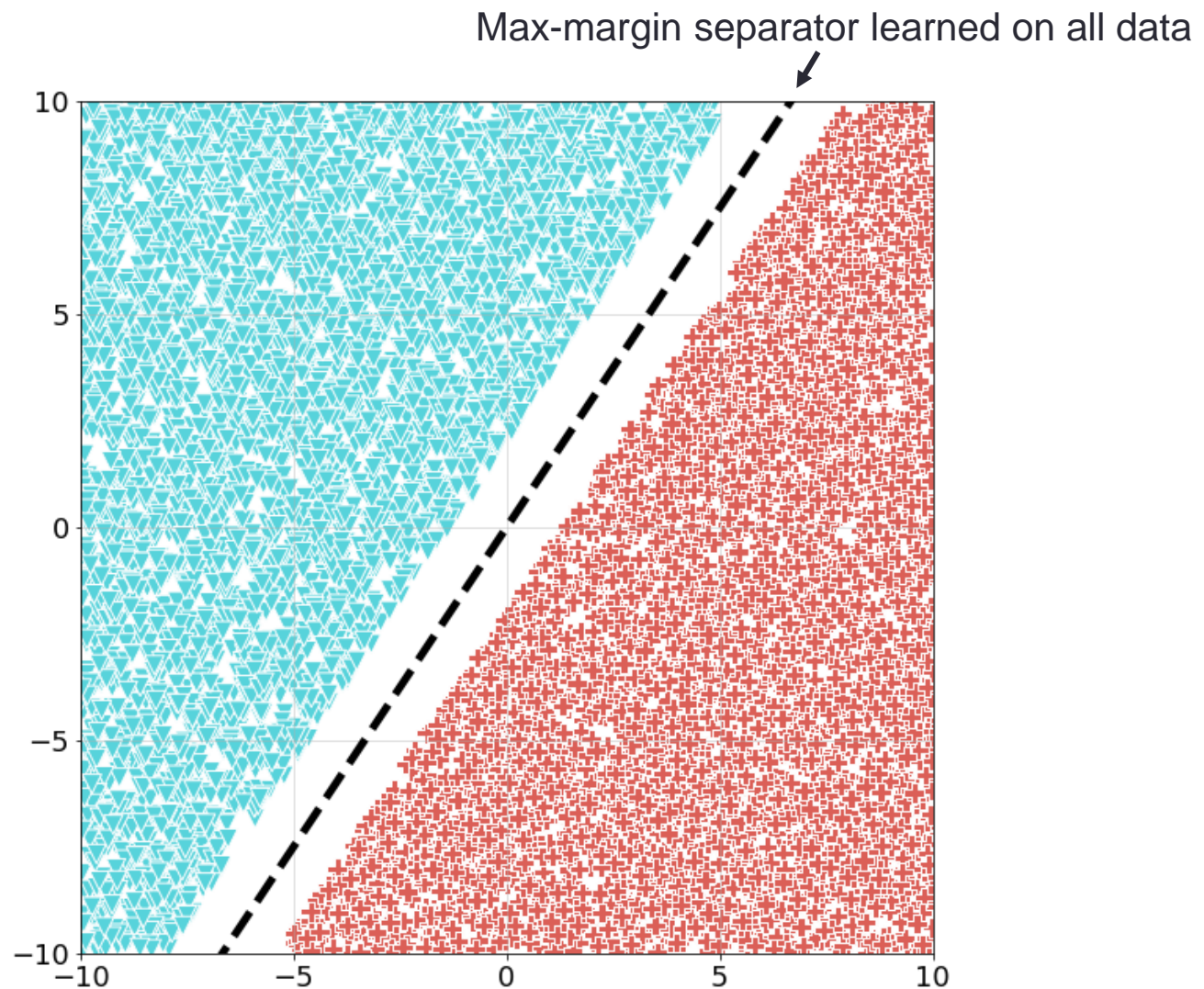# MODEL SELECTION

Itay Evron

# Outline

- Today's tutorial is different: more empirical than analytical

- Recap

  - Bias-variance error decomposition

  - Bias-variance tradeoff

- Demo I: Separable data

- Demo II: Inseparable data

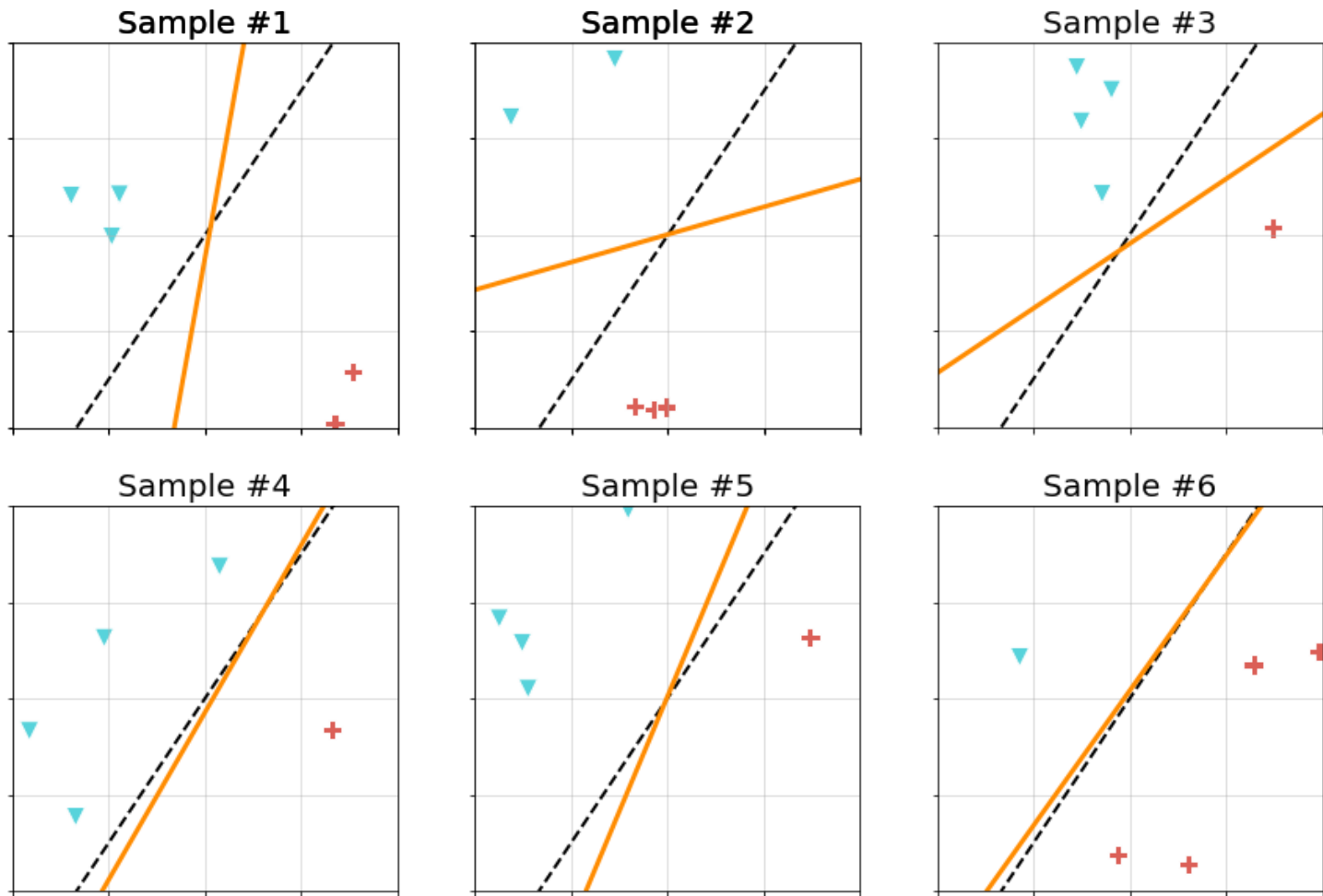- Model selection

# DEMO I:
# SEPARABLE DATA

# Linearly separable data

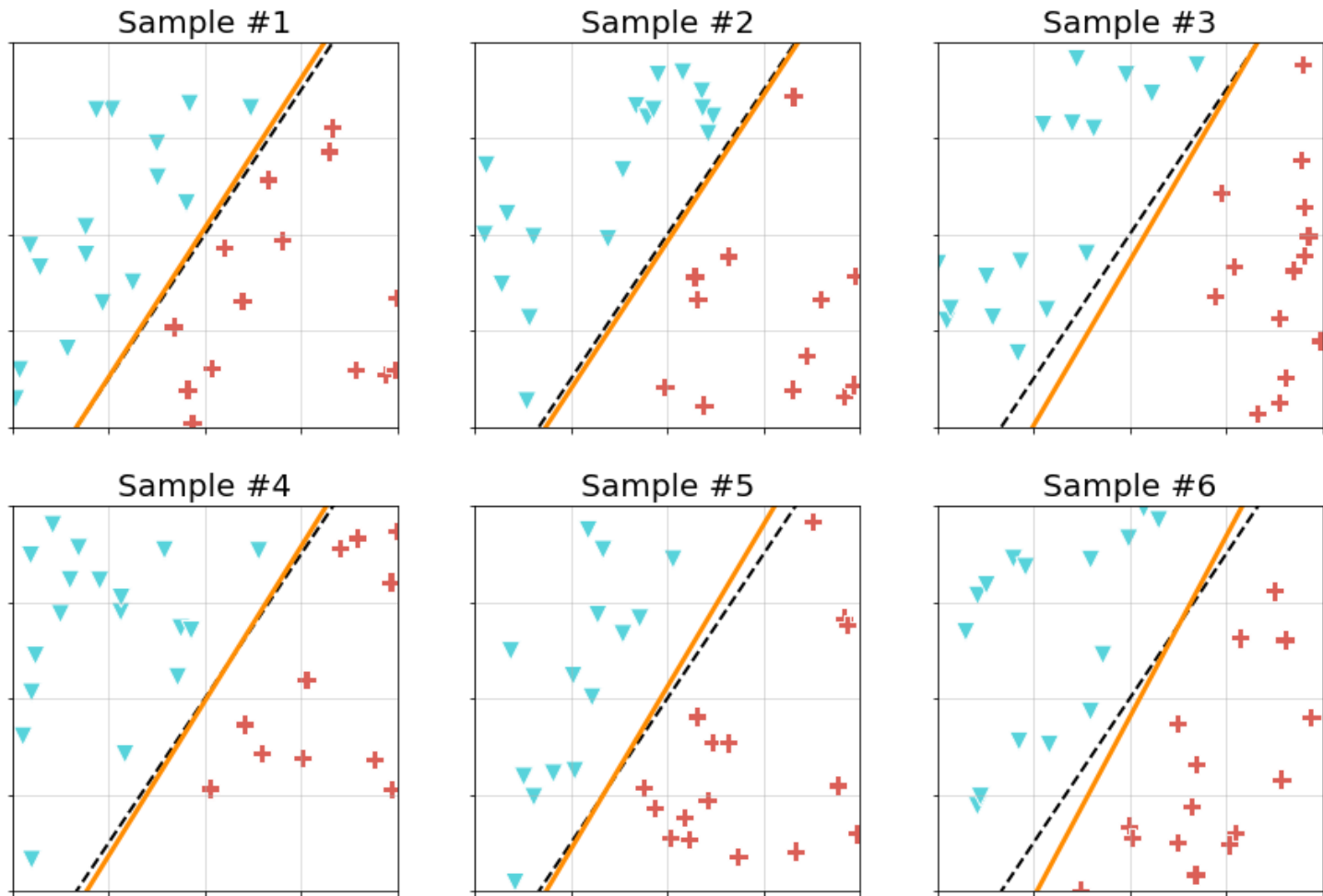Max-margin separator learned on all data
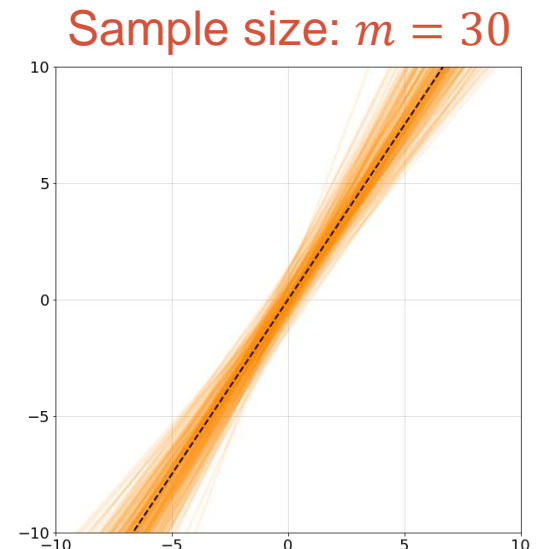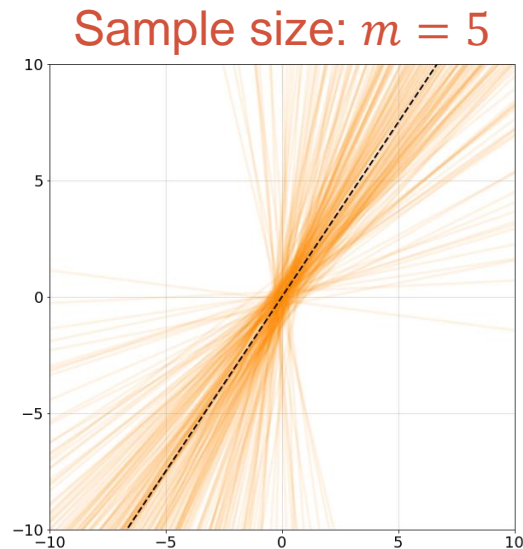


That's the entire distribution!
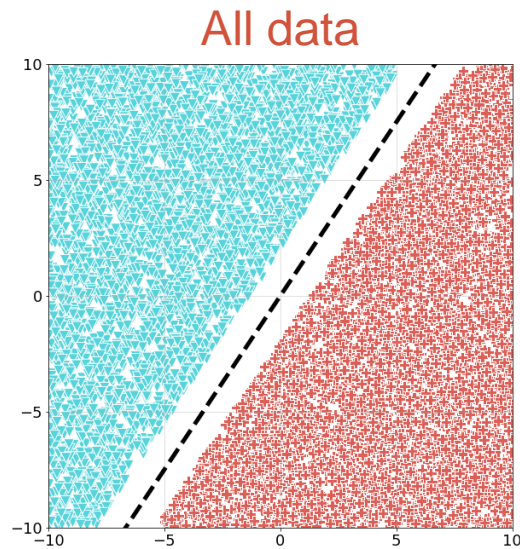
# Random samples when $m = 5$

# Random samples when $m = 30$

# More samples $\Rightarrow$ lower variance

### All data



### Sample size: $m = 5$



### Sample size: $m = 30$



### Everything is a random variable!

- The sample $S_i$
- The hypothesis $h_{S_i}$
- The loss $L_D(h_{S_i})$

High variance

Low bias



$L_D(h_{S_1})$

$L_D(h_{S_2})$

$\vdots$

$L_D(h_{S_k})$

# More samples $\Longrightarrow$ lower variance

All data

Sample size: $m = 5$
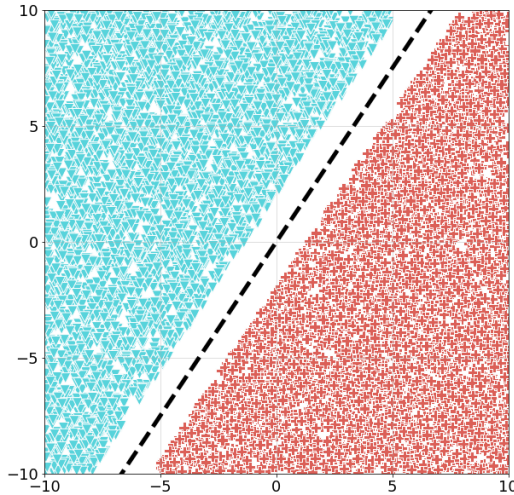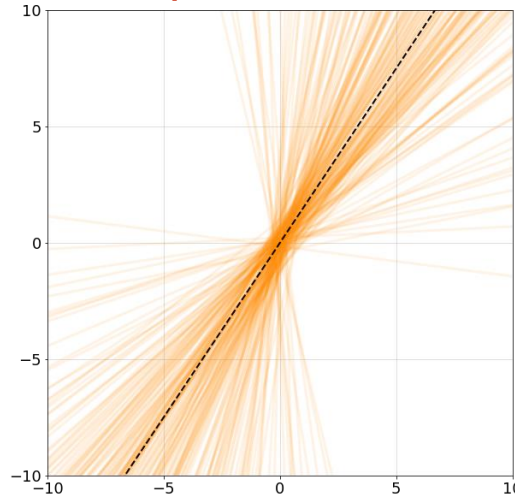
Sample size: $m = 30$

Everything is
a random variable!

- The sample $S_i$

- The hypothesis $h_{S_i}$

- The loss $L_D(h_{S_i})$

High variance

Low variance

Low bias

Low bias

# BIAS-VARIANCE ERROR DECOMPOSITION

# Recap: Bias-variance error decomposition

- Three interpretable sources of error:

$$\mathbb{E}_{S \sim D^m}\left[L_D^{sqr}(h_S)\right] = \mathbb{E}_x\left[\left(\bar{h}(x) - \bar{y}(x)\right)^2\right] + \mathbb{E}_{S,x}\left[\left(h_S(x) - \bar{h}(x)\right)^2\right] + \mathbb{E}_{x,y}\left[(\bar{y}(x) - y)^2\right]$$

***expected error***          ***bias²***          ***variance***          ***noise***

# Recap: Bias-variance error decomposition

- Three interpretable sources of error:

$$\mathbb{E}_{S \sim D^m}\left[L_D^{sqr}(h_S)\right] = \mathbb{E}_x\left[\left(\bar{h}(x) - \bar{y}(x)\right)^2\right] + \mathbb{E}_{S,x}\left[\left(h_S(x) - \bar{h}(x)\right)^2\right] + \mathbb{E}_{x,y}\left[(\bar{y}(x) - y)^2\right]$$

*expected error*         *bias²*        *variance*        *noise*

- **Noise:**

  - Property of data distribution (i.e., the statistical relation between $x$ and $y$)

  - In this tutorial: assume "realizability", i.e., $\exists f, \forall x$: $\underbrace{y = f(x)}_{\text{deterministic}}$ , i.e., **no noise**!

# Recap: Bias-variance error decomposition

- (without noise) ~~Three~~ Two interpretable sources of error:

$$\mathbb{E}_{S \sim D^m}\left[L_D^{sqr}(h_S)\right] = \mathbb{E}_{x,y}\left[\left(\bar{h}(x) - y\right)^2\right] + \mathbb{E}_{S,x}\left[\left(h_S(x) - \bar{h}(x)\right)^2\right] + \mathbb{E}_{x,y}\left[(\bar{y}(x) - y)^2\right]$$

**expected error**        **bias²**            **variance**        *noise*

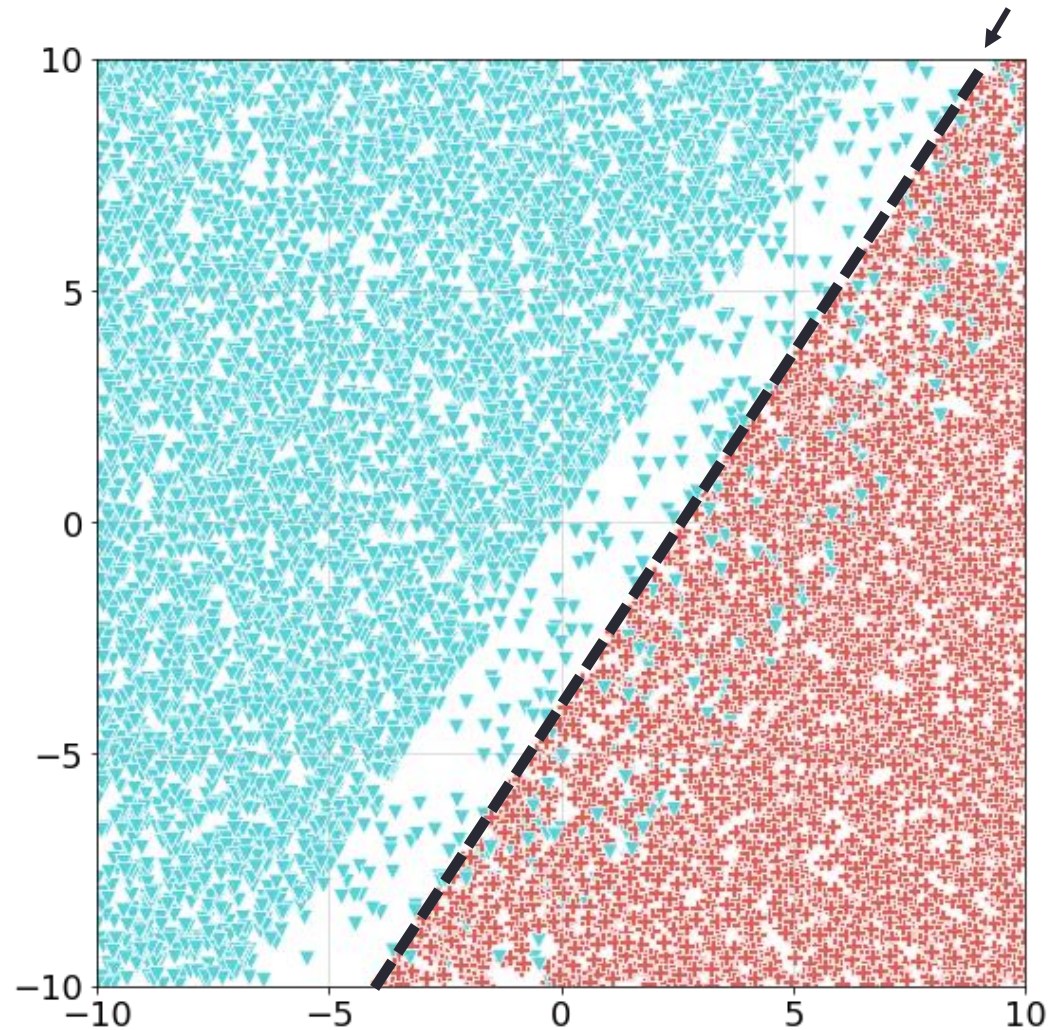- We will understand these quantities in the following slides.

# DEMO II: INSEPARABLE DATA

# Linearly <u>in</u>separable data

(linear) Separator with lowest generalization error



That's the entire distribution!

# Recap: Soft SVM

- Data is <u>not</u> linearly separable; hence we use Soft SVM

- Two conflicting objectives:

$$\underset{w\in\mathbb{R}^d, b\in\mathbb{R}}{\operatorname{argmin}} \; \lambda\|w\|_2^2 + \frac{1}{m}\sum_{i\in[m]} \max\{0, 1 - y_i(w^\top x_i + b)\}$$

$$\underset{w\in\mathbb{R}^d, b\in\mathbb{R}}{\operatorname{argmin}} \; \|w\|_2^2 + C\sum_{i\in[m]} \max\{0, 1 - y_i(w^\top x_i + b)\}$$

Complexity penalty (Regularization)
Prefer "simpler" models

Margin violation penalty
How much the $i^{th}$ example
*violates* the margin

# Recap: Kernel SVM

- To make things more interesting, we use an RBF kernel

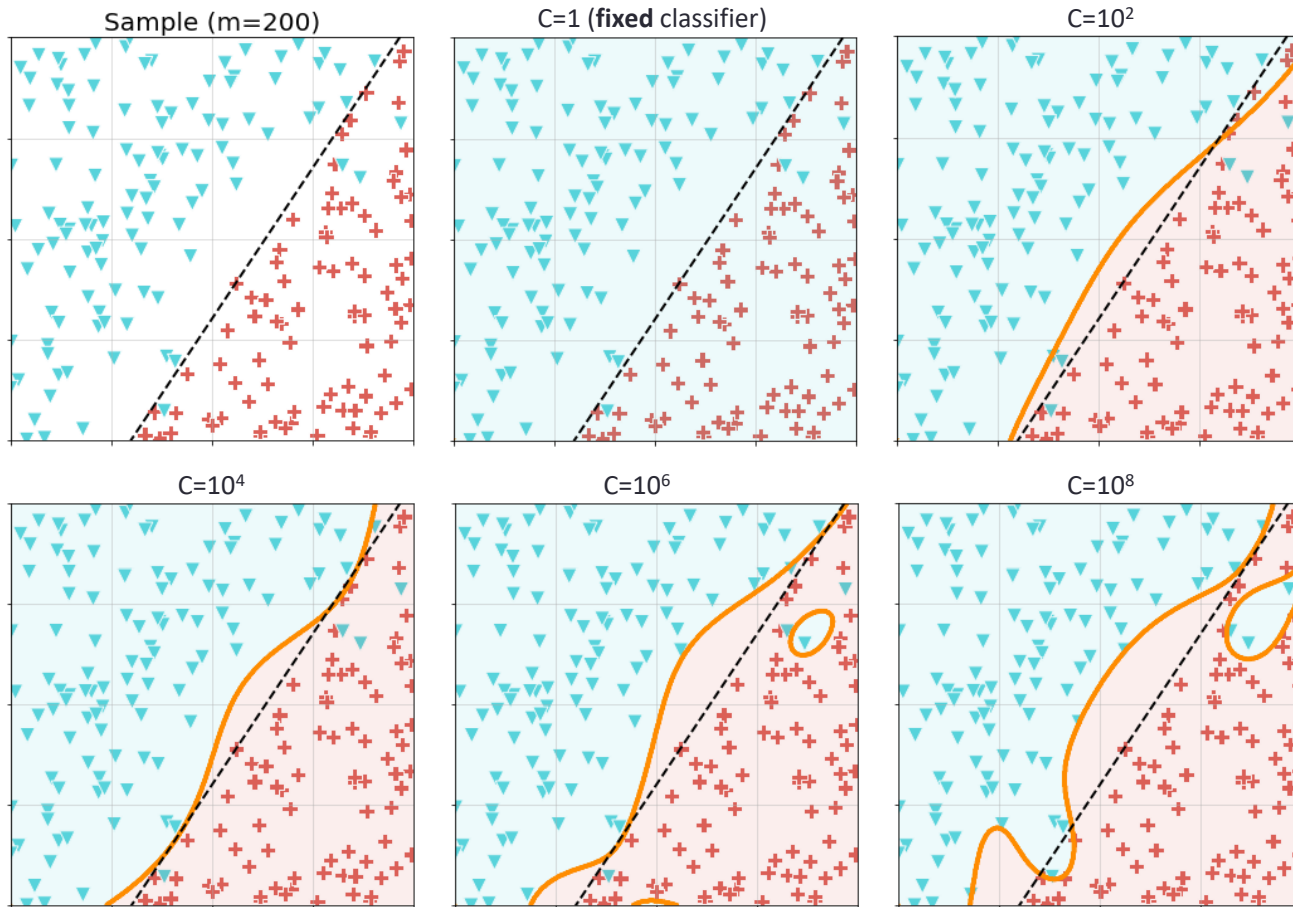- Solve an optimization problem equivalent to

$$\underset{w \in \mathbb{R}^{d'}, b \in \mathbb{R}}{\text{argmin}} \ \|w\|_2^2 + C \sum_{i \in [m]} \max\{0, 1 - y_i(w^\top \phi(x_i) + b)\}$$
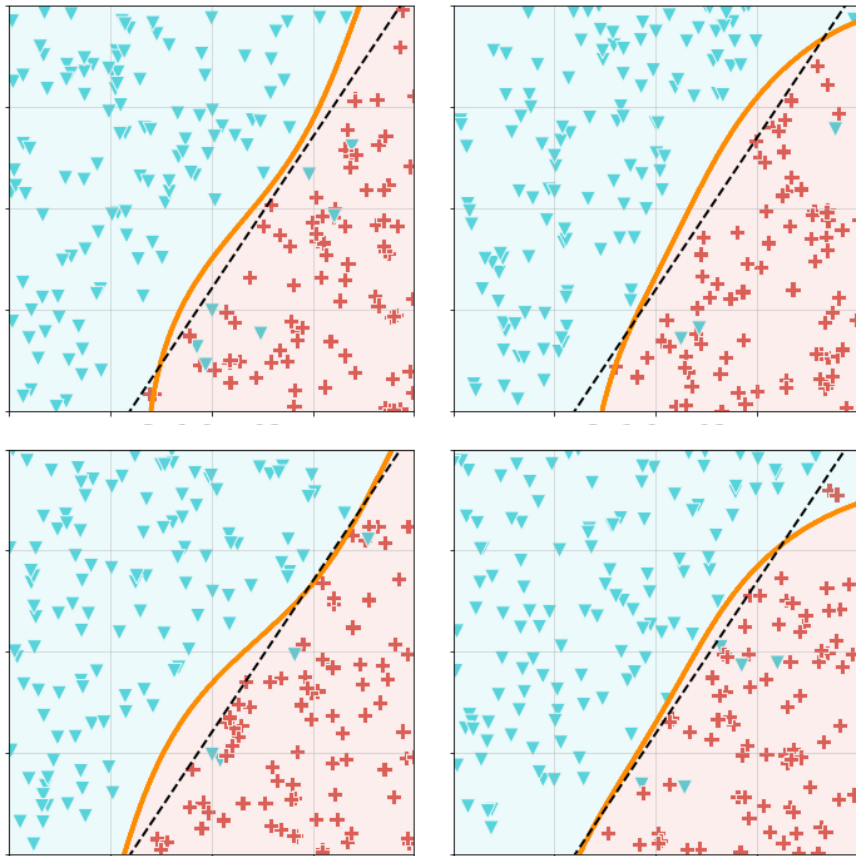
RBF feature mapping

# Larger $C \Rightarrow$ more complex models

$$\operatorname*{argmin}_{\boldsymbol{w}\in\mathbb{R}^{d'},b\in\mathbb{R}} \|\boldsymbol{w}\|_2^2 + C\sum_{i\in[m]} \max\{0, 1 - y_i(\boldsymbol{w}^\top \phi(\boldsymbol{x}_i) + b)\}$$
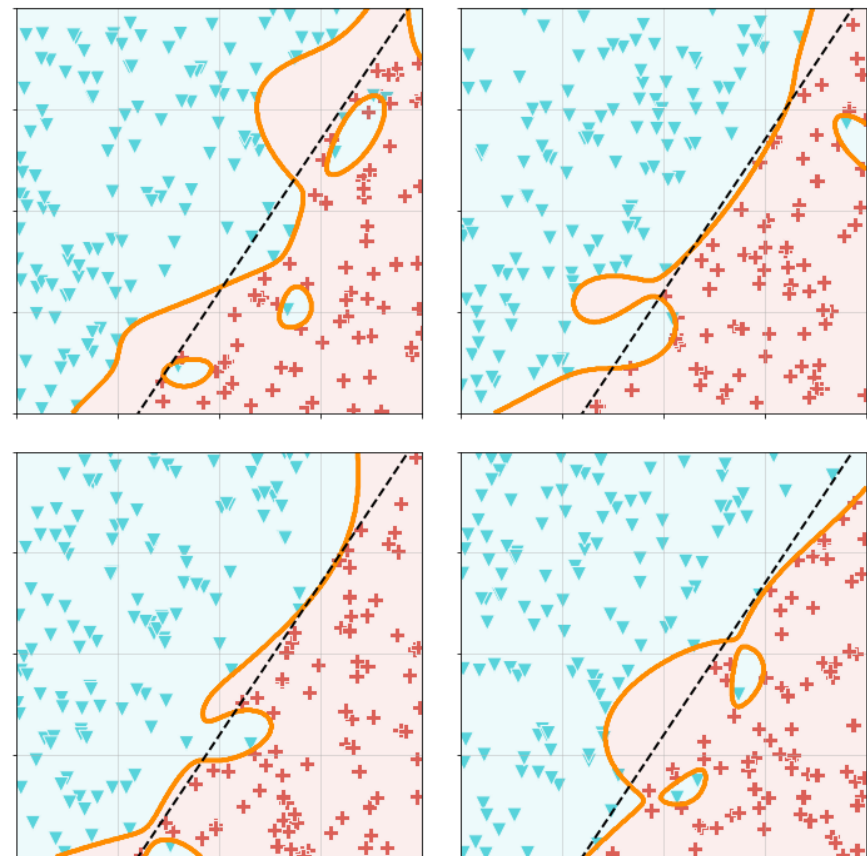
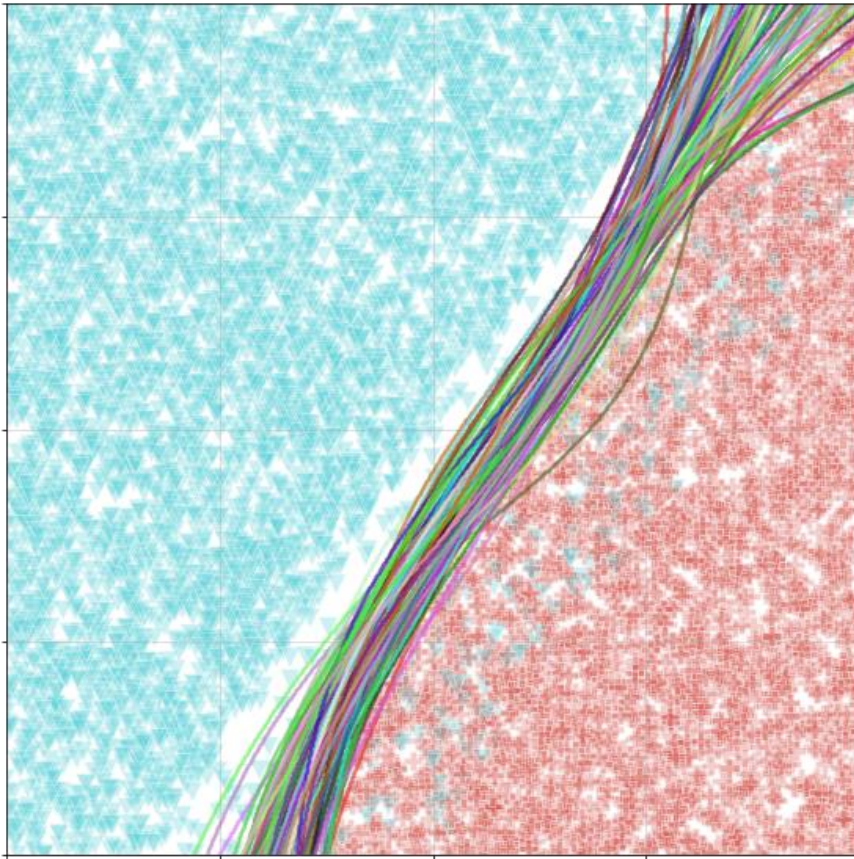# More complex models $\Rightarrow$ higher variance

Low complexity (C=$10^2$)

High complexity (C=$10^8$)

# More complex models $\Rightarrow$ higher variance

Low complexity ($C=10^2$)

High complexity ($C=10^8$)

# Formalizing the variance

- **Variance:** (of algorithm; w.r.t. $S$)

  - Measures how output hypotheses $h_S$ vary

    (how "sensitive" the learning algorithm is to changes in its input $S$)

  - Formally defined as: $\mathbb{E}_{S,x}\left[\left(h_S(x) - \bar{h}(x)\right)^2\right]$

  - Average hypothesis $\bar{h}$ as reference point (asks: relative to $\bar{h}$, how specialized is $h_S$ to $S$?)

- **The "average" hypothesis:** $\bar{h} = \mathbb{E}_{S \sim D^m}[h_S]$

One separator (C=10$^8$)  Many separators (C=10$^8$)  Average separator $\bar{h}$ (C=10$^8$)

# More complex models ⟹ higher variance



C=1 (**fixed** classifier)  C=$10^4$  C=$10^8$

**variance:** $\mathbb{E}_{S,x}\left[\left(h_S(x) - \bar{h}(x)\right)^2\right]$

Complexity (controlled by $C$)

# Formalizing the bias

- **Bias:**
  - Quantifies how well our hypothesis <u>class</u> fits the data (on average)
  - Formally defined as: $\mathbb{E}_{x,y}\left[\left(\bar{h}(x) - y\right)^2\right]$
  - Does not depend on sampled data (but does depend on data size)

One separator (C=$10^8$)   Many separators (C=$10^8$)   Average separator $\bar{h}$ (C=$10^8$)

# More complex models $\implies$ less bias



One separator
C=$10^8$

Many separators
C=$10^8$

Average separator $\bar{h}$
C=$10^8$

**bias²:** $\mathbb{E}_{x,y}\left[\left(\bar{h}(x) - y\right)^2\right]$

Complexity (controlled by $C$)

# More complex models $\Rightarrow$ less bias



Most separators
C=1 (**fixed** classifier)

**bias²:** $\mathbb{E}_{x,y}\left[\left(\bar{h}(x) - y\right)^2\right]$

Complexity (controlled by $C$)

# More complex models $\Rightarrow$ less bias



One separator

Many separators

Average separator $\bar{h}$

C=5

C=5

C=5

**bias²:** $\mathbb{E}_{x,y}\left[\left(\bar{h}(x) - y\right)^2\right]$

Complexity (controlled by $C$)

# The sweet spot of the expected error



One separator
C=10³

Many separators
C=10³

Average separator $\bar{h}$
C=10³

*sweet spot*

**expected error: bias² + variance**

**variance:** $\mathbb{E}_{S,x}\left[\left(h_S(x) - \bar{h}(x)\right)^2\right]$

**bias²:** $\mathbb{E}_{x,y}\left[\left(\bar{h}(x) - y\right)^2\right]$

Complexity (controlled by $C$)

# Cases we saw



C=1 (**fixed** classifier)

Low variance

High bias

**expected error: bias² + variance**

**variance:** $\mathbb{E}_{S,x}\left[\left(h_S(x) - \bar{h}(x)\right)^2\right]$

**bias²:** $\mathbb{E}_{x,y}\left[\left(\bar{h}(x) - y\right)^2\right]$

Complexity (controlled by $C$)

# Cases we saw



$C=10^8$

High variance

Low bias

**expected error: bias² + variance**

**variance:** $\mathbb{E}_{S,x}\left[\left(h_S(x) - \bar{h}(x)\right)^2\right]$

**bias²:** $\mathbb{E}_{x,y}\left[\left(\bar{h}(x) - y\right)^2\right]$

Complexity (controlled by $C$)

# Cases we saw



$C=10^3$

Low variance

Low bias

expected error: bias² + variance

variance: $\mathbb{E}_{S,x}\left[\left(h_S(x) - \bar{h}(x)\right)^2\right]$

bias²: $\mathbb{E}_{x,y}\left[\left(\bar{h}(x) - y\right)^2\right]$

Complexity (controlled by $C$)

# Bias-variance tradeoff

expected error: **bias² + variance**

variance: $\mathbb{E}_{S,x}\left[\left(h_S(x) - \bar{h}(x)\right)^2\right]$

bias²: $\mathbb{E}_{x,y}\left[\left(\bar{h}(x) - y\right)^2\right]$

*complexity (capacity)*

*large bias*          *"sweet spot"*          *large variance*
*(optimal complexity)*

# MODEL SELECTION

# Bias-variance error decomposition

- Three interpretable sources of error:

$$\mathbb{E}_{S \sim D^m}\left[L_D^{sqr}(h_S)\right] = \mathbb{E}_x\left[\left(\bar{h}(x) - \bar{y}(x)\right)^2\right] + \mathbb{E}_{S,x}\left[\left(h_S(x) - \bar{h}(x)\right)^2\right] + \mathbb{E}_{x,y}\left[(\bar{y}(x) - y)^2\right]$$

*expected error*                *bias²*                        *variance*                        *noise*

- In practice, we want to perform model selection:

  tune the model complexity to achieve a low expected error

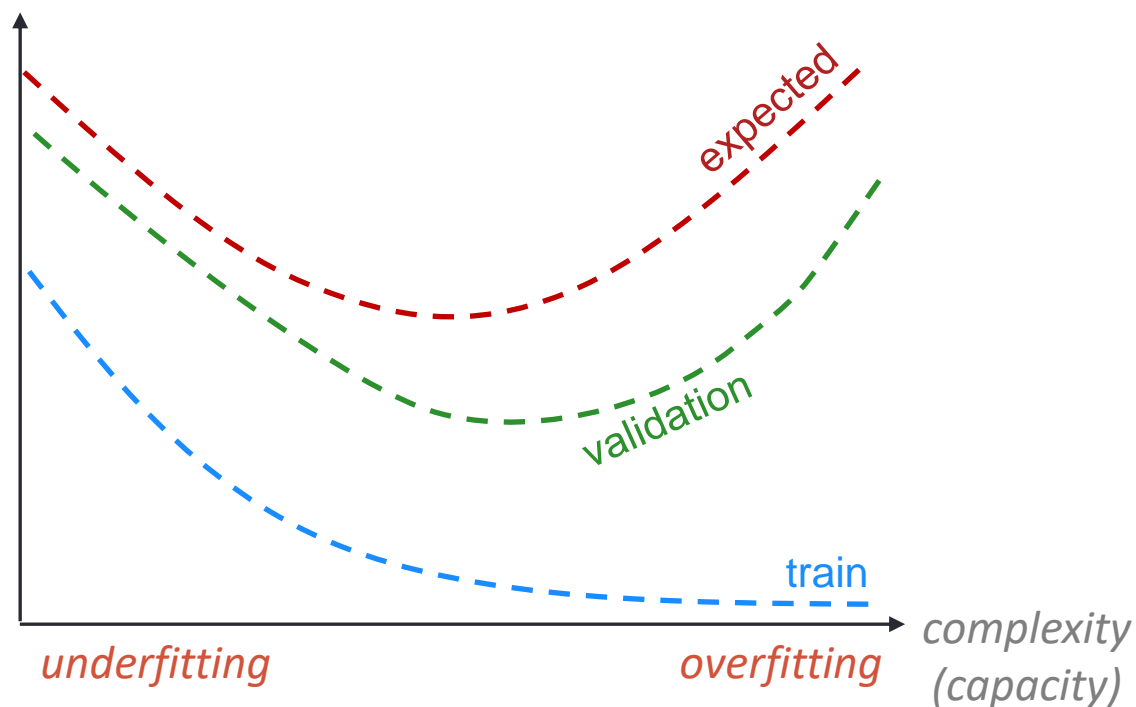- However, we cannot compute the actual expected error!

# Validation curve

- Here, we train on 100 examples and put 20 examples aside for validation

- The validation error is an estimator for the generalization error

# Error decomposition using validation

- Given a training set $S$, a validation set $V$, and a hypothesis $h_S$, the generalization error can be decomposed into:

$$L_D(h_S) = \big(L_D(h_S) - L_V(h_S)\big) + \big(L_V(h_S) - L_S(h_S)\big) + L_S(h_S)$$



*underfitting*         *overfitting*   *complexity (capacity)*

# Error decomposition using validation

- Given a training set $S$, a validation set $V$, and a hypothesis $h_S$, the generalization error can be decomposed into:

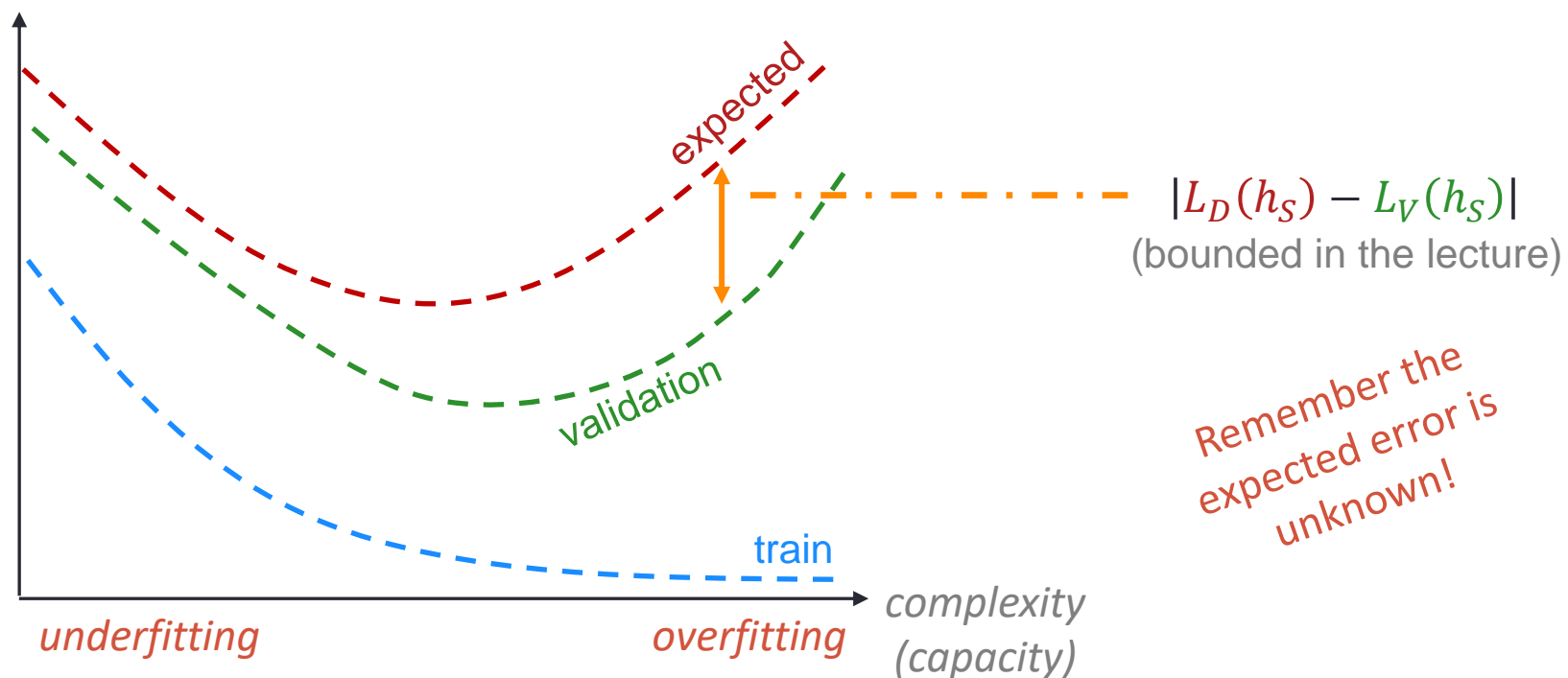$$L_D(h_S) = \left(L_D(h_S) - L_V(h_S)\right) + \left(L_V(h_S) - L_S(h_S)\right) + L_S(h_S)$$



$|L_D(h_S) - L_V(h_S)|$
(bounded in the lecture)

*expected*

*validation*

*train*

Remember the expected error is unknown!

*underfitting*          *overfitting*          *complexity (capacity)*
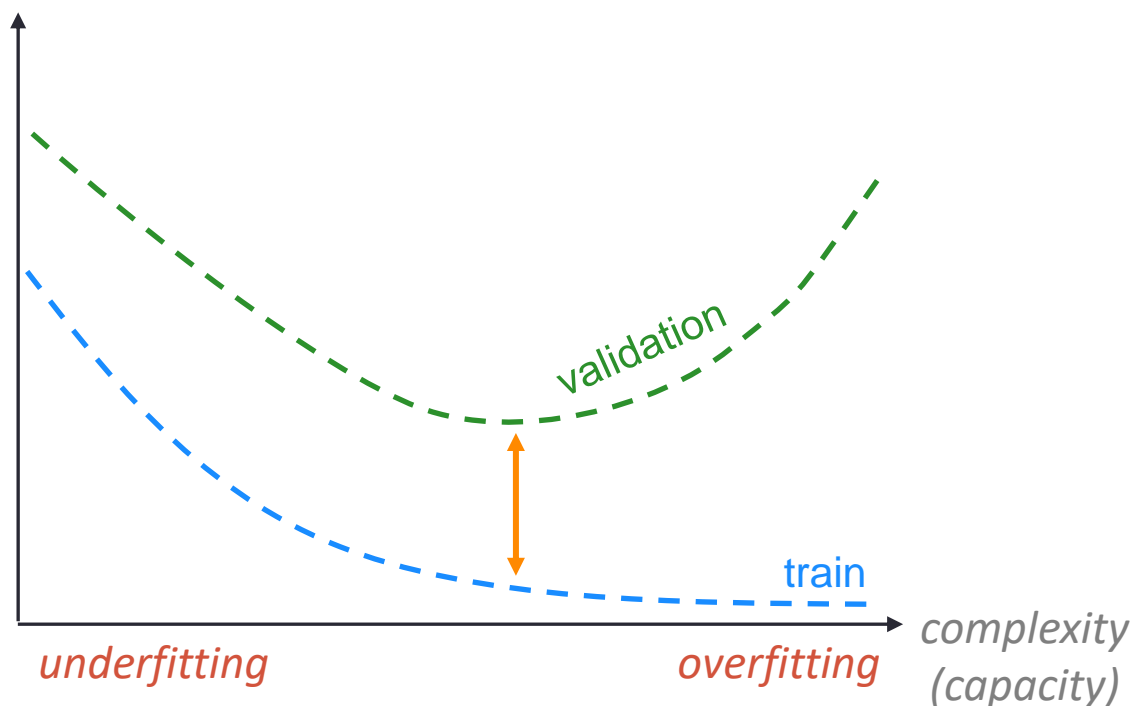
# Error decomposition using validation

- Given a training set $S$, a validation set $V$, and a hypothesis $h_S$, the generalization error can be decomposed into:

$$L_D(h_S) = \big(L_D(h_S) - L_V(h_S)\big) + \big(L_V(h_S) - L_S(h_S)\big) + L_S(h_S)$$

If this term is large, $h_S$ probably overfits.

Possible solutions:

- Get more samples

- Feature selection

- Lower the capacity

- Change regularization type



validation

train

*underfitting*          *overfitting*          *complexity (capacity)*
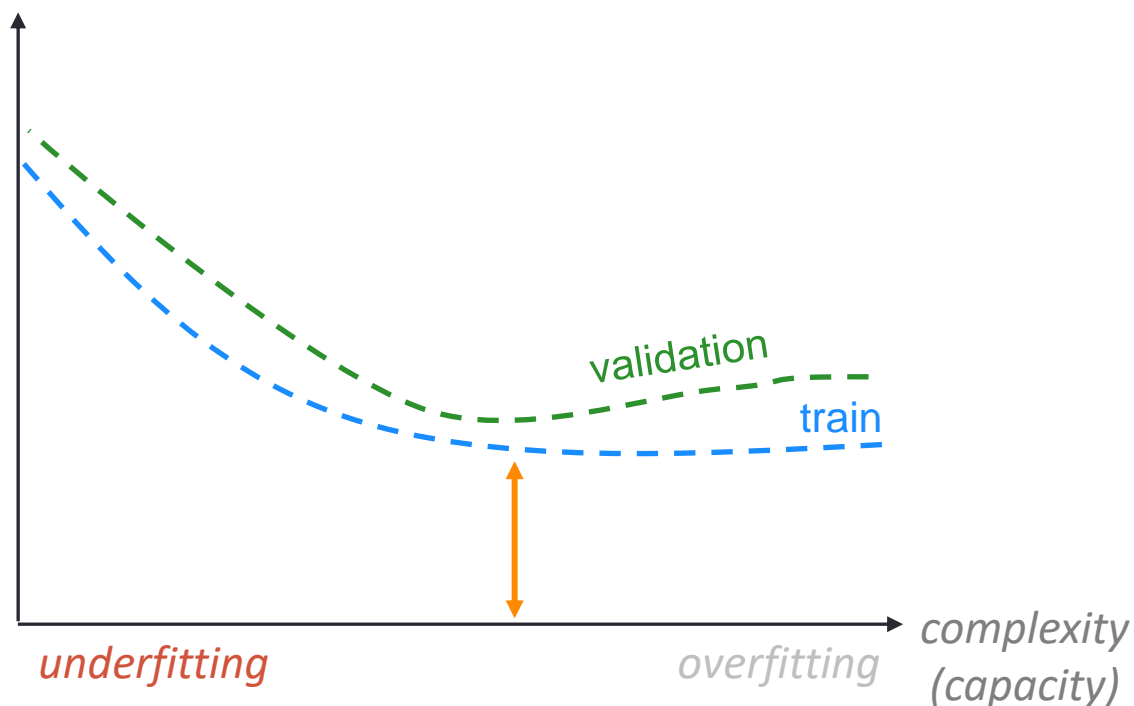
# Error decomposition using validation

- Given a training set $S$, a validation set $V$, and a hypothesis $h_S$, the generalization error can be decomposed into:

$$L_D(h_S) = \big(L_D(h_S) - L_V(h_S)\big) + \big(L_V(h_S) - L_S(h_S)\big) + L_S(h_S)$$

If this term is large, $h_S$ probably underfits.

Possible solutions:

- Increase the complexity

- Improve tuning

- Change feature mapping

- Change hypothesis class



validation

train

underfitting        overfitting        complexity (capacity)

# Summary

- The model complexity creates a tradeoff between bias and variance

- Model selection

  - Validation curves help tune hyperparameters (and model complexity)

  - Error could also be decomposed using validation

  - Further reading: Chapter 11 in Understanding ML: From Theory to Algorithms