# Introduction to Machine Learning

*SCP8084699 - LT Informatica*

ML System Design, Diagnoses and Learning Curves

Prof. Lamberto Ballan

DIPARTIMENTO **MATEMATICA**
Department of Mathematics "Tullio Levi-Civita"

# What we will learn today?

- A bit more on model selection and evaluation: advice for applying machine learning

- Diagnosing machine learning models

# Recap: Supervised Learning

- Classification (discrete) *vs* Regression (real-valued output)

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

| Size in feet² (x) | Price ($) in 1K's (y) |
|---|---|
| 2104 | 460 |
| 1416 | 232 |
| 1534 | 315 |
| 852 | 178 |
| ... | ... |

Training Set

Learning Algorithm

*Looking for a good* **Bias-Variance Tradeoff**

$x \longrightarrow \boxed{h} \longrightarrow y$

*Hypothesis Space H*

*h approximates the unknown target f*

$h \sim f \colon X {\to} Y$
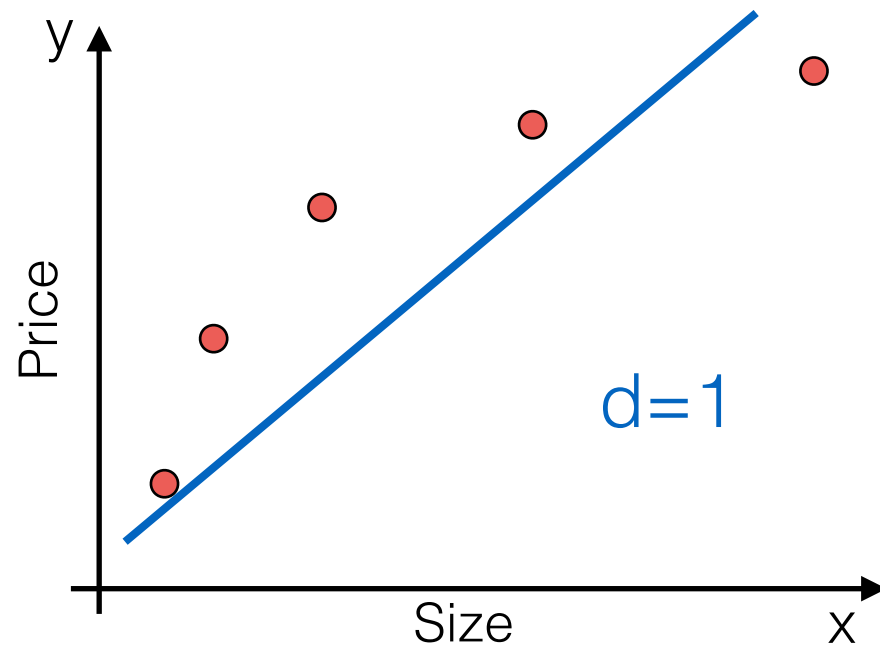
# Model Selection and Evaluation

- **Hold-out**: we keep a subset of $v$ samples from the training set (the validation set) to evaluate our model

  ‣ A classifier/regressor is trained on $m$-$v$ samples

  ‣ Parameters are optimized on the _training_-_validation_ sets: then you should evaluate performances on the _test_ set

  ‣ Size (cardinality) of training+validation sets should be greater than test set, e.g. 70%,15%,15%

- **$k$-fold cross validation**: iterate on $k$ disjoint subsets

- Given a task, pick the "right" evaluation metric
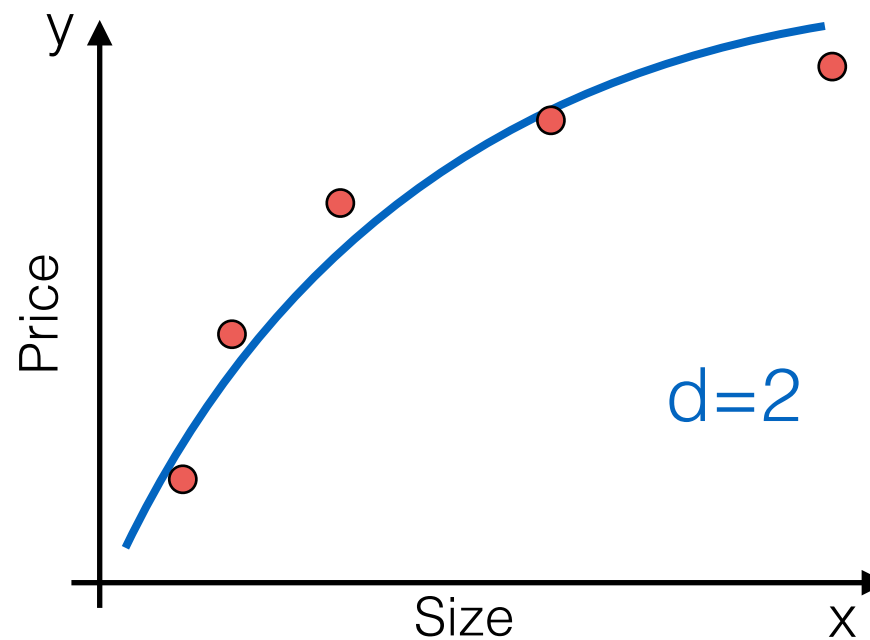
# Diagnosing bias *vs* variance

- If your learning model doesn't work as expected, almost all the time it will be because you have either a *high bias* problem or a *high variance* problem

  ‣ How to figure out what's happening (in practice)?

  ‣ What can we do to fix/alleviate the problem?

# Diagnosing bias vs variance
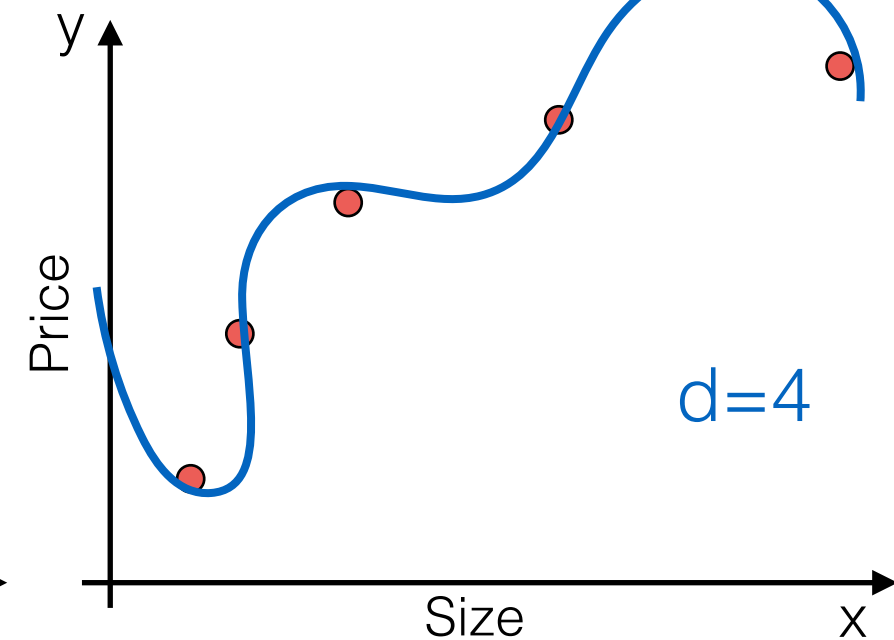
*Underfitting (high bias)*

*Overfitting (high variance)*



$$h_\theta(x) = \theta_0 + \theta_1 x$$

$$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

- We can now look again at this example taking into account hold-out and bias-variance tradeoff
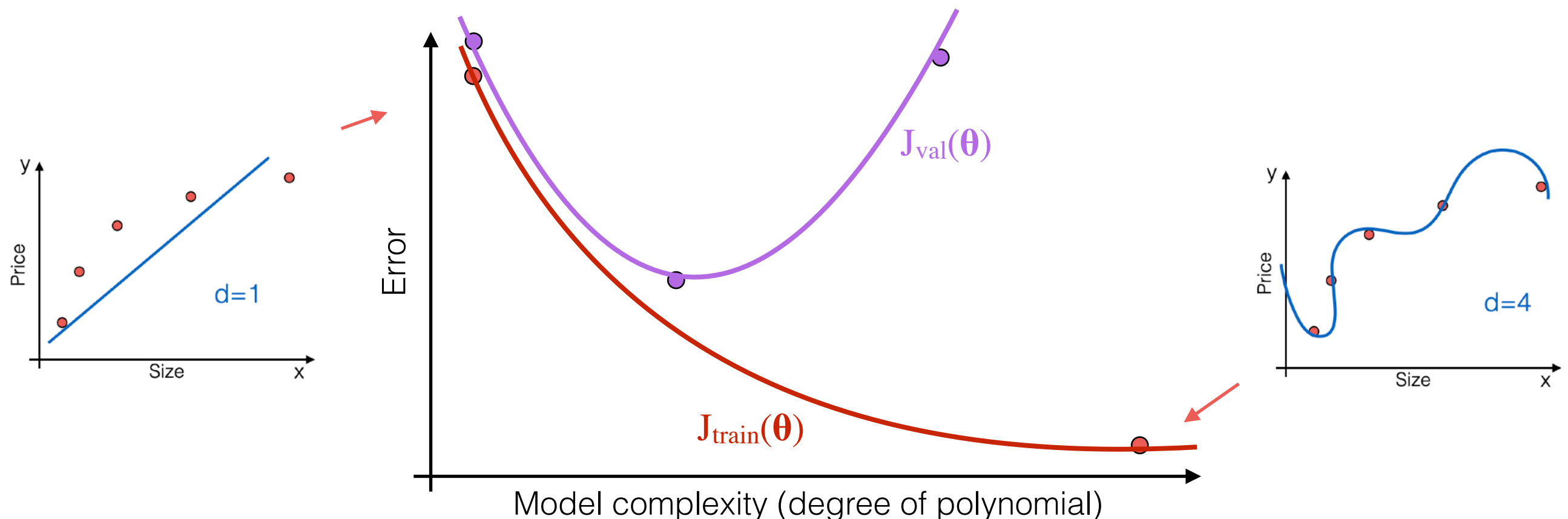
# Diagnosing bias vs variance
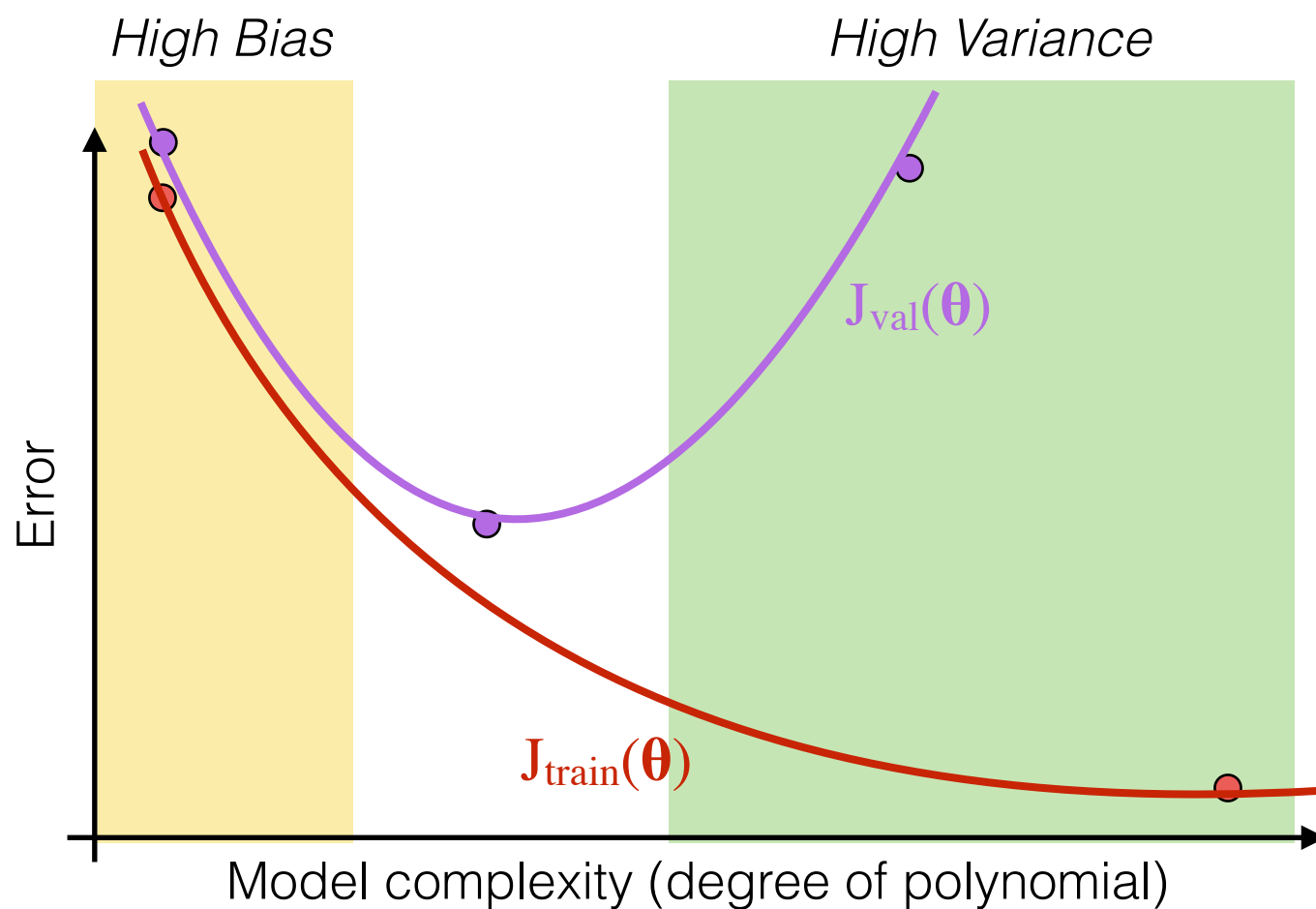
- "Measuring" bias vs variance:

  ‣ Training Error:   $J_{train}(\boldsymbol{\theta}) = \dfrac{1}{2m_t} \sum\limits_{i=1}^{m_t} (h_\theta(x^{(i)}) - y^{(i)})^2$

  ‣ Validation Error:  $J_{val}(\boldsymbol{\theta}) = \dfrac{1}{2m_v} \sum\limits_{i=1}^{m_v} (h_\theta(x^{(i)}) - y^{(i)})^2$

# Diagnosing bias vs variance

- Our learning model doesn't work as expected; is it a bias problem or a variance problem?



*High Bias*  *High Variance*

$J_{val}(\boldsymbol{\theta})$

$J_{train}(\boldsymbol{\theta})$

Error

Model complexity (degree of polynomial)

**High bias** (underfit):

$J_{train}(\boldsymbol{\theta})$ will be high

$J_{val}(\boldsymbol{\theta}) \approx J_{train}(\boldsymbol{\theta})$

**High variance** (overfit):

$J_{train}(\boldsymbol{\theta})$ will be low

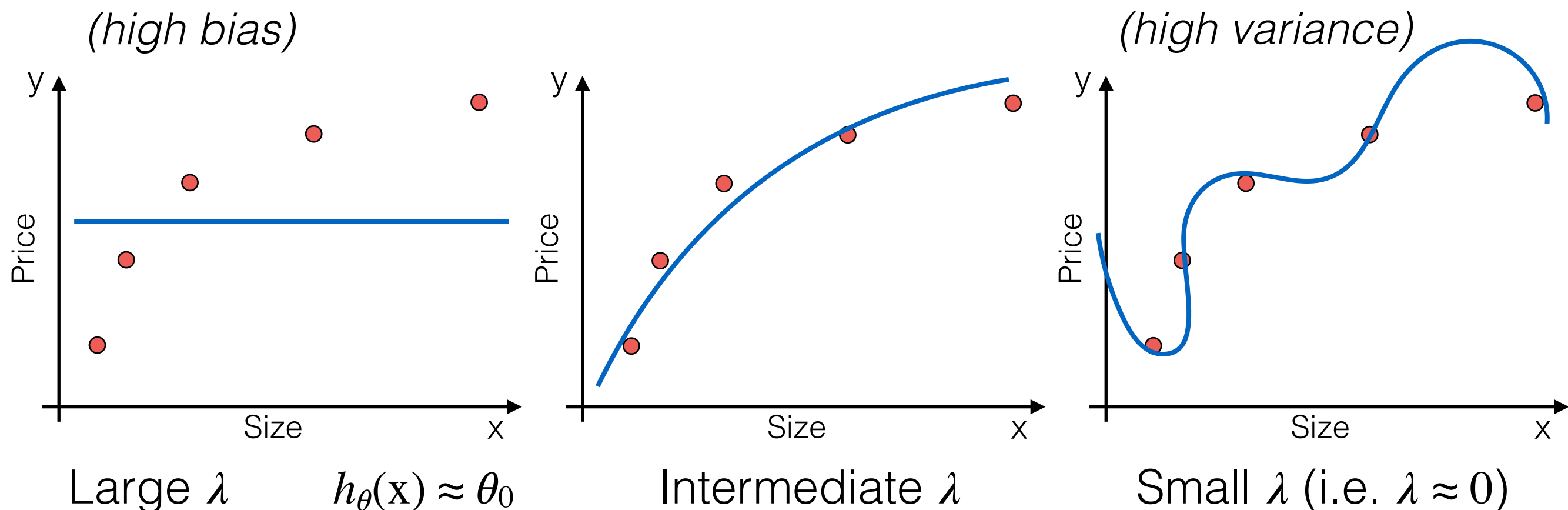$J_{val}(\boldsymbol{\theta}) \gg J_{train}(\boldsymbol{\theta})$

# Diagnosing bias vs variance

- What's the contribution of regularization?

$$\min_{\theta} \; \frac{1}{2m} \left[ \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^{n} \theta_j^2 \right]$$

$$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$



*(high bias)*

*(high variance)*

Large $\lambda$     $h_\theta(x) \approx \theta_0$       Intermediate $\lambda$       Small $\lambda$ (i.e. $\lambda \approx 0$)

# Diagnosing bias vs variance

- Choosing the regularization parameter $\lambda$:

$$\min_{\theta} \ \frac{1}{2m} \left[ \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^{n} \theta_j^2 \right]$$

$$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

- Note: our definition of $J_{train}$, $J_{val}$, $J_{test}$ don't change

  ‣ Training Error: $\quad J_{train}(\mathbf{\theta}) = \dfrac{1}{2m_t} \sum\limits_{i=1}^{m_t} (h_\theta(x^{(i)}) - y^{(i)})^2$

  ‣ Validation Error: $\ J_{val}(\mathbf{\theta}) = \dfrac{1}{2m_v} \sum\limits_{i=1}^{m_v} (h_\theta(x^{(i)}) - y^{(i)})^2$

  ‣ Test Error: $\quad J_{test}(\mathbf{\theta}) = \dfrac{1}{2m_e} \sum\limits_{i=1}^{m_e} (h_\theta(x^{(i)}) - y^{(i)})^2$

# Diagnosing bias vs variance

- Choosing the regularization parameter $\lambda$:

$$\min_{\theta} \frac{1}{2m} \left[ \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^{n} \theta_j^2 \right]$$

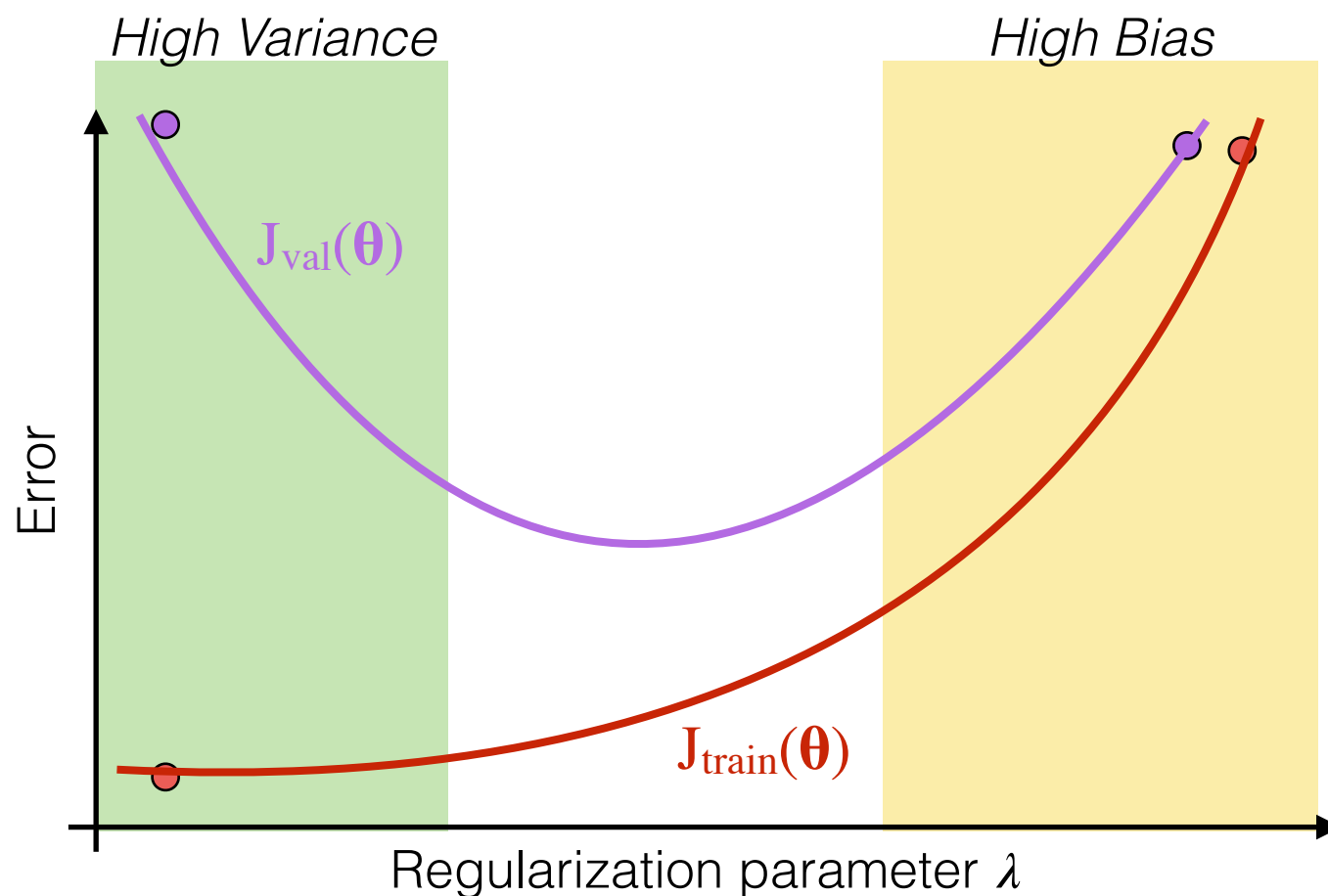Model: $h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$

*Model Selection*

1: try $\lambda = 0$ $\longrightarrow$ $\min_{\theta} J(\theta) = \theta^{(1)}$ $\longrightarrow$ $J_{val}(\theta^{(1)})$

2: try $\lambda = 0.01$ $\longrightarrow$ $\theta^{(2)}$ $\longrightarrow$ $J_{val}(\theta^{(2)})$

3: try $\lambda = 0.02$ $\longrightarrow$ $\theta^{(3)}$ $\longrightarrow$ $\boxed{J_{val}(\theta^{(3)})}$ *(lowest)*

4: try $\lambda = 0.04$ $\longrightarrow$ $\theta^{(4)}$ $\longrightarrow$ $J_{val}(\theta^{(4)})$

5: try $\lambda = 0.08$ $\longrightarrow$ $\theta^{(5)}$ $\longrightarrow$ $J_{val}(\theta^{(5)})$

$\vdots$

12: try $\lambda \approx 10$ $\longrightarrow$ $\theta^{(12)}$ $\longrightarrow$ $J_{val}(\theta^{(12)})$

# Diagnosing bias vs variance

- Bias/Variance as a function of the parameter $\lambda$:

$$\min_{\boldsymbol{\theta}} \ \frac{1}{2m} \left[ \sum_{i=1}^{m} (h_{\theta}(\mathbf{x}^{(i)}) - \mathbf{y}^{(i)})^2 + \lambda \sum_{j=1}^{n} \theta_j^2 \right]$$

*High Variance*          *High Bias*

$J_{val}(\boldsymbol{\theta})$

Error

$J_{train}(\boldsymbol{\theta})$

Regularization parameter $\lambda$

$$J_{train}(\boldsymbol{\theta}) = \frac{1}{2m_t} \sum_{i=1}^{m_t} (h_{\theta}(\mathbf{x}^{(i)}) - \mathbf{y}^{(i)})^2$$

$$J_{val}(\boldsymbol{\theta}) = \frac{1}{2m_v} \sum_{i=1}^{m_v} (h_{\theta}(\mathbf{x}^{(i)}) - \mathbf{y}^{(i)})^2$$
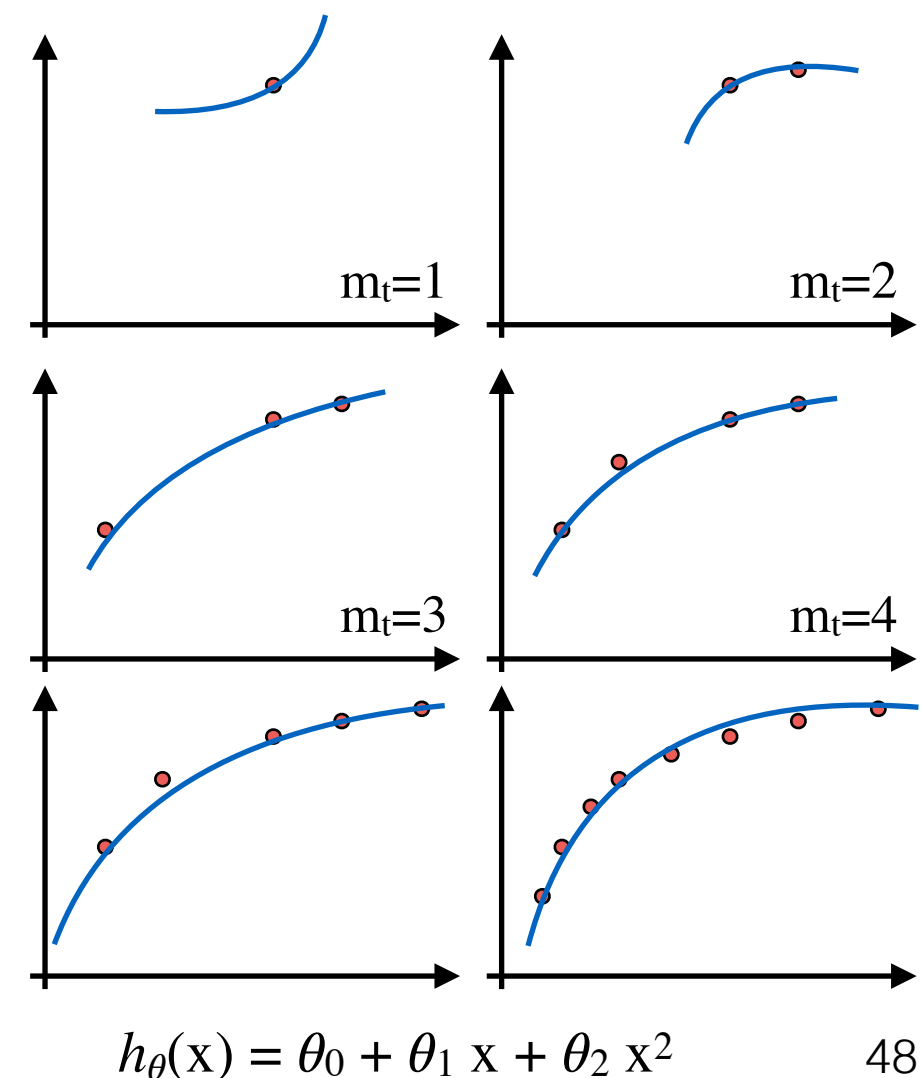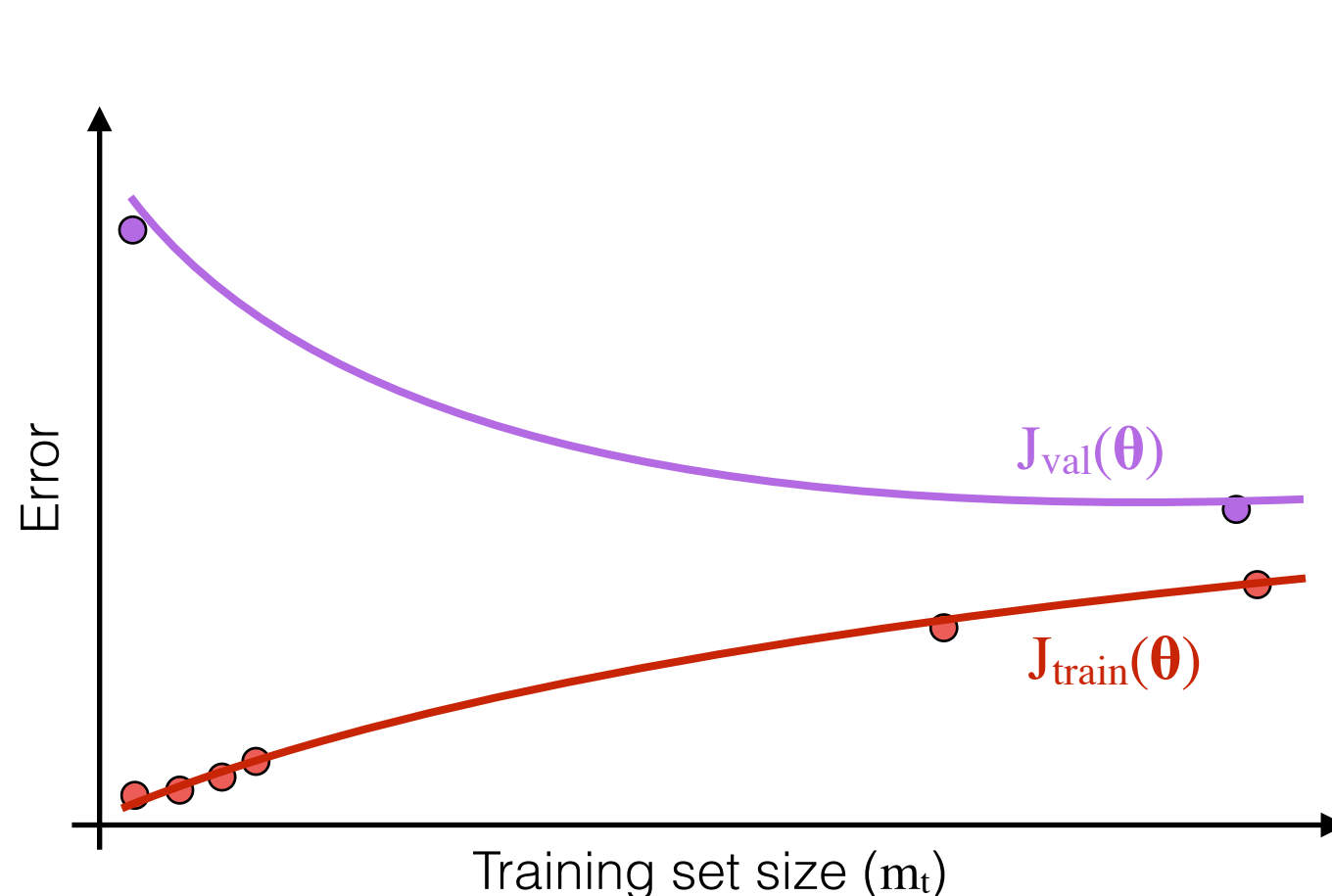
# Diagnosing bias vs variance

- By now you have seen bias and variance from a lot of different perspectives

- Let's now take all the insights we have gone through in order to build a "diagnostic tool" for ML systems

# Learning Curves

- Learning curves can be used to diagnose if a model may be suffering from bias, variance or a bit of both

$$J_{train}(\boldsymbol{\theta}) = \frac{1}{2m_t} \sum_{i=1}^{m_t} (h_\theta(x^{(i)}) - y^{(i)})^2 \qquad J_{val}(\boldsymbol{\theta}) = \frac{1}{2m_v} \sum_{i=1}^{m_v} (h_\theta(x^{(i)}) - y^{(i)})^2$$



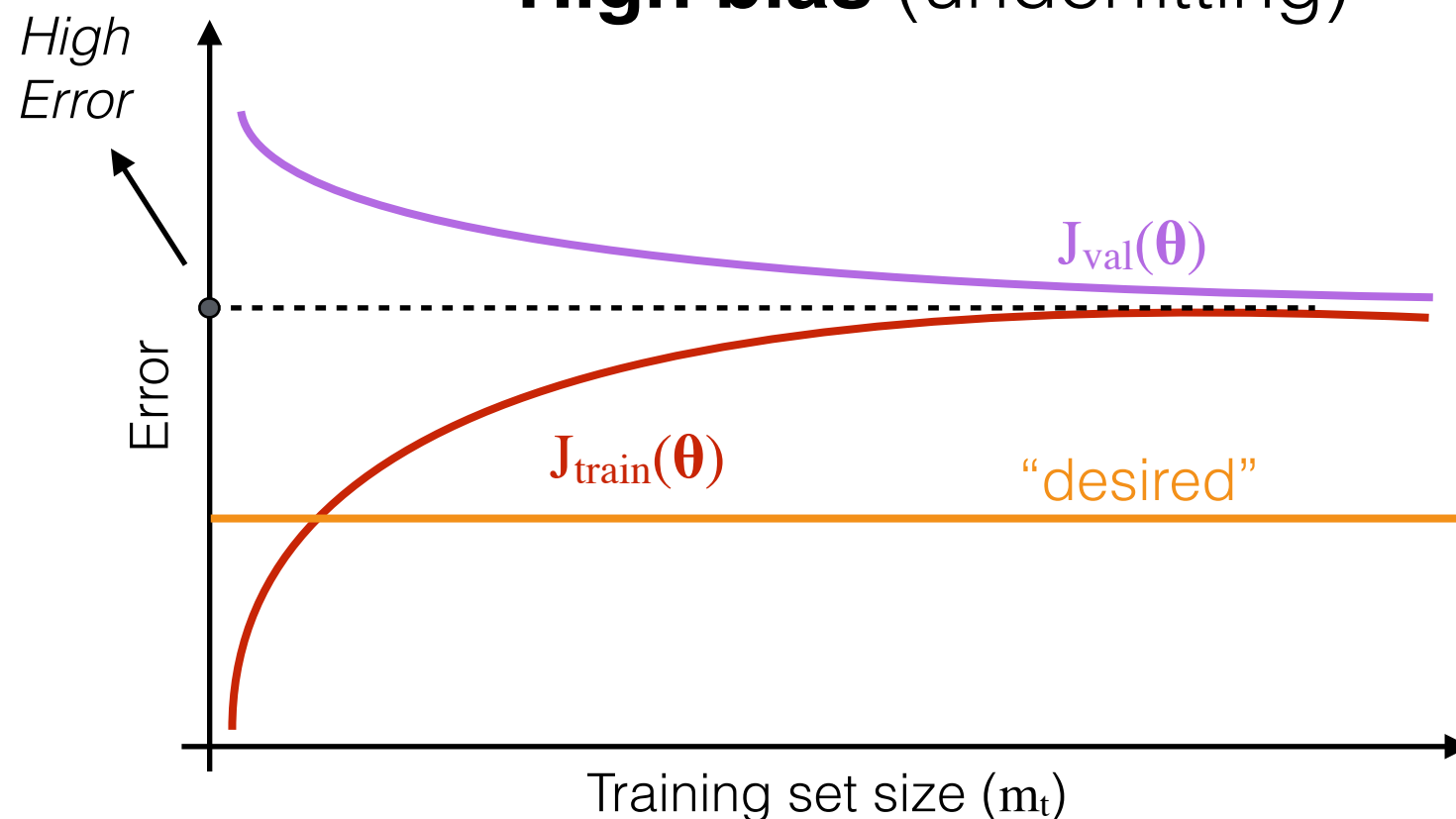$$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$
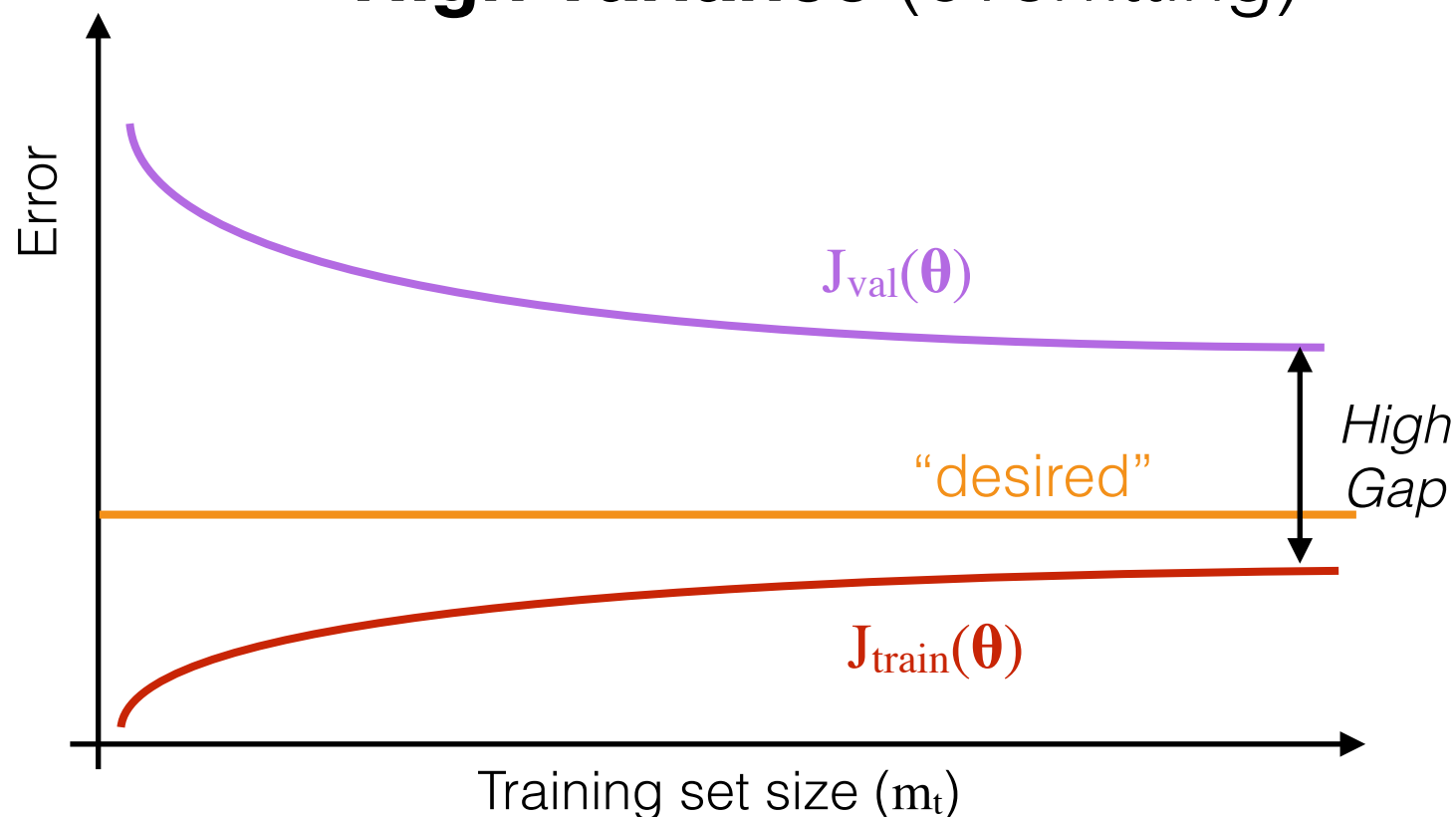
# Learning Curves

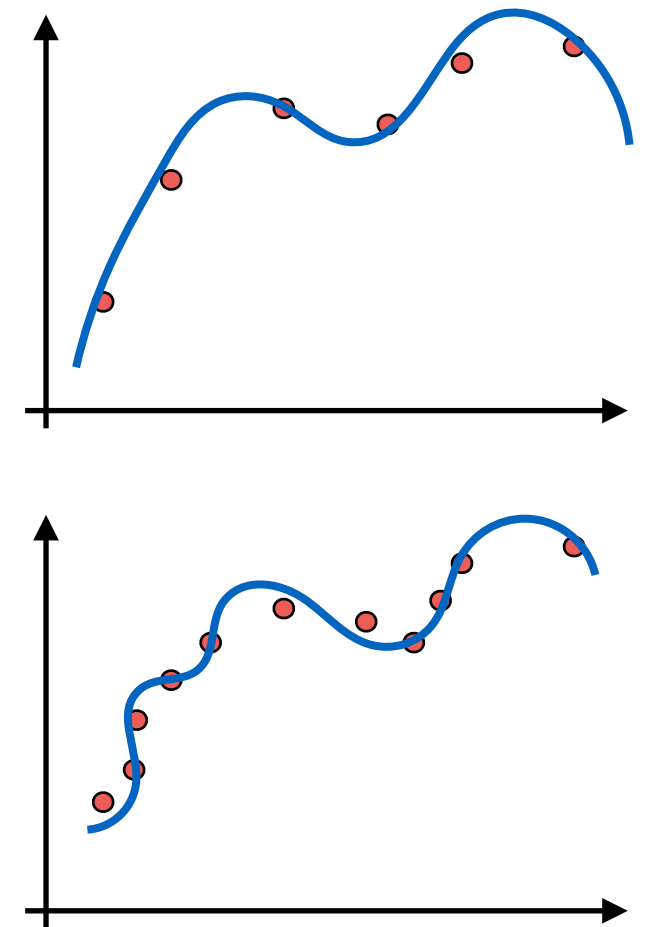- That's the general intuition… but what's about bias and variance problems?

**High bias** (underfitting)

$$h_\theta(\mathrm{x}) = \theta_0 + \theta_1\,\mathrm{x}$$

*High Error*

Error

$J_{val}(\boldsymbol{\theta})$

$J_{train}(\boldsymbol{\theta})$

"desired"

Training set size ($m_t$)

*Note: in case of high bias, getting more training data will not help much*

# Learning Curves

- That's the general intuition… but what's about bias and variance problems?

**High variance** (overfitting)



$h_\theta(\mathrm{x}) = \theta_0 + \theta_1\,\mathrm{x} + \ldots + \theta_{50}\,\mathrm{x}^{50}$

(*and small* $\lambda$)

$J_{val}(\mathbf{\theta})$

"desired"

$J_{train}(\mathbf{\theta})$

*High Gap*

Error

Training set size ($\mathrm{m_t}$)
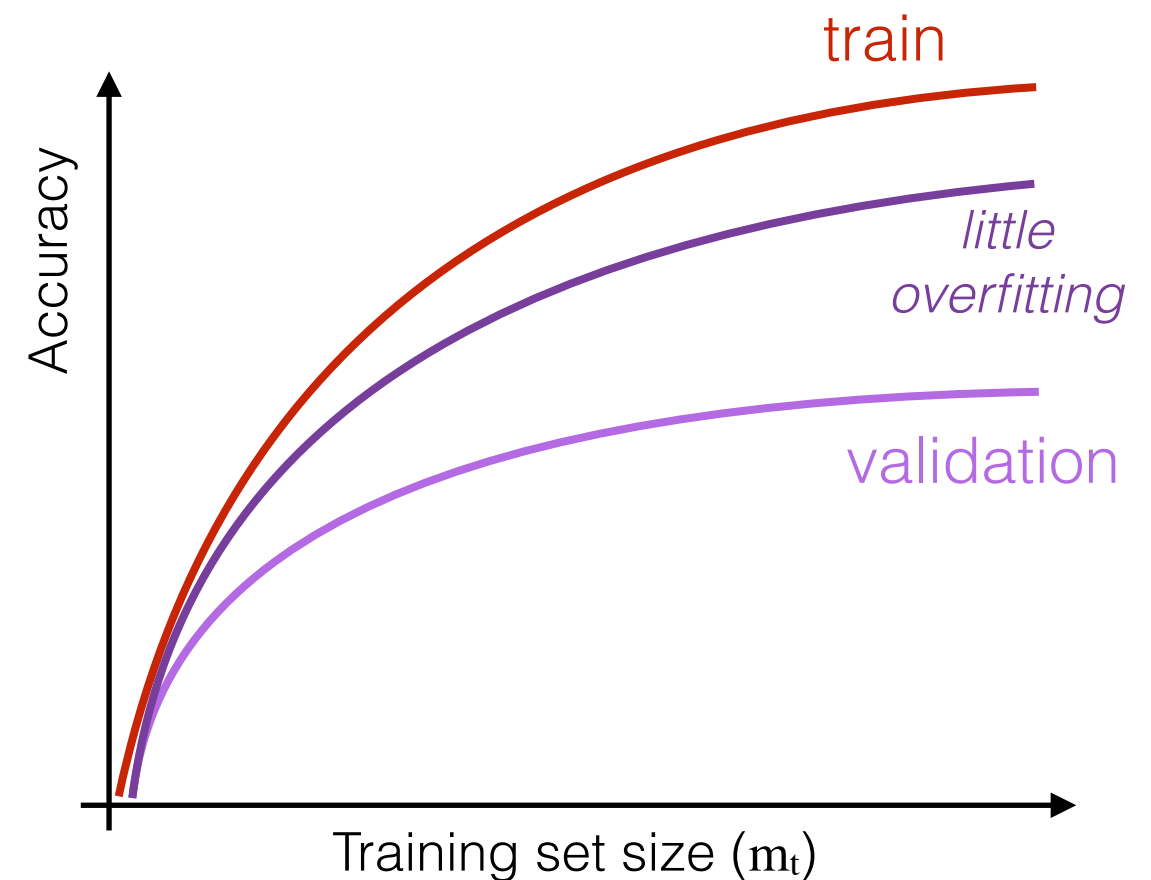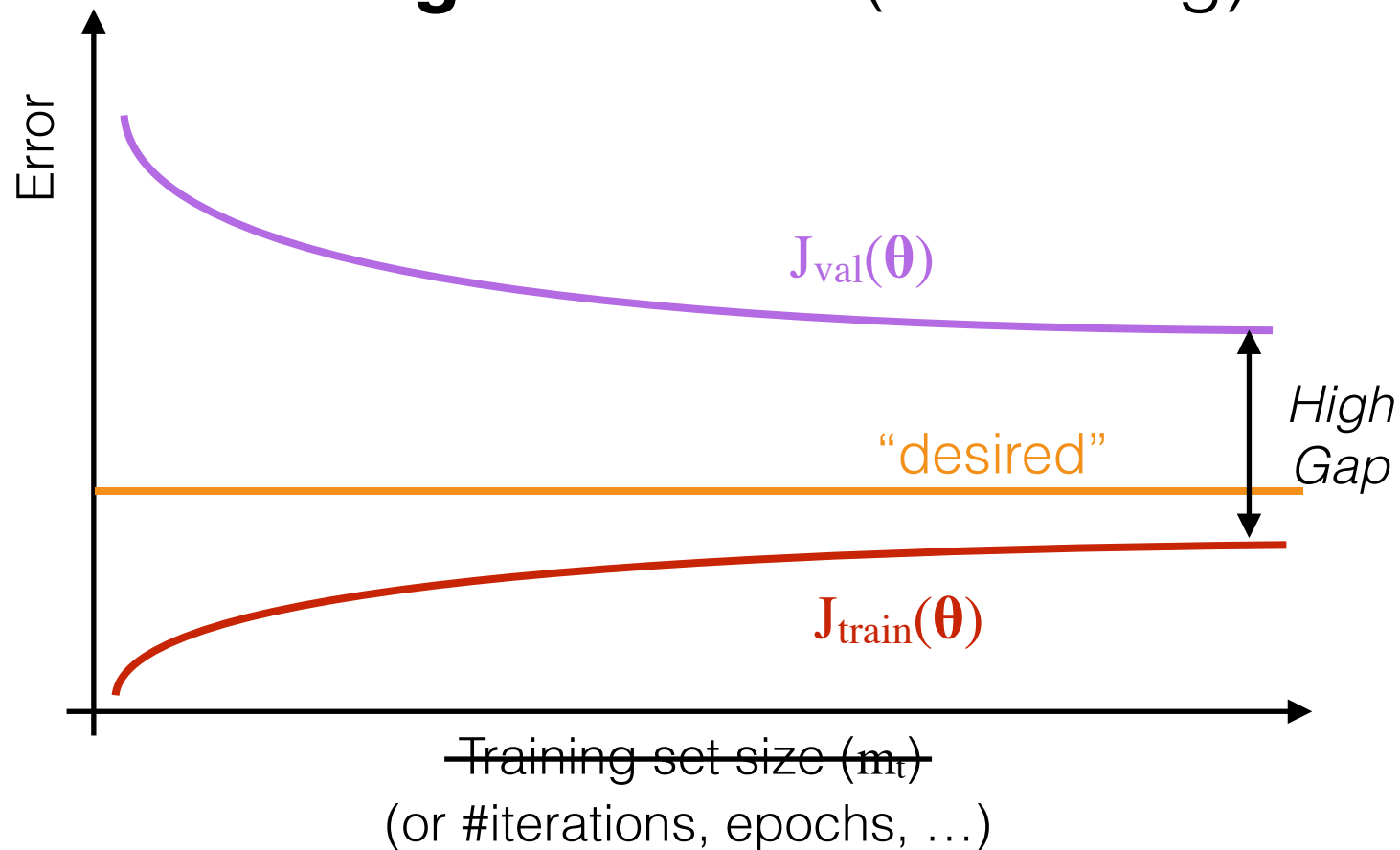
*Note: in case of high variance, getting more training data is likely to help*

# Learning Curves

- You can compute learning curves w.r.t. different "dimensions" (e.g. evaluation measures, no. samples)

**High variance** (overfitting)



$J_{val}(\boldsymbol{\theta})$

"desired"

*High Gap*

$J_{train}(\boldsymbol{\theta})$

~~Training set size ($m_t$)~~
(or #iterations, epochs, …)

train

*little overfitting*

validation

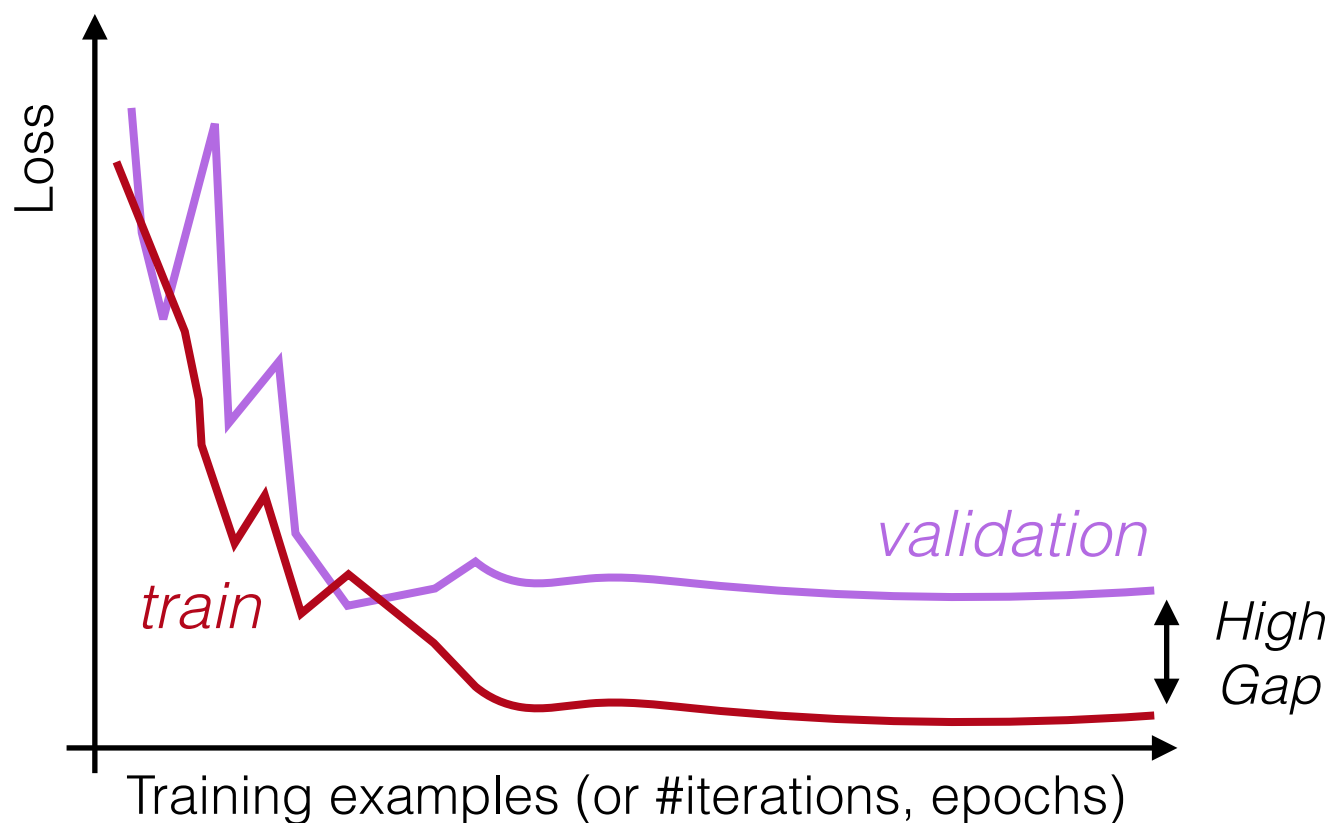Accuracy

Training set size ($m_t$)

# What to do next

- Debugging (and babysitting) a learning algorithm:

  ‣ Suppose you have implemented a regularized linear regression model for predicting housing prices

  ‣ It doesn't work on new data; what should you do next?

  ‣ You can get more training data → *Fixes high variance*

  ‣ Try smaller set of features → *Fixes high variance*

  ‣ Try getting more features → *Fixes high bias*

  ‣ Try adding complexity to the model (e.g. polynomial features) → *Fixes high bias*

  ‣ Try decreasing $\lambda$ → *Fixes high bias*

  ‣ Try increasing $\lambda$ → *Fixes high variance*

# Diagnosing our datasets

- Learning curves can be also used to diagnose the quality of our training/validation sets
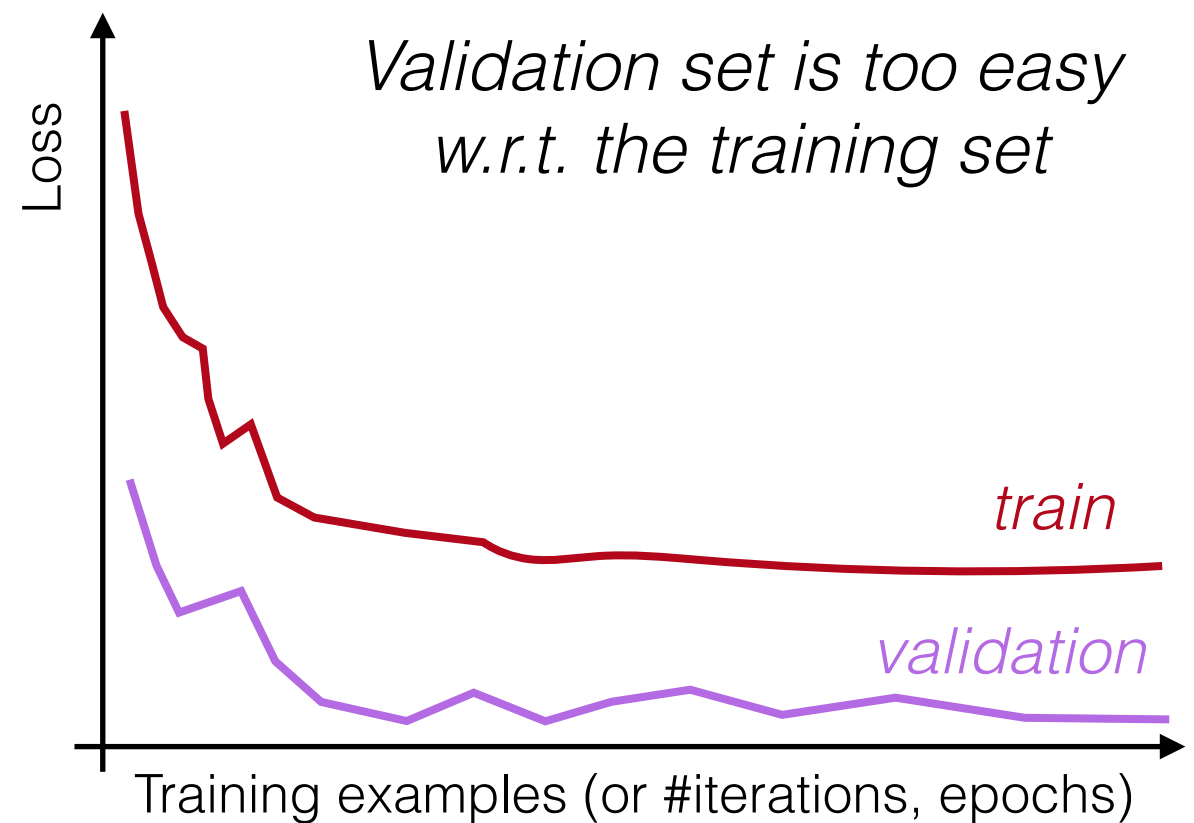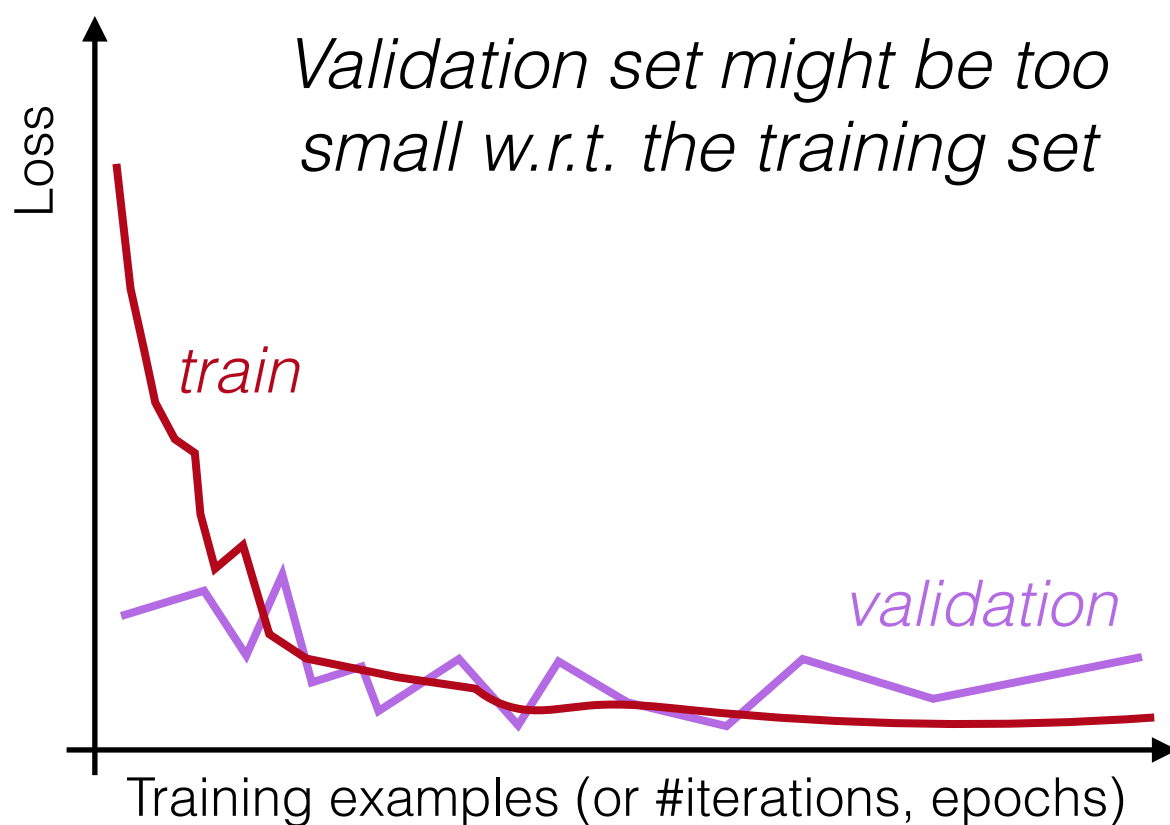
**Unrepresentative Training Set**



- ‣ The training set does not provide sufficient information to learn the problem

- ‣ It may occur if the training set has too few examples as compared to the validation set

# Diagnosing our datasets

- Learning curves can be also used to diagnose the quality of our training/validation sets

**Unrepresentative Validation Set**



*Validation set might be too small w.r.t. the training set*

*Validation set is too easy w.r.t. the training set*

# Contact

- **Office:** Torre Archimede 3CD, room 320

- **Office hours** (ricevimento)**:** Monday 11:00-13:00

✉ lamberto.ballan@unipd.it

🏠 http://www.lambertoballan.net

🏠 http://vimp.math.unipd.it

@ twitter.com/lambertoballan