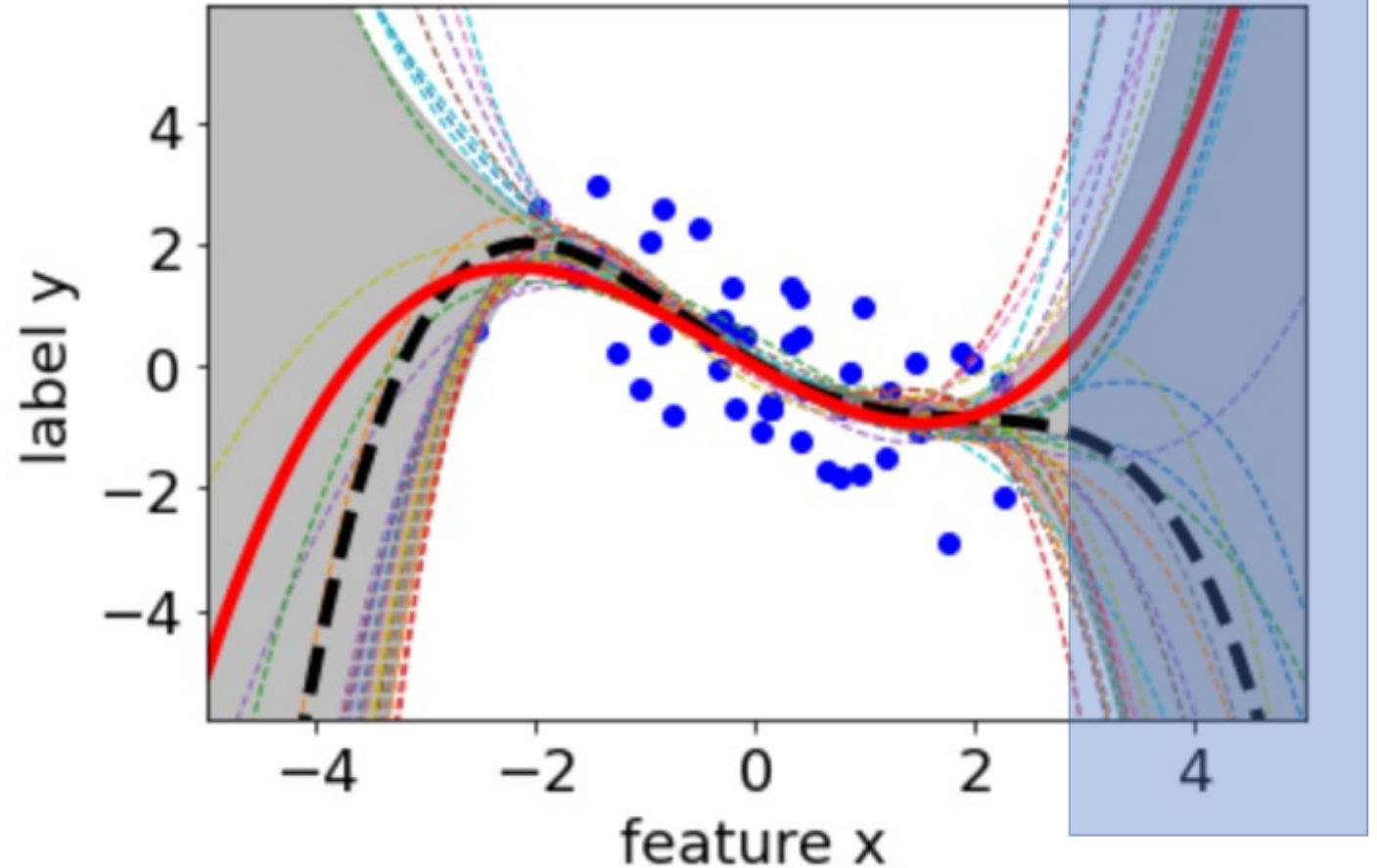bias reduction typically incurs variance increase and vice versa

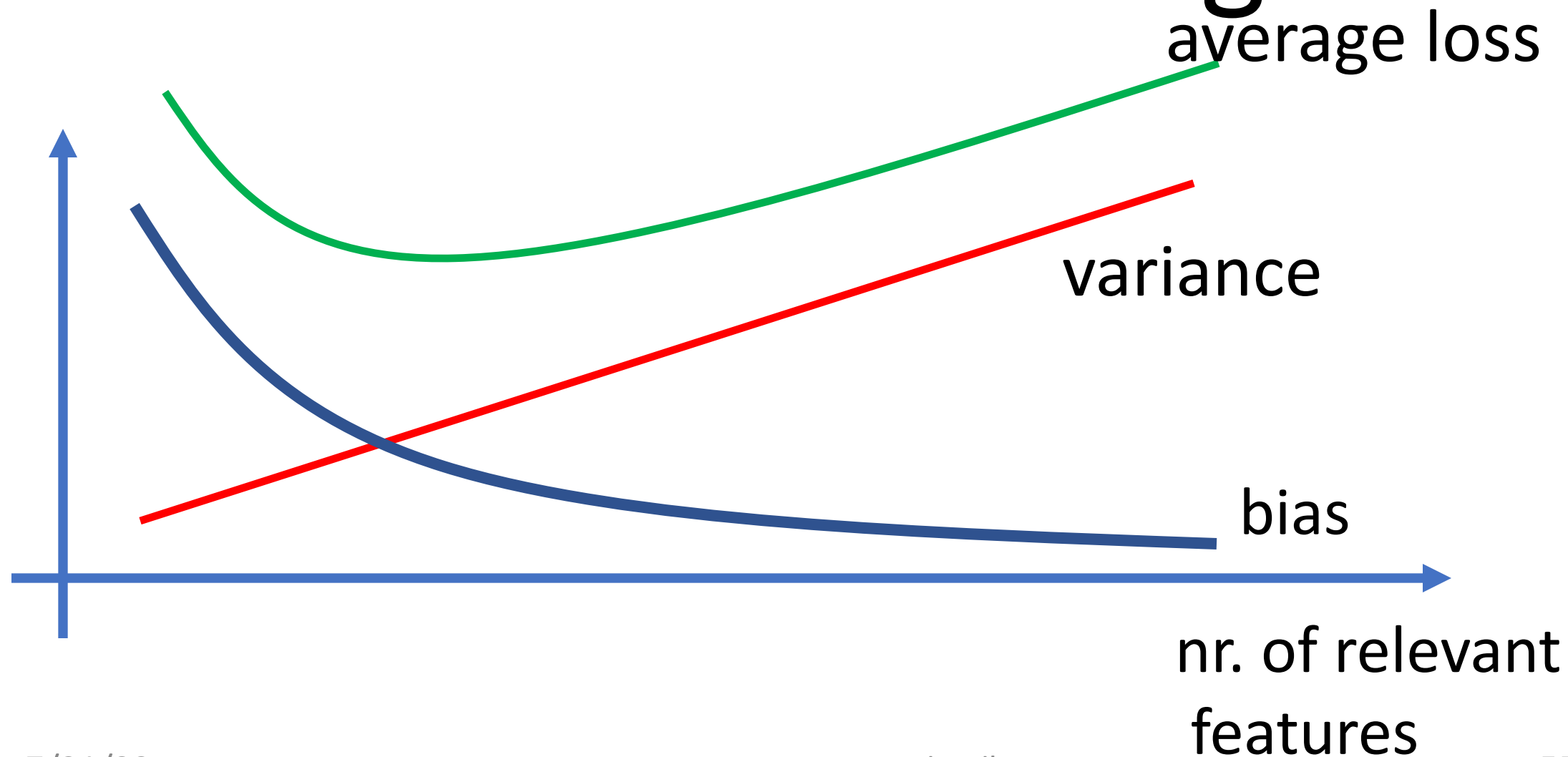# Smaller Model (Poly.Degree)

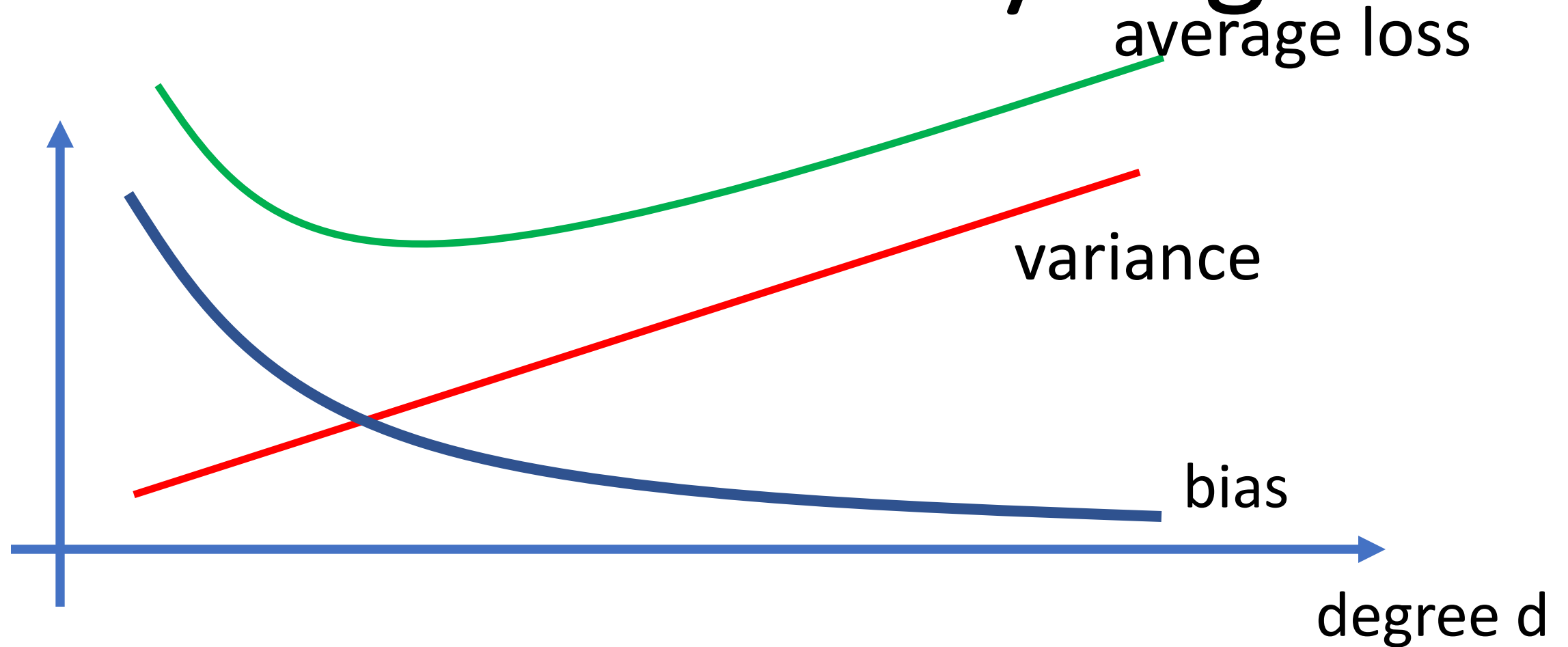- small variance

- large bias

# Larger Model (Poly. Degree)

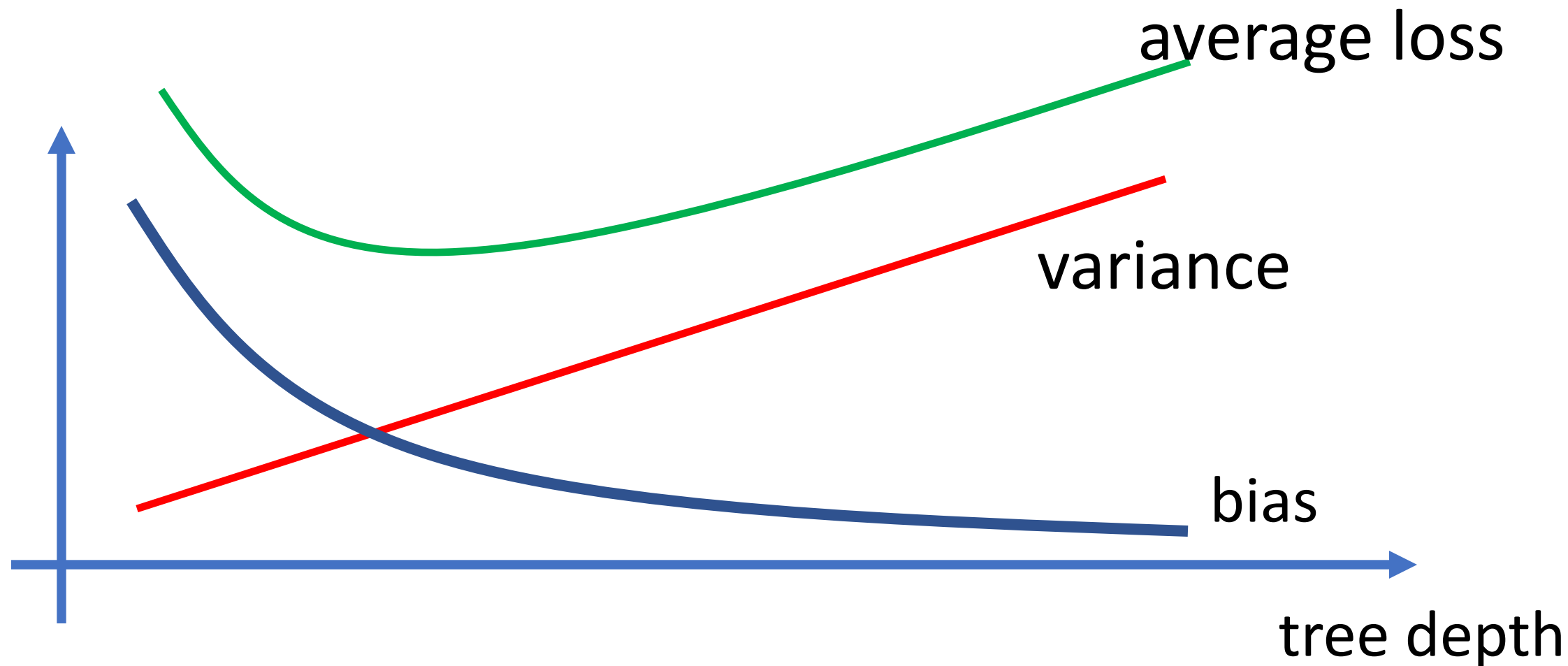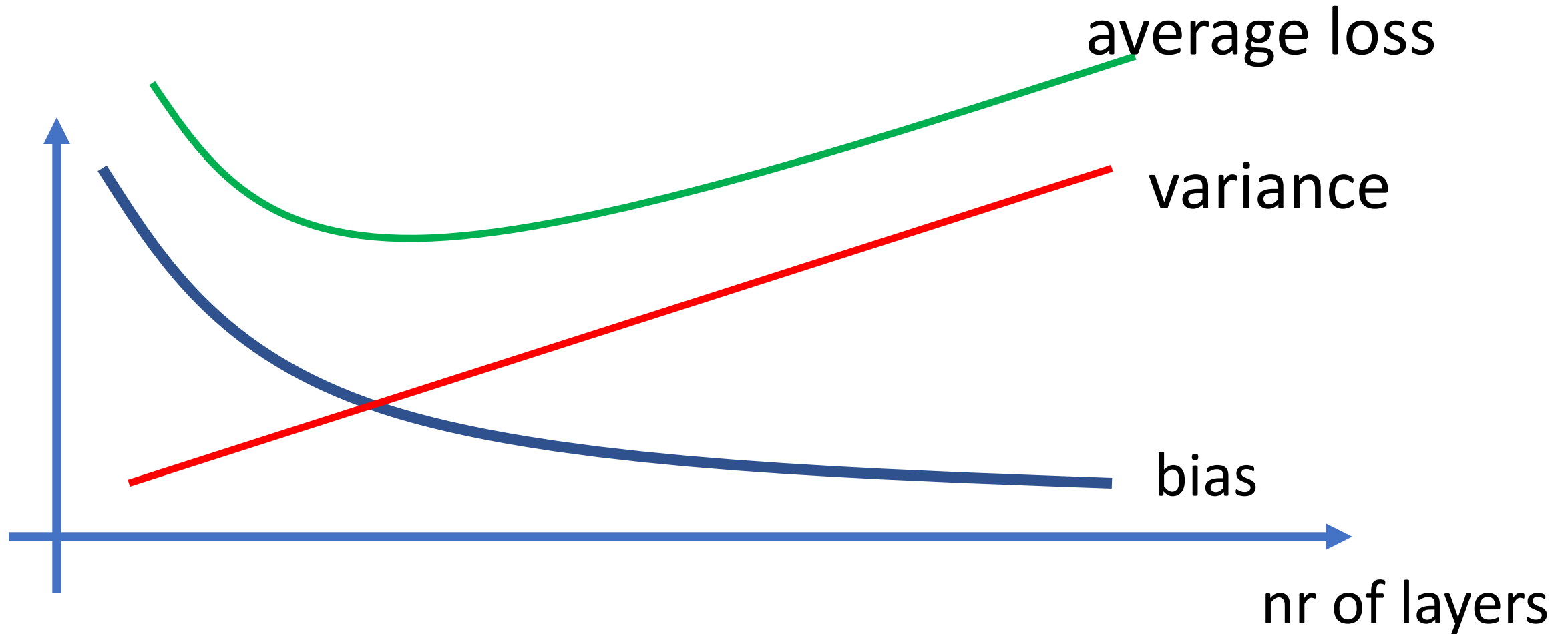- large variance

- small bias

# Bias vs. Variance Lin.Reg.
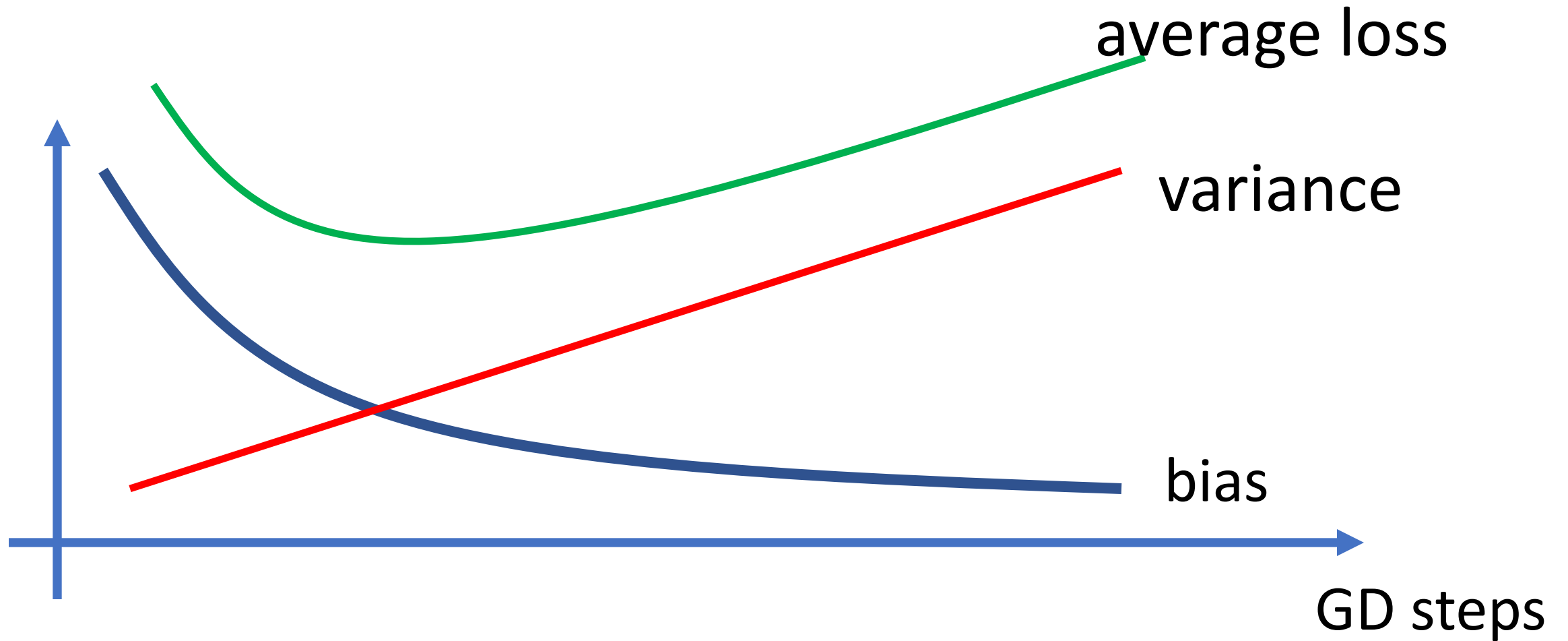
# Bias vs. Variance Poly.Reg.

# Bias vs. Variance Dec. Tree.

# Bias vs. Variance Deep Learning

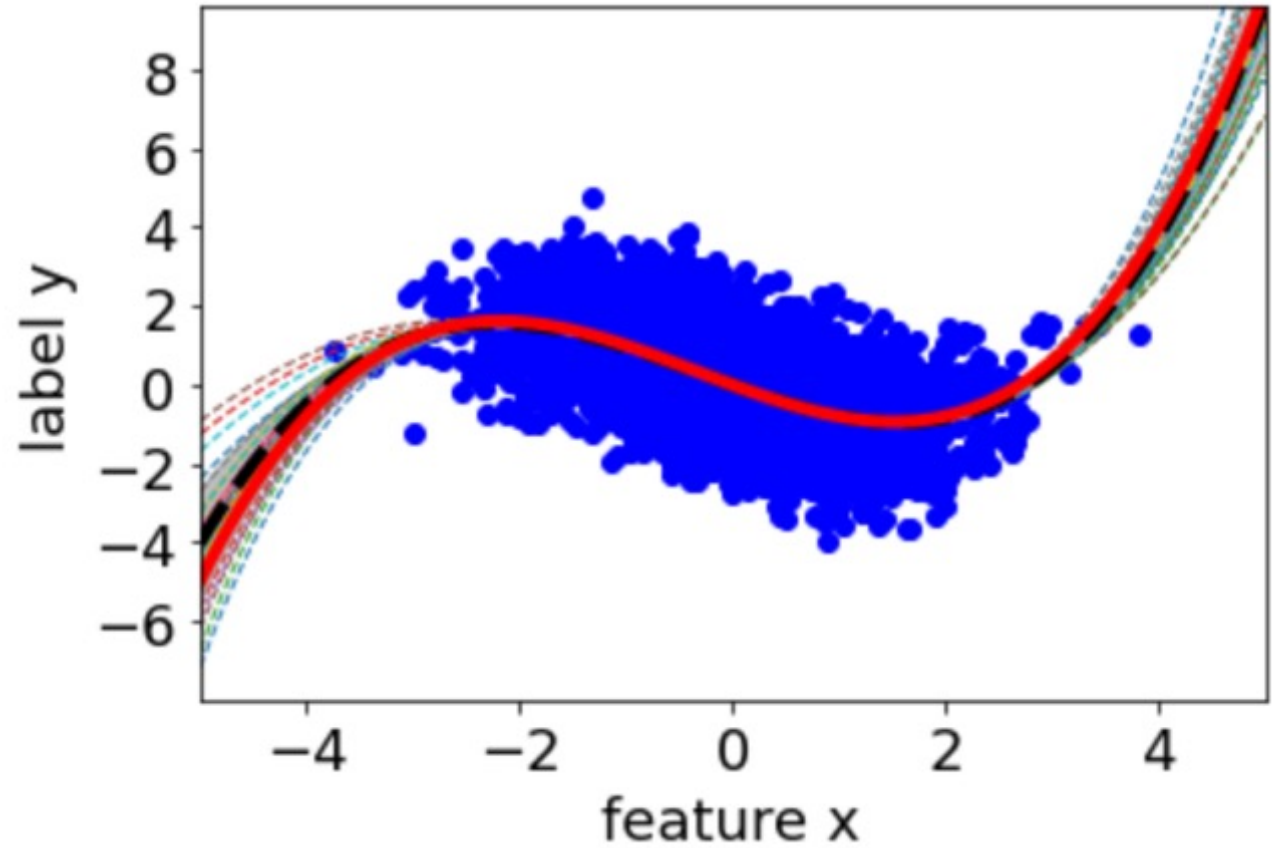# Bias vs. Variance Grad. Desc.



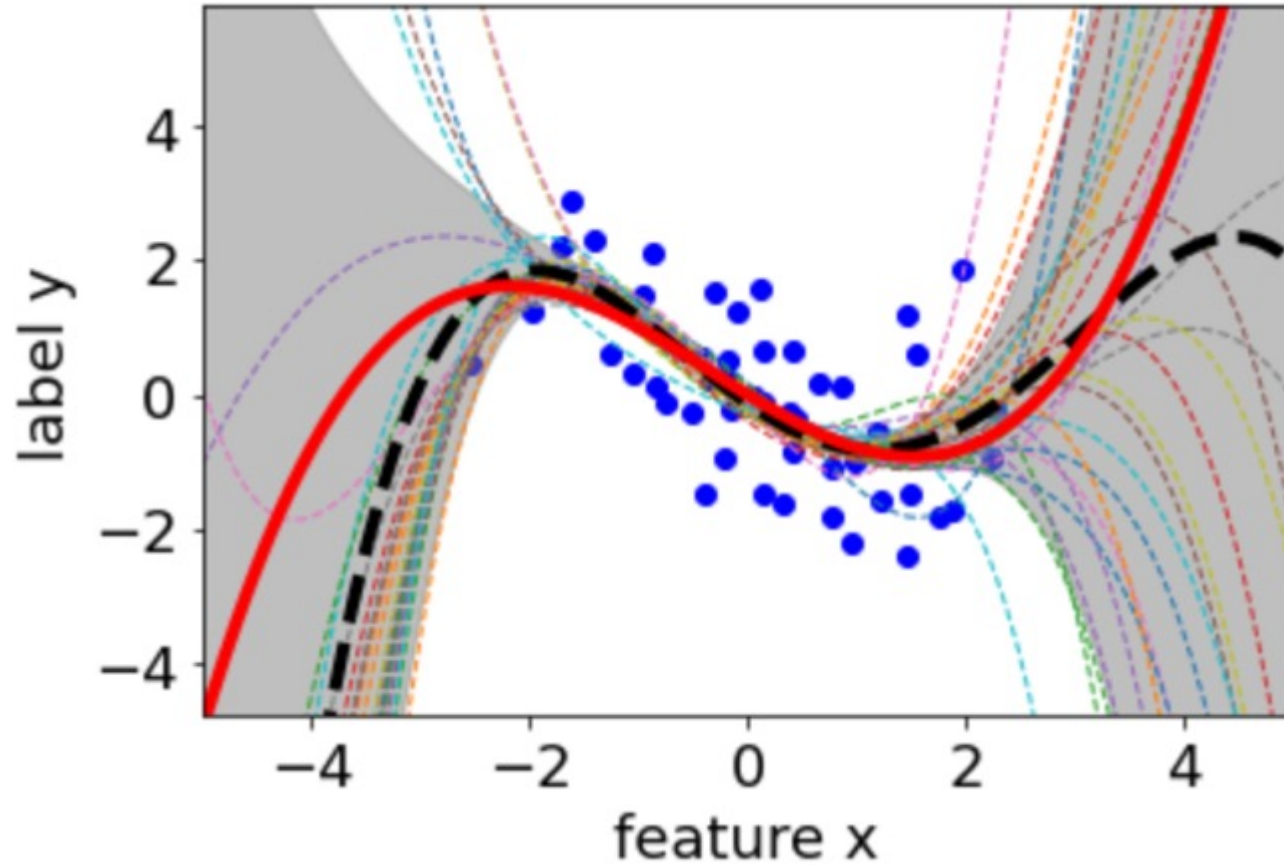average loss

variance

bias

GD steps

# More Data
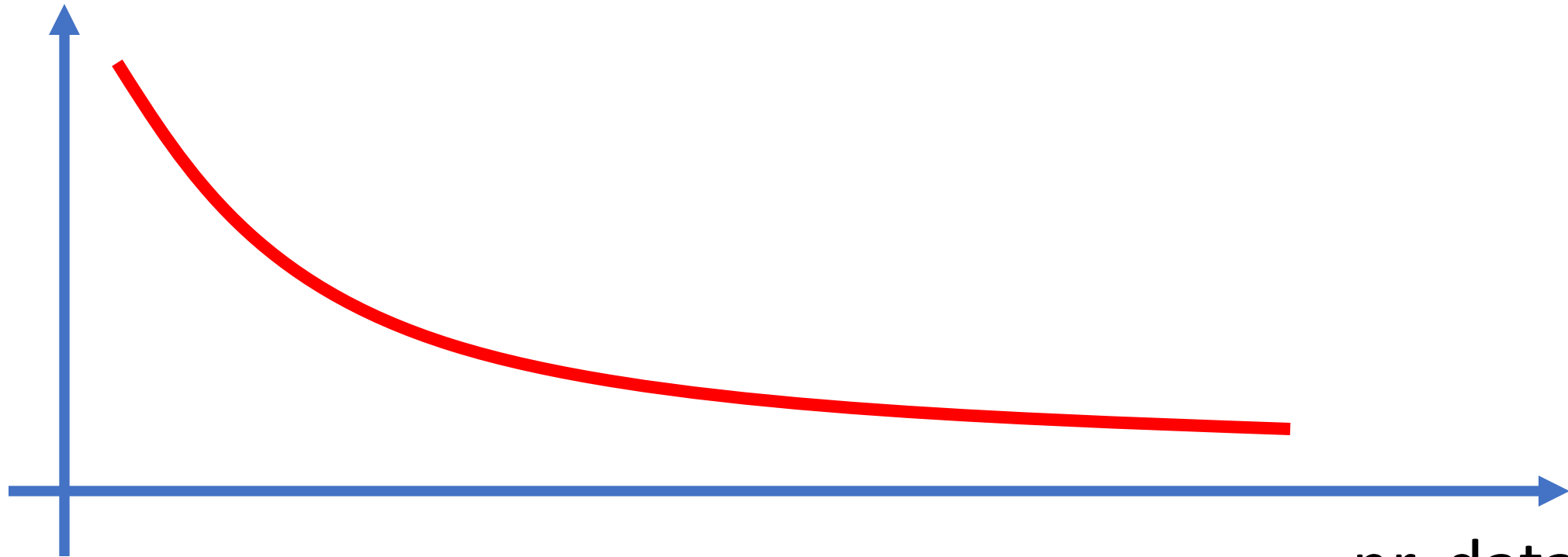
- small variance

# Less Data

- large variance

# Learning Curve



variance

nr. data points

# Alex' Rule of Thumb

effective number of training data points

\>

10 * nr. tunable effective model parameters

stretch the term "effective" as much as possible !

# ML Diagnosis

A. Jung HCML Summer School'22

# Simple Recipe

- consider ML method with some hypothesis space
- learn hypothesis by min. average loss on train.set
- training error = average loss of learnt hypothesis
- compute validation error
- compare val err, train err with a baseline
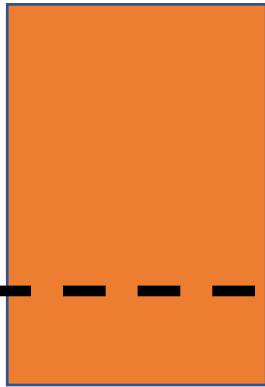
# Benchmark/Baseline

could be obtained from

- probabilistic models

- domain expertise

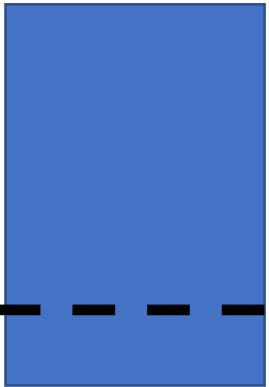- existing ML methods

- human performance

- ...

training
error

validation
error

- small train error -> hypothesis space is large
- large val err -> overfitting

- Workaround ?
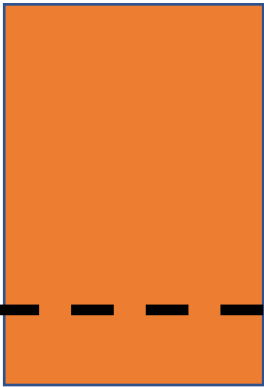
benchmark

training
error

validation
error

- large train error -> no good hypothesis found
- Workaround ?

training
error

validation
error

• Case Solved !

# Take Home Messages

- large models (e.g. deep nets) often overfit

- small training error does not mean much!

- diagnosis by comparing train/val err

- bias/variance analysis can guide model improvement

# Thank You !

A. Jung HCML Summer School'22