



# Statistical Machine Learning

## *Lecture 5 – Cross-validation and the bias-variance trade-off*



UPPSALA  
UNIVERSITET

**Johan Wågberg**

Division of Systems and Control

Department of Information Technology

Uppsala University

`johan.wagberg@it.uu.se`

## Summary of Lecture 4 (I/III)

---

**Linear Discriminant Analysis (LDA)** models the conditional class probabilities as

$$p(y = m | \mathbf{x}) = \frac{p(\mathbf{x} | y = m)p(y = m)}{\sum_{j=1}^M p(\mathbf{x} | y = j)p(y = j)}.$$

where

- $p(y = m) = \pi_m$  is the *prior* probability of class  $m$ .
- $p(\mathbf{x} | y = m) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_m, \boldsymbol{\Sigma})$  is the probability density of  $\mathbf{x}$  for an observation that comes from the  $m$ th class.

The **parameters** are:  $\pi_1, \dots, \pi_M, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_M, \boldsymbol{\Sigma}$

## Summary of Lecture 4 (II/III)

---

The parameters are estimated as the class frequencies and (within class) sample means and covariances, respectively,

$$\begin{aligned}\hat{\pi}_m &= \frac{n_m}{n} & m = 1, \dots, M, \\ \hat{\boldsymbol{\mu}}_m &= \frac{1}{n_m} \sum_{i:y_i=m} \mathbf{x}_i, & m = 1, \dots, M, \\ \hat{\boldsymbol{\Sigma}} &= \frac{1}{n - M} \sum_{m=1}^M \sum_{i:y_i=j} (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_m)(\mathbf{x}_i - \hat{\boldsymbol{\mu}}_m)^\top.\end{aligned}$$

The LDA classifier assigns a test input  $\mathbf{x}_*$  to class  $m$  with the maximum predicted probability  $p(y = m | \mathbf{x}_*)$ . LDA is a **linear classifier**.

## Summary of Lecture 4 (III/III)

---

**Parametric models** are specified using a fixed-dimensional vector of parameters.

**Non-parametric models** allow the flexibility of the model to grow with data.

---

One non-parametric model is the ***k*-nearest neighbour classifier**.

Given training data  $\mathcal{T} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ , for a test input  $\mathbf{x}_*$ ,

1. Identify the  $k$  training inputs  $\mathbf{x}_i$  closest to  $\mathbf{x}_*$ .
2. Classify  $\mathbf{x}_*$  according to a majority vote amongst these  $k$  training samples.



# Evaluating a supervised machine learning method

# Evaluating performance

---

**Training/learning:** Minimizes the average loss on training data. Loss functions differ between methods and the values can not directly be compared.

We would like to evaluate **already trained** models with respect to hyperparameters (eg. regularization parameter  $\gamma$  or  $k$  in  $k$ -NN) but also different methods (LDA, logistic regression, neural net).

---

## Error function

$$E(\hat{y}(\mathbf{x}), y) = \begin{cases} \mathbb{I}\{\hat{y}(\mathbf{x}) = y\} & \text{(classification)} \\ (\hat{y}(\mathbf{x}) - y)^2 & \text{(regression)} \end{cases}$$

*Note: Error function not necessarily the same as loss function.*

# Evaluating performance

Let  $\mathcal{T} = \{(x_i, y_i)\}_{i=1}^n$  be training data.

## Training error

$$E_{\text{train}} = \frac{1}{n} \sum_{i=1}^n E(\hat{y}(\mathbf{x}_i), y_i)$$

Measures how well a predictor  $\hat{y}$  performs on **training data**, but we are interested in **new data**. Let  $p(\mathbf{x}, y)$  be the joint distribution over data.

## Expected new data error

$$E_{\text{new}} = \mathbb{E}_{\star} [E(\hat{y}(x_{\star}), y_{\star})] = \int E(\hat{y}(x_{\star}), y_{\star}) p(x_{\star}, y_{\star}) dx_{\star} dy_{\star}$$

Impossible to compute since  $p(\mathbf{x}, y)$  is unknown, but minimizing  $E_{\text{new}}$  is our ultimate goal.

Maybe we can learn it?

# Approximating integrals

By the law of large numbers, we can approximate integrals using samples:

$$\mathbb{E}[h(\mathbf{x})] = \int h(\mathbf{x})p(\mathbf{x})d\mathbf{x} \approx \frac{1}{n} \sum_{i=1}^n h(\mathbf{x}_i), \quad \mathbf{x}_i \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}_i), \quad i = 1, \dots, n$$

With samples from  $p(\mathbf{x}, y)$ , we can estimate  $E_{\text{new}}$ !

*Note: Important that samples come from real world, “in-production” distribution.*

Training data are (or should be)  $n$  samples from  $p(\mathbf{x}, y)$ .

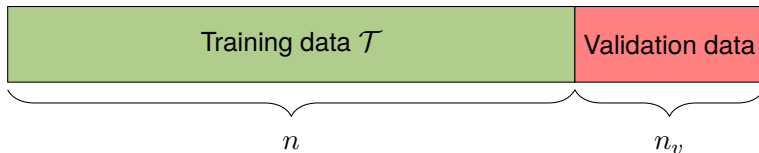
Approximate  $E_{\text{new}} \stackrel{?}{\approx} E_{\text{train}}$ ? **NO!**

Training data is part of the predictor  $\hat{y}(\mathbf{x}; \mathcal{T})$ !



# Estimating $E_{\text{new}}$ : hold-out validation data

Split data in **training data** and **validation data**.



**Hold-out validation error:**

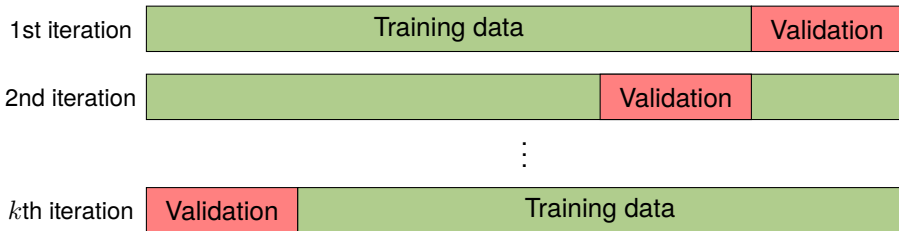
$$E_{\text{new}} \approx E_{\text{hold-out}} = \frac{1}{n_v} \sum_{i=1}^{n_v} E(\hat{y}(\mathbf{x}_i; \mathcal{T}), y_i)$$

- + Simple
- ! Good estimate of  $E_{\text{new}}$  requires the validation set to be large
- ! Good predictor  $\hat{y}$  requires the training set to be large
- Not all data is used for learning

Always split **randomly** between training and validation data!

# Estimating $E_{\text{new}}$ : $k$ -fold cross-validation

Split data in  $k$  batches and hold out batch  $\ell$  when estimating model. Use batch  $\ell$  to estimate  $E_{\text{new}}$  and average over all  $k$  estimates. Estimate final model using whole dataset.



## $k$ -fold cross-validation error

$$E_{\text{new}} \approx E_{k\text{-fold}} = \frac{1}{k} \sum_{\ell=1}^k E_{\text{hold-out}}^{(\ell)}$$

- + Gives a better estimate of  $E_{\text{new}}$
- Computationally more demanding than the hold-out data approach

## Using a test set

---

An important use of  $E_{k\text{-fold}}$  is to choose between models or select hyperparameters ( $k$  in  $k$ -NN or  $\lambda$  in regularization). If a good estimate of  $E_{\text{new}}$  is important:

We can no longer use  $E_{k\text{-fold}}$  to estimate  $E_{\text{new}}$ !

Set aside a test set and use **only** to estimate  $E_{\text{new}}$ .

# Flavors of cross-validation

---

- $k$ -fold cross-validation: Typically  $k \approx 10$
- Leave-one-out cross-validation:  $k$ -fold cross-validation with  $k = n$
- Monte Carlo cross-validation: Random selection of validation set at each iteration

# Understanding $E_{\text{new}}$

---

Collecting data is a random process where  $\mathcal{T}$  is sampled from  $p(\mathbf{x}, y)$ . Because  $\mathcal{T}$  is random, so is our learned model  $\hat{y}(\mathbf{x}; \mathcal{T})$ .

To better understand the behavior of  $E_{\text{new}}$ , we need to introduce

$$\begin{aligned}\bar{E}_{\text{train}} &= \mathbb{E}_{\mathcal{T}} [E_{\text{train}}], \\ \bar{E}_{\text{new}} &= \mathbb{E}_{\mathcal{T}} [E_{\text{new}}],\end{aligned}$$

where  $\mathbb{E}_{\mathcal{T}} [\cdot]$  is the average over **training data**  $\mathcal{T}$ .

*Note:  $k$ -fold cross-validation estimates  $\bar{E}_{\text{new}}$  rather than  $E_{\text{new}}$ .*

It usually holds that

$$\bar{E}_{\text{train}} < \bar{E}_{\text{new}}.$$



On average, a method usually performs better on training data than new data.



# Model complexity

Since  $\bar{E}_{\text{train}} < \bar{E}_{\text{new}}$ , define generalization gap:

$$\bar{E}_{\text{new}} = \bar{E}_{\text{train}} + \text{generalization gap}$$

A models ability to adapt to patterns in the data, we call the **model complexity**<sup>1</sup>.

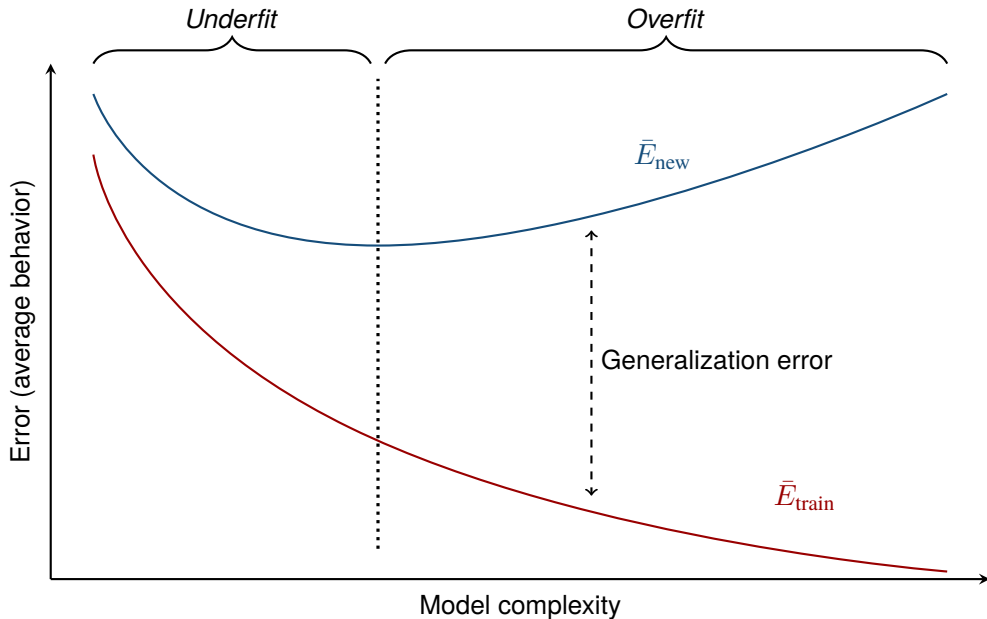
Model complexity ↗  
  $\bar{E}_{\text{train}}$   
 Generalization gap

Model complexity ↘  
  $\bar{E}_{\text{train}}$   
 Generalization gap

$\bar{E}_{\text{new}}$  usually attains a minimum at some intermediate complexity.

<sup>1</sup>A formal definition is possible. We will not do that in this course.

# Model complexity, $\bar{E}_{\text{train}}$ and $\bar{E}_{\text{new}}$







# Example: cross-validation for model selection

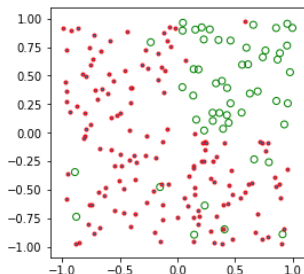
Binary classification,  $p = 2$ .

Evaluate the following methods

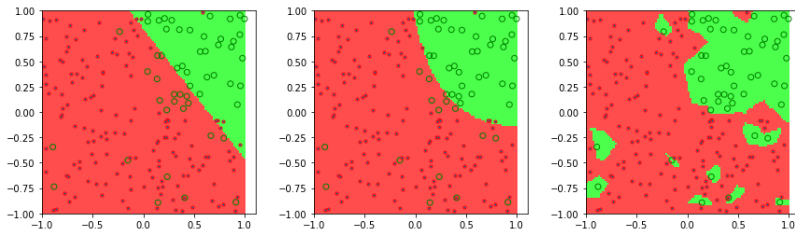
- ▶ logistic regression
- ▶ QDA
- ▶  $k$ -NN

We would like to pick the method with lowest  $E_{\text{new}}$ .

Since we only have the data, we can only *estimate*  $E_{\text{new}}$ .



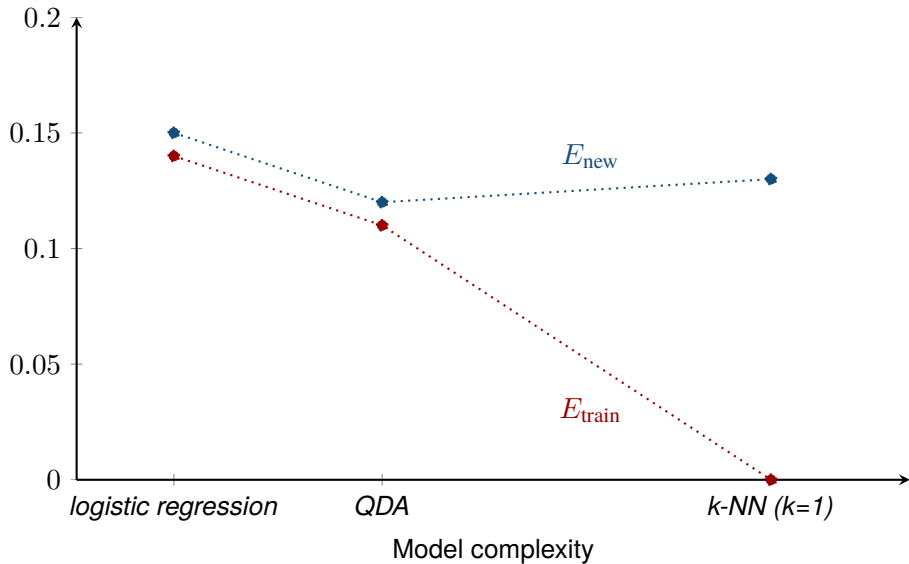
# Example: cross-validation for model selection



	logistic regression	QDA	$k$ -NN ( $k = 1$ )
$E_{\text{train}}$	0.14	0.11	0.0
$E_{\text{new}}^2$	0.15	0.12	0.13

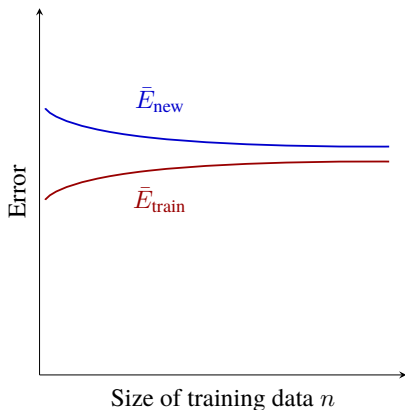
<sup>2</sup>Estimated using cross-validation

# Example: cross-validation for model selection

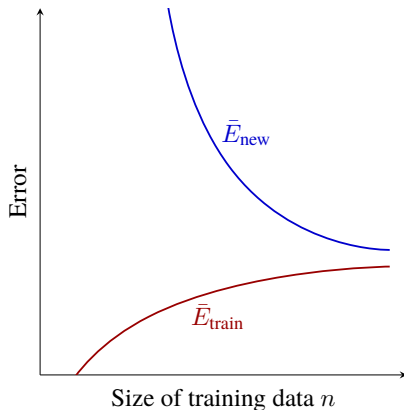


# Generalization error and training data size

In general, more training data  $n$  decreases the generalization error.  
 However,  $\bar{E}_{\text{train}}$  typically increases with  $n$ .



Low model complexity



High model complexity

# Digging even deeper to understand $E_{\text{new}}$

---

Let us focus on the regression setting.  
Assume that "the real world" works as

$$y = f_0(\mathbf{x}) + \epsilon, \text{ where } \begin{cases} \epsilon \text{ random, independent of } \mathbf{x} \\ \mathbb{E}[\epsilon] = 0, \\ \mathbb{E}[\epsilon^2] = \sigma^2. \end{cases}$$

Denote the **average trained model**

$$\bar{f}(\mathbf{x}) \triangleq \mathbb{E}_{\mathcal{T}} [\hat{y}(\mathbf{x}; \mathcal{T})].$$

The average model if we could re-train the model on new data an infinite number of times.

# Bias-variance decomposition

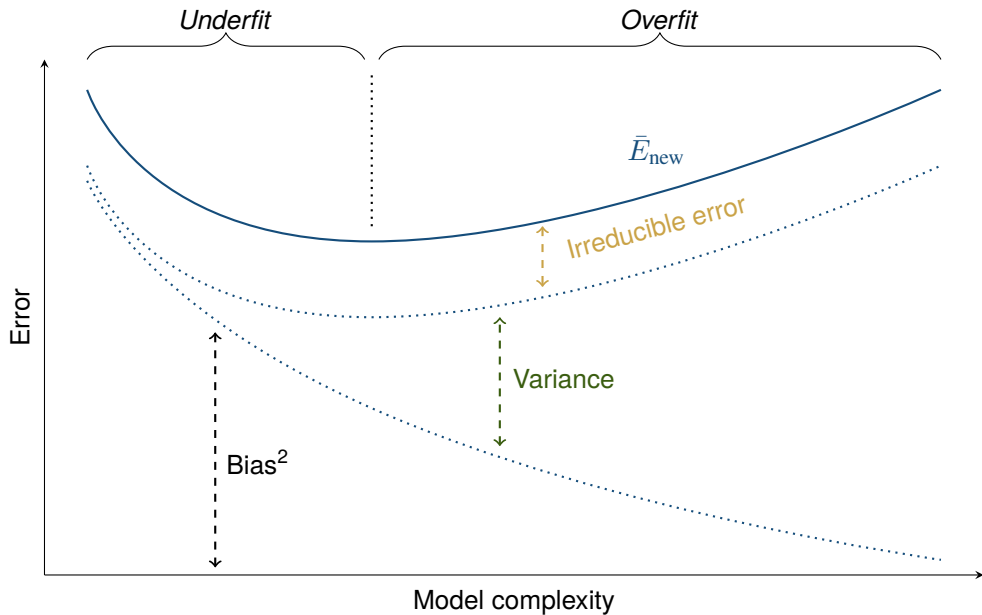
$$\begin{aligned}
 \bar{E}_{\text{new}} &= \mathbb{E}_{\mathcal{T}} \left[ \mathbb{E}_{\star} \left[ \left( \hat{y}(\mathbf{x}_{\star}; \mathcal{T}) - y_{\star} \right)^2 \right] \right] = \mathbb{E}_{\star} \left[ \mathbb{E}_{\mathcal{T}} \left[ \left( \hat{y}(\mathbf{x}_{\star}; \mathcal{T}) - f_0(\mathbf{x}_{\star}) - \epsilon \right)^2 \right] \right] \\
 &= \mathbb{E}_{\star} \left[ \mathbb{E}_{\mathcal{T}} \left[ \left( \hat{y}(\mathbf{x}_{\star}; \mathcal{T}) - \bar{f}(\mathbf{x}_{\star}) + \bar{f}(\mathbf{x}_{\star}) - f_0(\mathbf{x}_{\star}) - \epsilon \right)^2 \right] \right] \\
 &= \underbrace{\mathbb{E}_{\star} \left[ \mathbb{E}_{\mathcal{T}} \left[ \left( \hat{y}(\mathbf{x}_{\star}; \mathcal{T}) - \bar{f}(\mathbf{x}_{\star}) \right)^2 \right] \right]}_{\text{Variance}} + \underbrace{\mathbb{E}_{\star} \left[ \left( \bar{f}(\mathbf{x}_{\star}) - f_0(\mathbf{x}_{\star}) \right)^2 \right]}_{\text{Bias}^2} + \underbrace{\sigma^2}_{\text{Irreducible error}}
 \end{aligned}$$

*Technical interpretation:*

- **Bias<sup>2</sup>**  $\mathbb{E}_{\star} \left[ \left( \bar{f}(\mathbf{x}_{\star}) - f_0(\mathbf{x}_{\star}) \right)^2 \right]$ : The part of  $\bar{E}_{\text{new}}$  that is due to the fact that the model cannot represent the true  $f_0$ .
- **Variance**  $\mathbb{E}_{\star} \left[ \mathbb{E}_{\mathcal{T}} \left[ \left( \hat{y}(\mathbf{x}_{\star}; \mathcal{T}) - \bar{f}(\mathbf{x}_{\star}) \right)^2 \right] \right]$ : The part of  $\bar{E}_{\text{new}}$  that is due to the variability in the training dataset.

*Intuitive interpretation:*

- **Bias**: The inability of a method to describe the complicated patterns we would like it to describe. Low model complexity.
- **Variance**: How sensitive a method is to the training data. High model complexity.



Finding a balanced fit (neither over- nor underfit) is called the **the bias-variance tradeoff**.

# Regression example

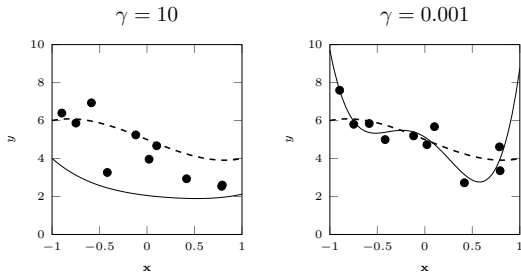
The data ( $n = 10$ ) comes from

$$y = 5 - 2x + x^3 + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, 1),$$

and our regression model is

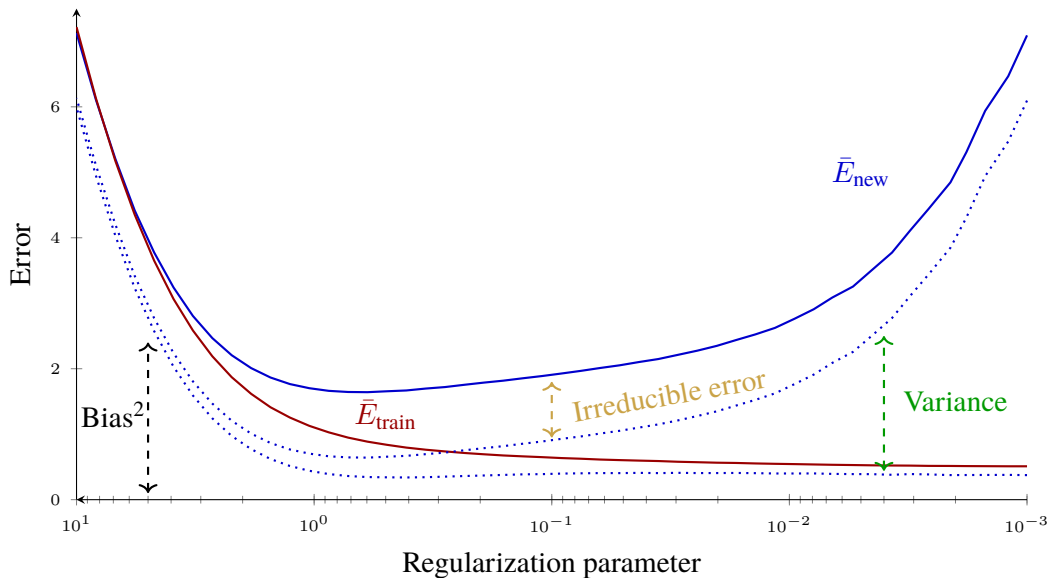
$$y = \theta_0 + \theta_1x + \theta_2x^2 + \theta_3x^3 + \theta_4x^4 + \varepsilon.$$

We use ridge regression to tune model complexity/bias-variance.

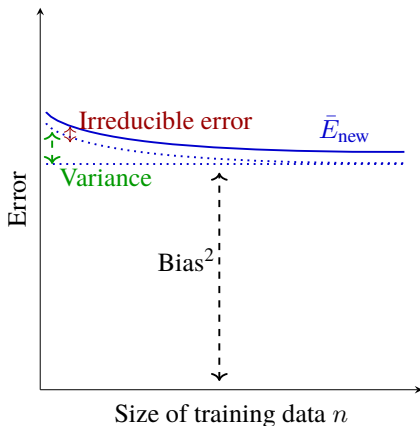




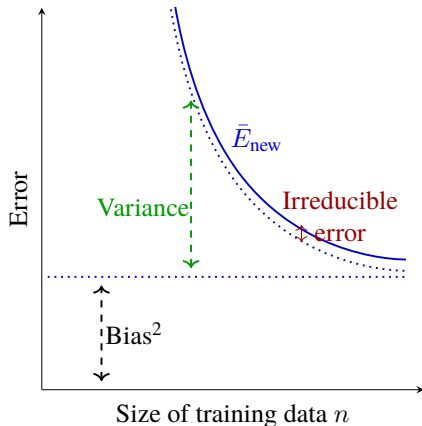
# Regression example



# Bias, variance and training data size



Low model complexity



High model complexity

# Make the tradeoff

---

Some methods for **decreasing the model complexity/increasing the bias/ decreasing the variance**:

- Increase  $k$  in  $k$ -NN
- Regularization
- Bagging
- Early stopping (for methods trained using optimization, notably deep learning)
- Dropout (deep learning)

Warning!  $\theta_0$  in linear regression (and later deep learning) is sometimes called “bias term”. That is **completely unrelated** to bias in this context.

# A few concepts to summarize lecture 5

---

$E(y, \hat{y})$ : Error function which compares predictions  $\hat{y}$  to true output  $y$ : MSE for regression, misclassification for classification.

$E_{\text{train}}$ : The training data error ( $E_{\text{train}}$  small = the method fits the training data well).

$E_{\text{new}}$ : The expected new data error; how well a method will perform when faced with an endless stream of new data.

**Cross-validation**: A method for estimating  $E_{\text{new}}$  using the training data.

**Model complexity**: How prone a method is to adapt to complicated patterns in the training data.

**Overfitting**: When a given method yields a smaller  $E_{\text{train}}$  and larger  $E_{\text{new}}$  than a model with lower model complexity would have done. That happens because the method/model is capturing patterns in the training data caused by random chance rather than true properties of the underlying function.

**Bias**: The inability of a method to describe the true patterns in the classification or regression problem. Low model complexity.

**Variance**: Sensitivity to random effects (noise) in the training data. High model complexity.