# Machine Learning I
## 80-629

# Apprentissage Automatique I
## 80-629

Machine Learning fundamentals
— Week #2

# Today: what's a machine learning problem

- Core concepts

  - Modeling and parameters

  - Bias/Variance

  - Overfitting

  - Representing uncertainty

- Types of learning problems

  - Supervised learning (generative and discriminative models)

  - Unsupervised learning

  - Reinforcement learning

# Capsules

1. Machine Learning problem

2. Types of Learning problems

3. A first Supervised Model

4. Model Evaluation

5. Regularization

6. Bias/Variance

- I will follow the exposition of Chapter 5 in "Deep Learning".

  - "Operational" approach vs. a decision-theoretic/ probabilistic approach

# The components of a learning problem

- I will follow the exposition of Chapter 5 in "Deep Learning"

  - "Operational" approach vs. a decision-theoretic/ probabilistic approach

# Three main components

- Task (T)

- Performance measure (P)

- Experience (E)

"A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P, if its performance at tasks in T, measured by P, improves with experience E."

–Tom Mitchell (1997)

# Task (T)

- The end goal(s). The question you are answering.

- For example:

  - Self-driving

  - Differentiate cats from dogs

  - Recommend movies of interest to users
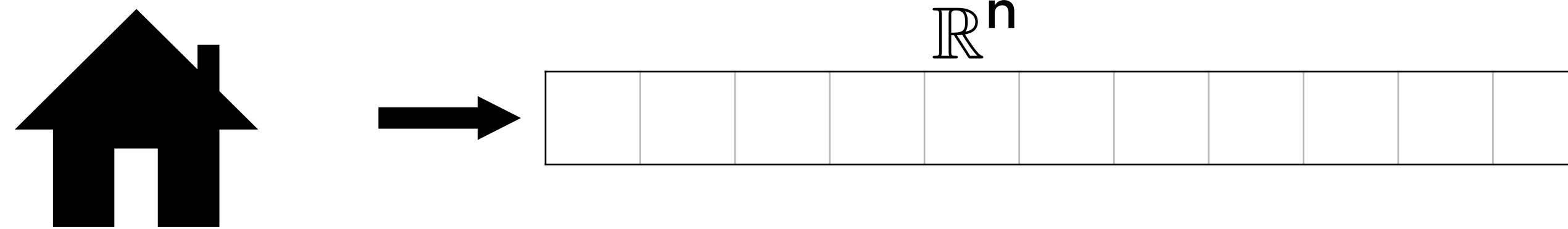
  - Select a good portfolio of stocks

# Task (T)

# Task (T)

- **Determine the price of houses?**
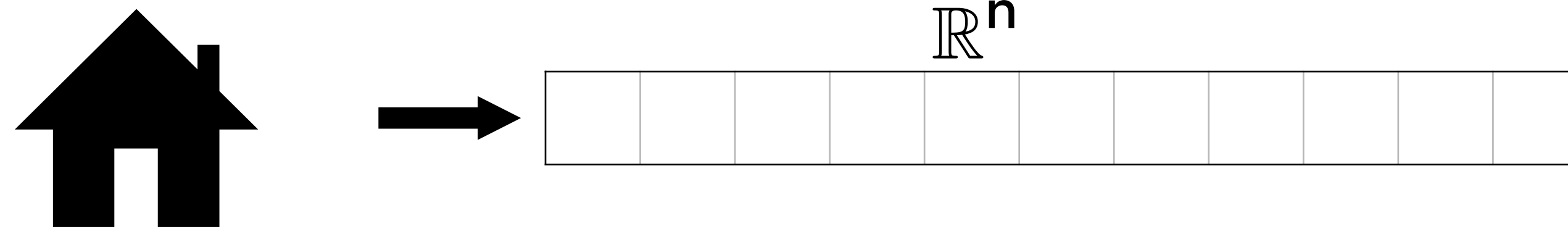
# Task (T)

- Determine the price of houses?

  - Encode houses into a set of features

    - area, number of rooms (bedrooms, bathrooms), municipal evaluation, neighborhood, etc.

$$\mathbb{R}^n$$

# Task (T)

- Determine the price of houses?

  - Encode houses into a set of features

    - area, number of rooms (bedrooms, bathrooms), municipal evaluation, neighborhood, etc.

$$\mathbb{R}^n$$

- Function from feature to house price

$$f : \boxed{\phantom{xxxxxx}} \rightarrow \text{price}$$

$$f : \mathbb{R}^n \rightarrow \mathbb{R}^+$$

# Example tasks

- **Regression: Assign a real value to an example**

$$f : \mathbb{R}^n \to \mathbb{R}$$

- **Classification: Classify instances in one of k classes**

$$f : \mathbb{R}^n \to \{1, \ldots, k\}$$

- **Clustering: Assign each instance to a cluster**

$$f : \mathbb{R}^n \to \{1, \ldots, k\}$$

# More examples

- **Transcription (e.g., document classification)**

$$f : \mathbb{R}^{n \times m} \to \mathbb{R}^k$$

- **Multi-label classification (e.g., tag prediction)**

$$f : \mathbb{R}^n \to \{0, 1\}^m$$

- **Translation (e.g., sentence from French to English)**

$$f : \mathbb{R}^n \to \mathbb{R}^m$$

# Model

- functions f are examples of models

  - Model is a simpler representation of the world

  - Has parameters (w)

    Model #1: $\qquad f(x; w) = w_1 x_1 + w_2 x_2 + \ldots + w_n x_n$

$f : \mathbb{R}^n \to \mathbb{R}$

    Model #2: $\qquad f(x; w) = w_1 x_1 + w_2 x_1^2 + w_3 x_2 + \ldots + w_{n+1} x_n$

# The performance measure (P)

- How to quantify the "goodness" of f? How to compare models?

# The performance measure (P)

- How to quantify the "goodness" of f? How to compare models?

  - Specific to the task at hand

    - E.g., linear regression: squared error

# The performance measure (P)

- How to quantify the "goodness" of f? How to compare models?

  - Specific to the task at hand

    - E.g., linear regression: squared error

$$\frac{1}{P} := \sum_i \left( f(x_i) - y_i \right)^2$$

$f_1$ is better than $f_2$ if $P_1 > P_2$

# The performance measure (P)

- How to quantify the "goodness" of f? How to compare models?

  - Specific to the task at hand

    - E.g., linear regression: squared error

    $$\frac{1}{P} := \sum_i \left( f(x_i) - y_i \right)^2$$

    $f_1$ is better than $f_2$ if $P_1 > P_2$

  - Encodes knowledge of what's important

    - A model with perfect performance behaves perfectly

- Examples: accuracy, error rate, log-probability, F score.

# Experience (E)

- What data does f experience?

  - (Focus on algorithms that experience whole datasets)

  - Unsupervised. Examples alone.

  $$\{\mathbf{x}_i\}_{i=0}^n$$

  - Supervised. Examples come with labels.

  $$\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=0}^n$$

# Different types of experience

1. Unsupervised Learning

2. Supervised Learning

3. Reinforcement Learning

# Experience (E)

- **What data does f experience?**

  - **(Focus on algorithms that experience whole datasets)**

  - **Unsupervised. Examples alone.**

$$\{\mathbf{x}_i\}_{i=0}^{n}$$

  - **Supervised. Examples come with labels.**

$$\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=0}^{n}$$

# 1. Unsupervised

# 1. Unsupervised

- **Experience examples alone**

$$\{x_i\}_{i=0}^n$$

# 1. Unsupervised

- **Experience examples alone**

$$\{\mathbf{x}_i\}_{i=0}^{n}$$

- **Learn "useful properties of the structure of the data"**

# 1. Unsupervised

- Experience examples alone

$$\{x_i\}_{i=0}^n$$
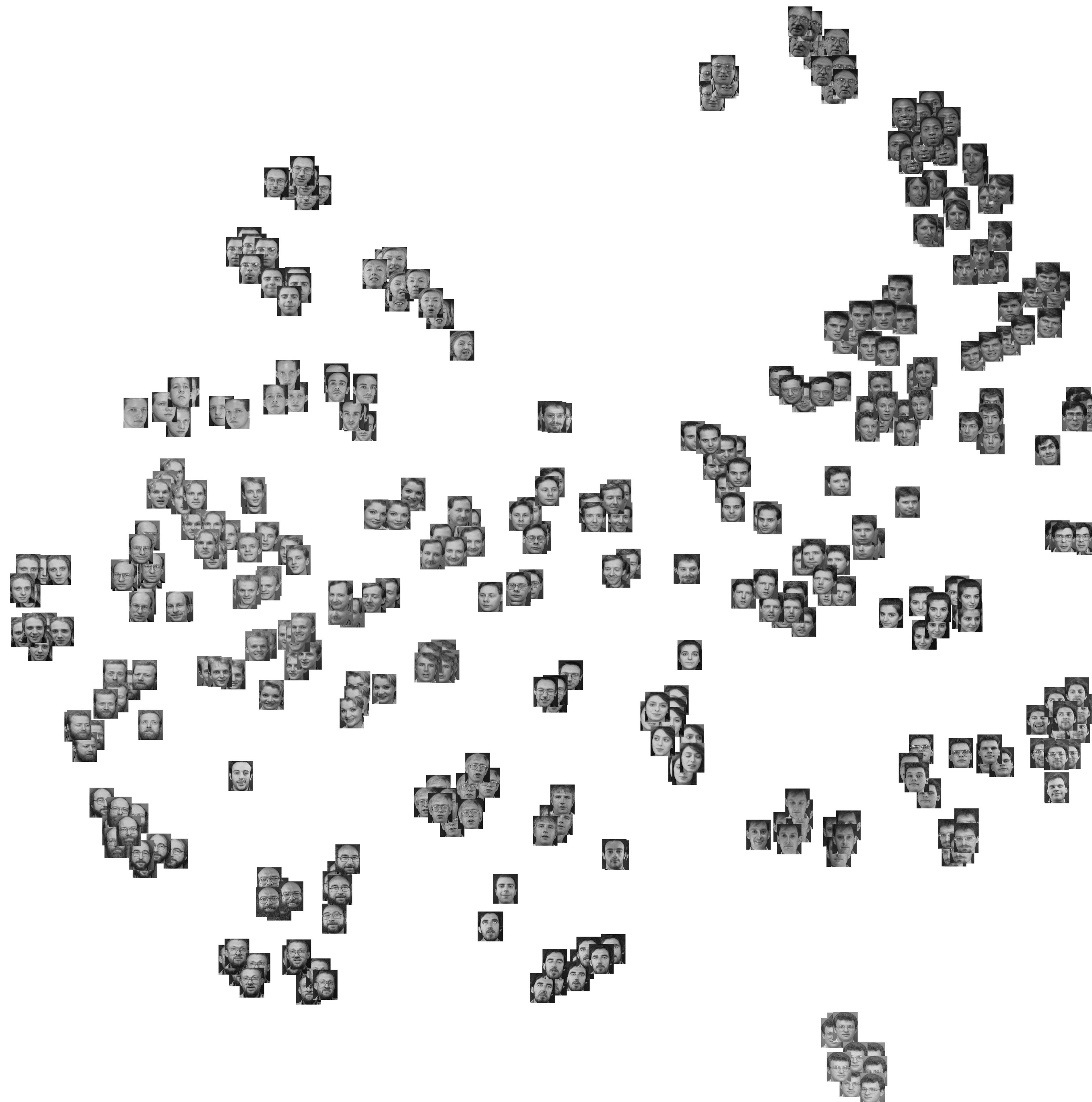
- Learn "useful properties of the structure of the data"

  - E.g., Clustering



https://home.deib.polimi.it/matteucc/Clustering/tutorial_html/

# 1. Unsupervised

- Experience examples alone

$$\{x_i\}_{i=0}^n$$

- Learn "useful properties of the structure of the data"

  - E.g., Clustering

- Probabilistic models

  - Density modeling p(x), PCA, FA.

https://home.deib.polimi.it/matteucc/Clustering/tutorial_html/

Example from a non-linear
dimensionality reduction
technique (tsne)
https://lvdmaaten.github.io/tsne/

Laurent Charlin — 80-629

# 2. Supervised

- **Experience examples and their label(s)**

$$\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=0}^{n}$$

- **Given an example (x) predict its label (y)**

$$\mathbf{f} : \mathbf{X} \rightarrow \mathbf{Y} \qquad \mathbf{p(y|x)}$$

- **E.g., regression, classification**



[https://en.wikipedia.org/wiki/Regression_analysis]



[https://jaxenter.com/machine-learning-an-introduction-for-programmers-122135.html]

# Distinction can be blurry

- Supervised data modeled jointly:

$$(\mathbf{x}, \mathbf{y}), \quad \mathbf{p}(\mathbf{y} \mid \mathbf{x}) = \frac{\mathbf{p}(\mathbf{x}, \mathbf{y})}{\sum_{\mathbf{y}'} \mathbf{P}(\mathbf{x}, \mathbf{y}')}$$

- Unsupervised data modeled as supervised data:

$$\mathbf{x} \in \mathbb{R}^{\mathbf{n}}, \quad \mathbf{p}(\mathbf{x}) = \prod_{i=1}^{n} \mathbf{p}(\mathbf{x_i} \mid \mathbf{x_1}, \ldots, \mathbf{x_{i-1}}) \mathbf{P}(\mathbf{x_0})$$

> Conditional model: $\mathbf{P}(\mathbf{y} \mid \mathbf{x})$
>
> Generative model: $\mathbf{P}(\mathbf{y}, \mathbf{x})$

# Semi-supervised learning

- **Idea: Can we augment a supervised dataset with unsupervised data**

$$(\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=0}^n, \{\mathbf{x}_j\}_{j=0}^m)$$

  - **Unlabelled data are cheap (images on the web). Labeled data are expensive**

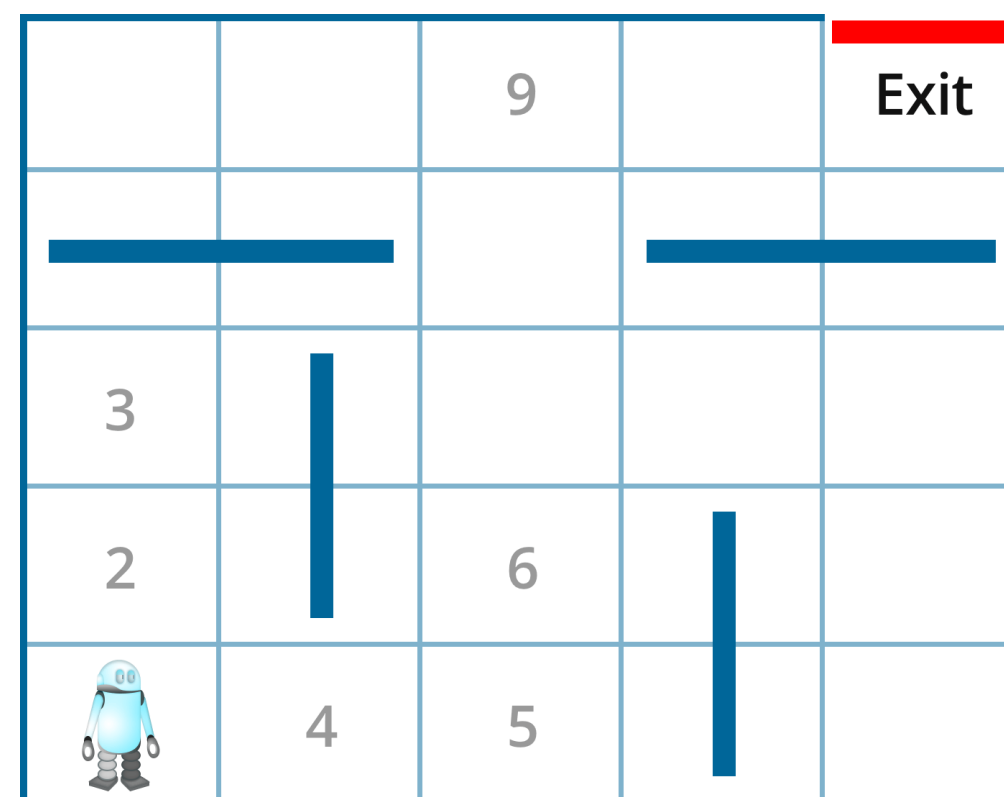  - **Can the unlabelled data help model x and yield a better model for y?**

23

[Dog]                    [ ]

# 3. Reinforcement learning

- **The algorithm interacts with the environment**

  - **The algorithm observes its environment. Prototypical example is a robot navigating a maze**

  - **The resulting dataset depends on the algorithm's choices**

# A first supervised example

# Dataset(s)

- Each instance:  $x_i$

- A dataset is a set of instances:  $\{x_i\}$

- Often denoted as a matrix $X$, (design matrix)

|  | Nb.bed. | Area | Neigh. | . | . |  | Price |
|---|---|---|---|---|---|---|---|
| $x_0$ | 1 | 0 | 0 | 0 | 0 | $y_0$ | 125000 |
| $x_1$ | 1 | 100 | 1 | .2 | .5 | $y_1$ | 150000 |
| $x_2$ | 3 | 200 | 0 | .1 | .2 | $y_2$ | 350000 |
| $x_3$ | 1 | 150 | 1 | .4 | .1 | $y_3$ | 275000 |
| $x_4$ | 2 | 210 | 2 | .5 | 1.1 | $y_4$ | 225000 |

  - rows are instances

  - columns are features

  - (assumes that all instances can be encoded using a fixed-size vector)

# Concrete full example: Linear regression

$$y_i = w_0 x_{i0} + w_1 x_{i1} + w_2 x_{i2} + \cdots + w_p x_{ip} \qquad w \in \mathbb{R}^p$$

$$= \sum_{j=0}^{p} w_j x_{ij} = w^\top x_i$$

- **w is a vector of parameters (weights)**

- **$w_j$ represents the effect of feature j on y**

- **Task: predict y from x using $w^\top x_i$**

# Concrete full example: Linear regression

$$y = w^\top x$$

- **Task: predict** $y_i$ **from** $x_i$ **using** $w^\top x_i$

- **Performance: mean squared error**

$$\text{MSE} := \frac{1}{n} \sum_i (y_i - w^\top x_i)^2$$

$$= \frac{1}{n} \sum_i (y_i - \hat{y}_i)^2$$

- **Experience:** $\{(x, y)\}$

# How do you find optimal parameters?

- Optimize the MSE with respect to the parameters of the model

$$L(\mathbf{w}) = \frac{1}{n} \sum_i (\mathbf{y}_i - \mathbf{w}^\top \mathbf{x}_i)^2$$

- Take the gradient of the MSE. Set equation to 0 and solve.

## Take gradient of the Loss wrt its parameters

$$\nabla_{\mathbf{w}} L(\mathbf{w}) = \nabla_{\mathbf{w}} \frac{1}{n} \sum_{i} (y_i - \mathbf{w}^{\top} x_i)^2$$

$$= \frac{1}{n} \sum_{i} \nabla_{\mathbf{w}} (y_i - \mathbf{w}^{\top} x_i)^2$$

$$= \frac{1}{n} \sum_{i} 2(y_i - \mathbf{w}^{\top} x_i)(-x_i)$$

## Set gradient to 0.

$$0 = \frac{1}{n} \sum_{i} (-2x_i y_i + 2(\mathbf{w}^{\top} x_i) x_i)$$

$$\Rightarrow \sum_{i} 2(\mathbf{w}^{\top} x_i) x_i = \sum_{i} 2 x_i y_i$$

$$\Rightarrow \mathbf{w}(X^{\top} X) = X^{\top} Y$$

$$\Rightarrow \mathbf{w} = (X^{\top} X)^{-1} X^{\top} Y$$

# $\mathbf{y = w_1 x_1}$



Linear regression example

Optimization of $\boldsymbol{w}$

[Figure 5.1, Chapter 5, Deep Learning]

# Model evaluation

# The goal of ML

- **Predict on new inputs ($X^{new}$)**

  - **Given the price of houses this year (X, Y), I want to predict the price of houses next year ( $X^{new}$ )**

| | Nb.bed. | Area | Neigh. | · | · | | Price |
|---|---|---|---|---|---|---|---|
| $x_0$ | 1 | 0 | 0 | 0 | 0 | $y_0$ | 125000 |
| $x_1$ | 1 | 100 | 1 | .2 | .5 | $y_1$ | 150000 |
| $x_2$ | 3 | 200 | 0 | .1 | .2 | $y_2$ | 350000 |
| $x_3$ | 1 | 150 | 1 | .4 | .1 | $y_3$ | 275000 |
| $x_4$ | 2 | 210 | 2 | .5 | 1.1 | $y_4$ | 225000 |

X — Y

| | Nb.bed. | Area | Neigh. | · | · | | Price |
|---|---|---|---|---|---|---|---|
| $x_0$ | 1 | 0 | 0 | 0 | 0 | $y_0$ | ? |
| $x_1$ | 2 | 50 | 1 | .3 | .8 | $y_1$ | ? |
| $x_2$ | 1 | 100 | 1 | .5 | 1.4 | $y_2$ | ? |
| $x_3$ | 4 | 170 | 0 | .7 | .4 | $y_3$ | ? |
| $x_4$ | 1 | 120 | 3 | .9 | .5 | $y_4$ | ? |

$X^{new}$ — $Y^{new}$

- **We can use our estimated ($\widehat{w}$): $Y^{new} = \widehat{w}^\top X^{new}$**

# Generalization

Price

- Loss$^{(X,Y)}$ : The one you can evaluate

- Loss$^{(X^{new},Y^{new})}$ : The one that you care about

$y_0$, $y_1$, $y_2$, $y_3$, $y_4$ with values ?, ?, ?, ?, ? — $Y^{new}$

- In general minimizing the former will not yield the best loss on the latter:

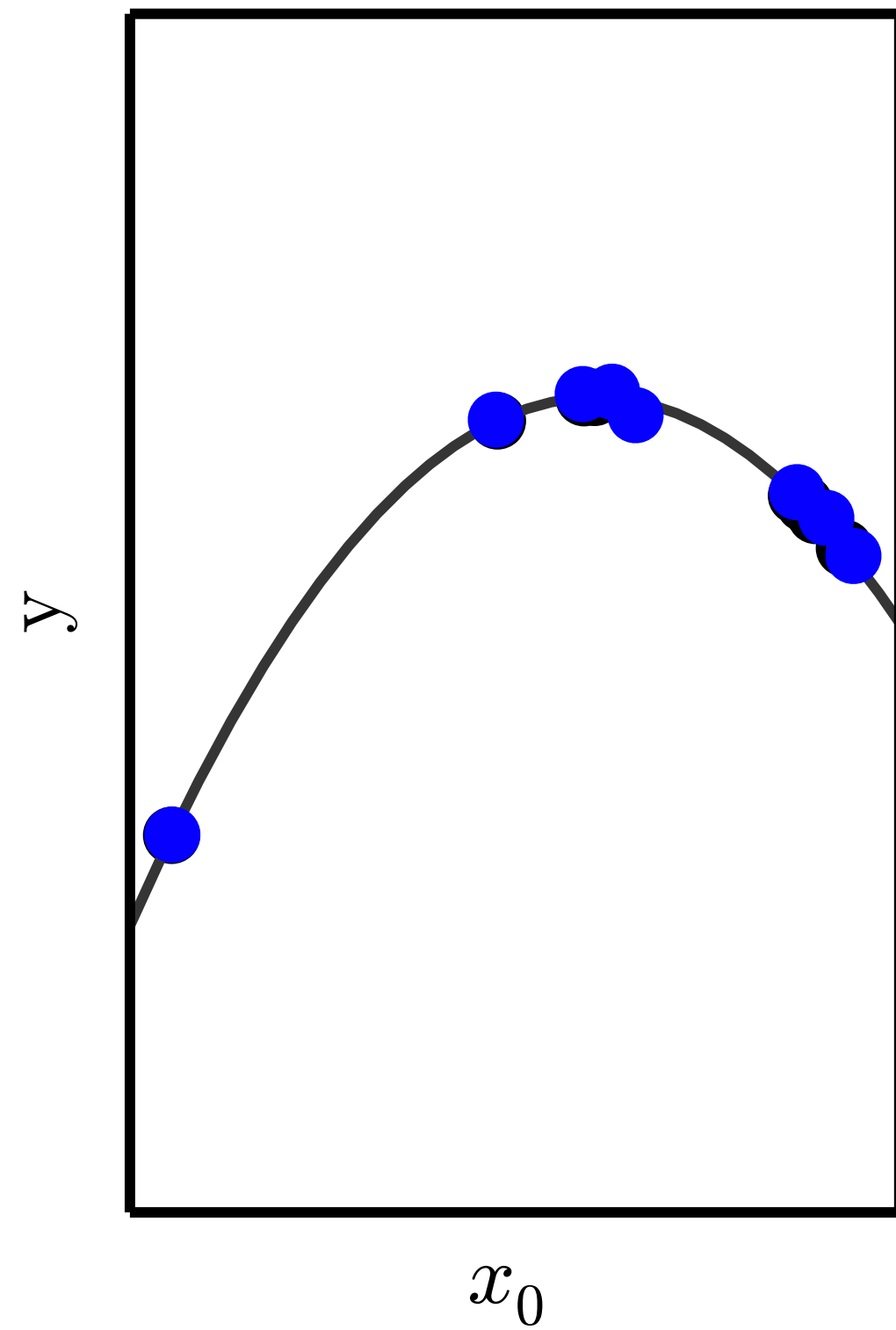$$\arg\min_{W} \mathbf{Loss}^{(X,Y)} \neq \arg\min_{W'} \mathbf{Loss}^{(X^{new},Y^{new})}$$



Optimization of $w$

MSE vs $w_1$

$$\mathbf{Loss}_1^{(X,Y)} \quad > \quad \mathbf{Loss}_2^{(X,Y)} \quad > \quad \mathbf{Loss}_3^{(X,Y)}$$

**Loss $\equiv$ MSE**



$y$      $x_0$      $y$      $x_0$      $y$      $x_0$

[Figure 5.2, Chapter 5, Deep Learning]

$$\text{Loss}_1^{(X,Y)} > \text{Loss}_2^{(X,Y)} > \text{Loss}_3^{(X,Y)}$$

**Loss ≡ MSE**

$$\text{Loss}_1^{(X,Y)} \quad > \quad \text{Loss}_2^{(X,Y)} \quad > \quad \text{Loss}_3^{(X,Y)}$$
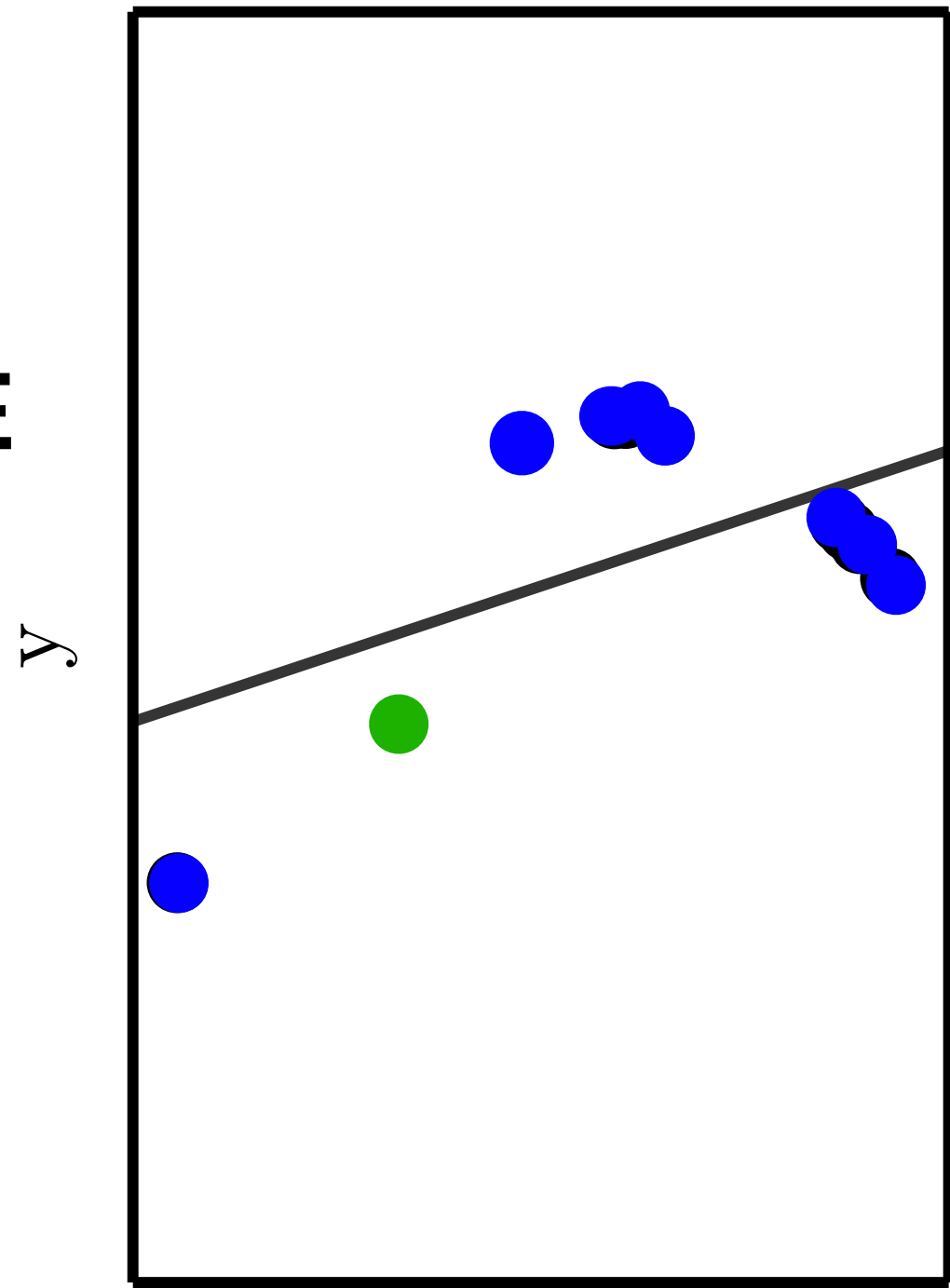
**Loss ≡ MSE**



$$\text{Loss}_1^{(X^{new},Y^{new})} \quad < \quad \text{Loss}_2^{(X^{new},Y^{new})} \quad < \quad \text{Loss}_3^{(X^{new},Y^{new})}$$

[Figure 5.2, Chapter 5, Deep Learning]

# Some Terminology

# Some Terminology

- (X,Y) : training set

- $(X^{new}, Y^{new})$ :  test set

- $Loss^{(X,Y)}$ : train loss (error)

- $Loss^{(X^{new}, Y^{new})}$ : test/generalization loss (error)

# Some Terminology

- (X,Y) : training set

- $(X^{new}, Y^{new})$ :  test set

- $Loss^{(X,Y)}$ : train loss (error)

- $Loss^{(X^{new}, Y^{new})}$ : test/generalization loss (error)

**Our goal in ML is (to train the model) to obtain small generalization error**

# Some Terminology
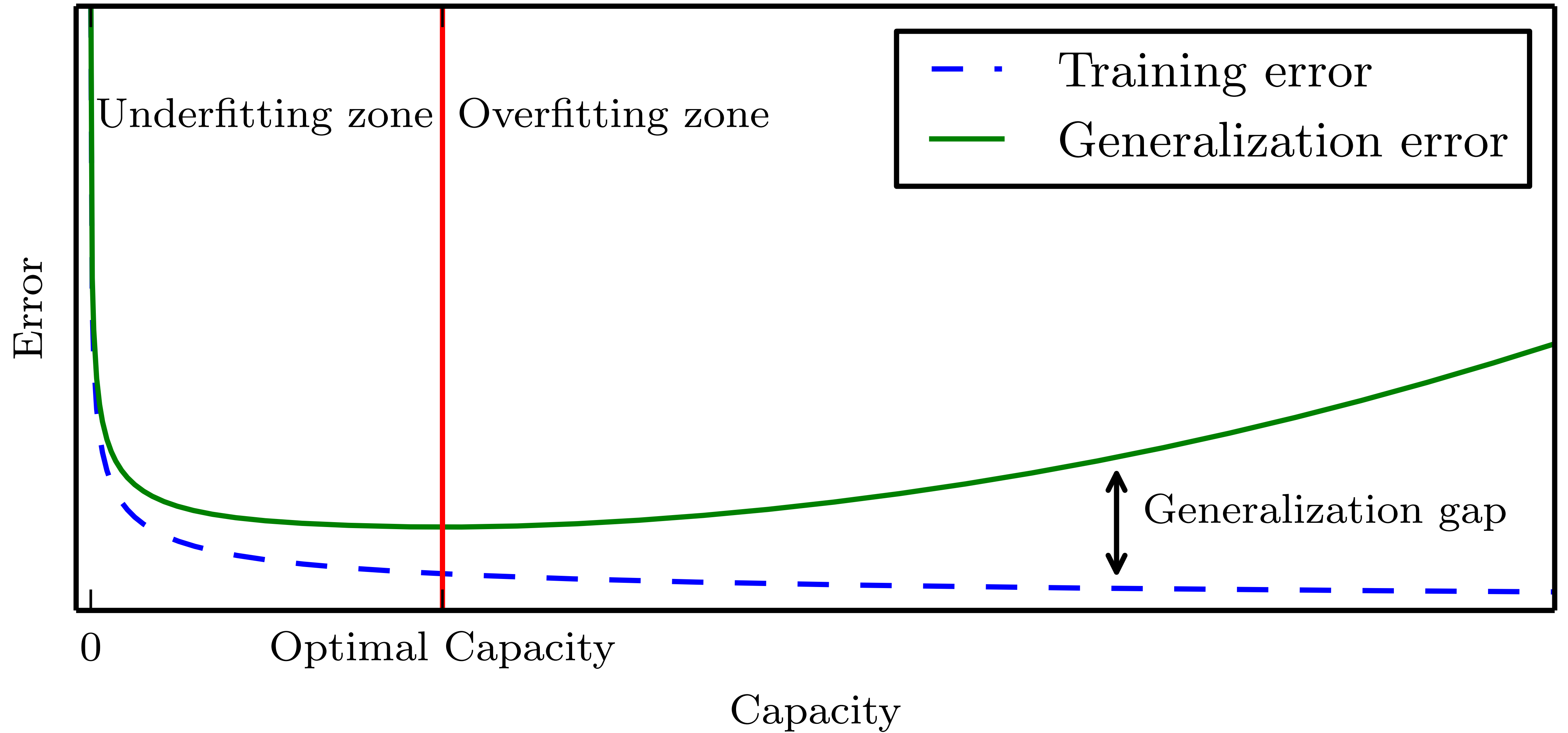
- (X,Y) : training set

- $(X^{new}, Y^{new})$ :  test set

- $\text{Loss}^{(X,Y)}$ : train loss (error)

- $\text{Loss}^{(X^{new}, Y^{new})}$ : test/generalization loss (error)

**Our goal in ML is (to train the model) to obtain small generalization error**

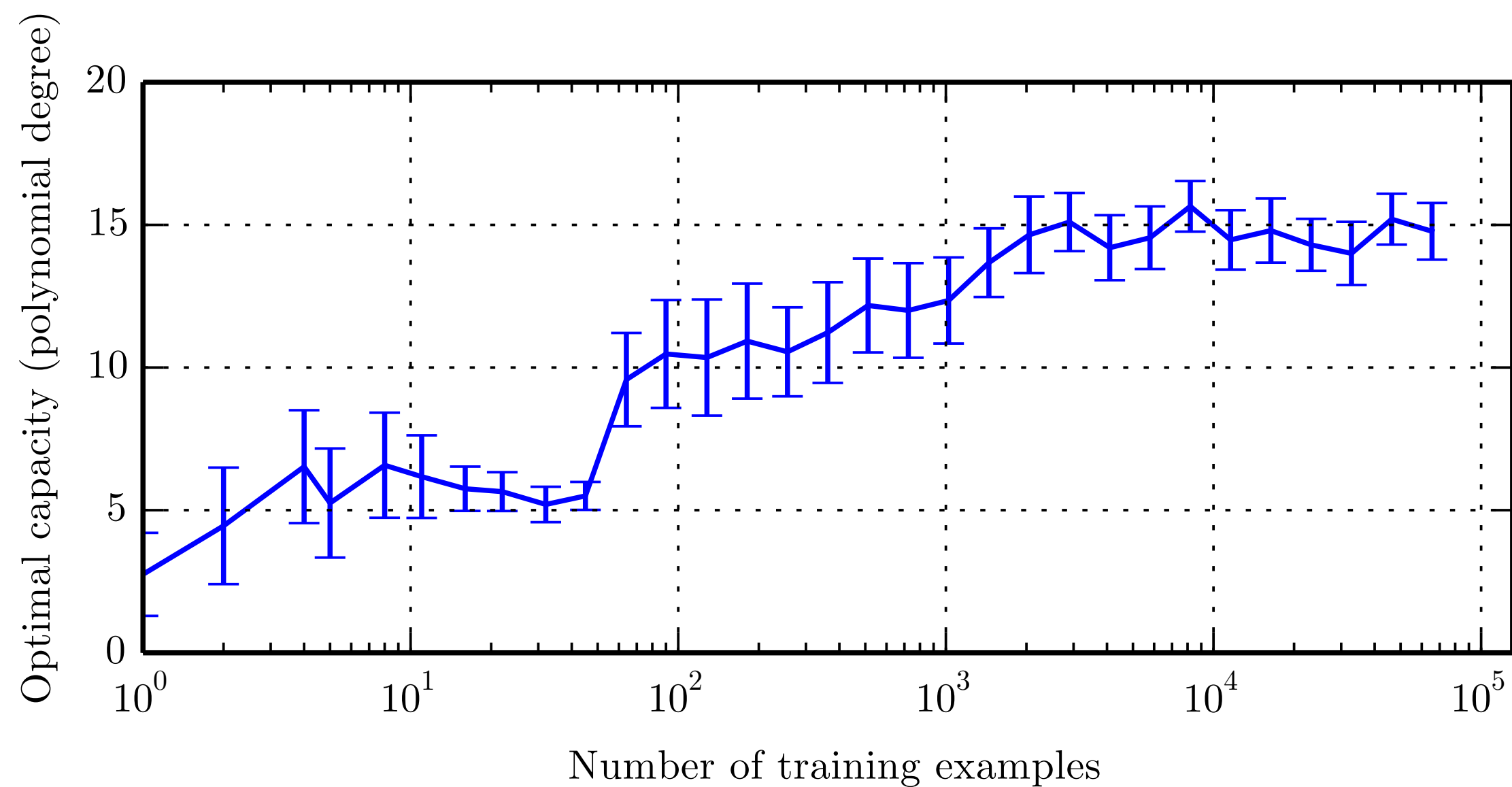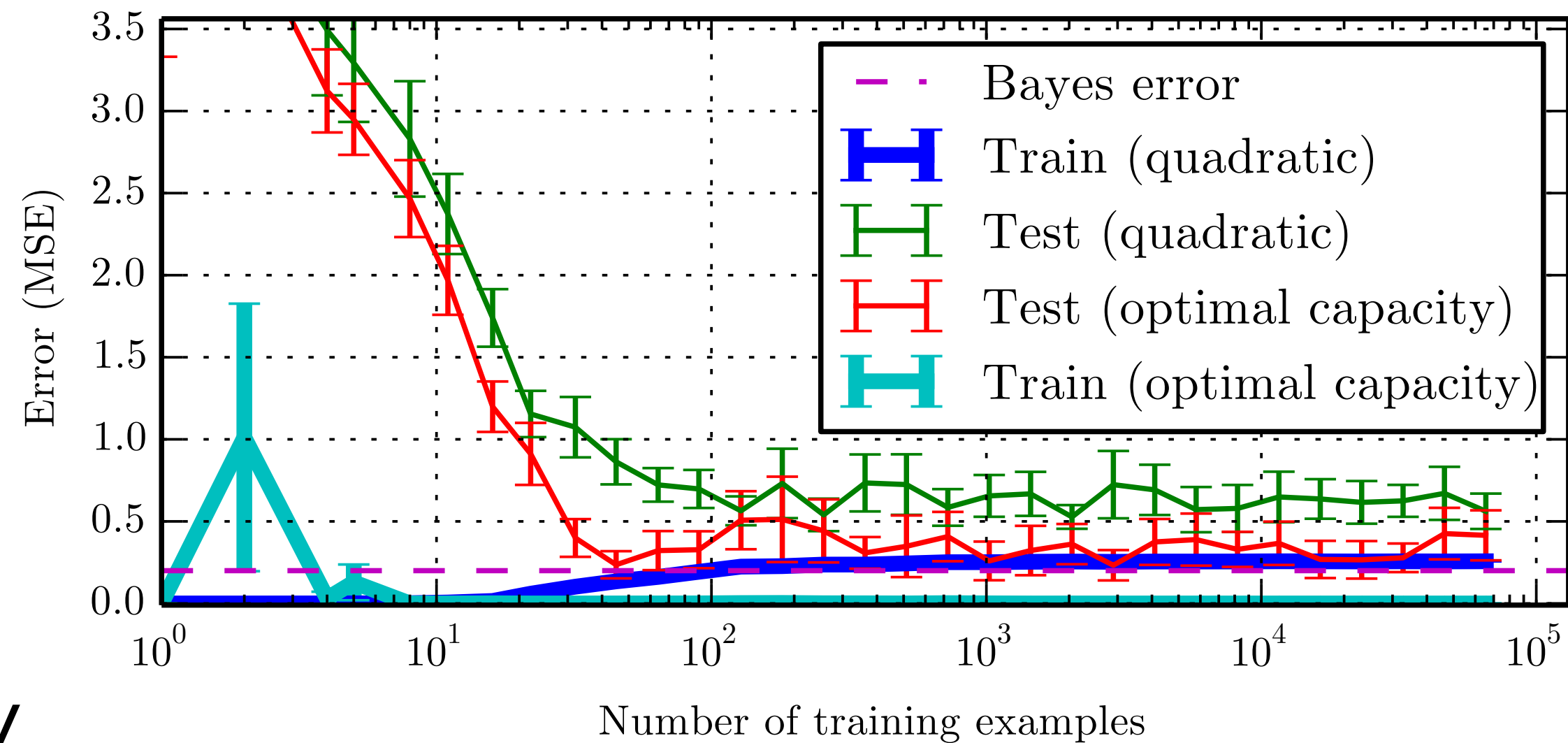- Capacity: "The ability of a model to fit a variety of functions" [DL]

$$\text{capacity}(w_2^\top x^2 + w_1^\top x) > \text{capacity}(w_1^\top x)$$

Underfitting zone | Overfitting zone

Legend: Training error (blue dashed), Generalization error (green solid)

Error (y-axis), Capacity (x-axis)

0 | Optimal Capacity

Generalization gap

[Figure 5.3, Chapter 5, Deep Learning]

Synthetic data is generated
using a degree 5 polynomial

$$y = w_5x^5 + w_4x^4 + w_3x^3 + w_2x^2 + w_1x^1$$

Training set size
also plays an important
role in a model's capacity
to generalize

[Figure 5.4, Chapter 5, Deep Learning]

# Formal learning guarantees

- It is possible to bound the generalization gap

  - Bounds involve:

    - the size of the training set

    - the capacity of the learning model

# Informally

- Larger datasets (train) are helpful

  - Allow you to better fit models and/or fit more complex models

- Larger capacity models can be better but (all being equal) they will require more data

# Regularization

# Regularization

- Can affect a model's effective capacity

  - Instead of changing the model (reminder: polynomials)

  - Focusses on particular (good) solutions
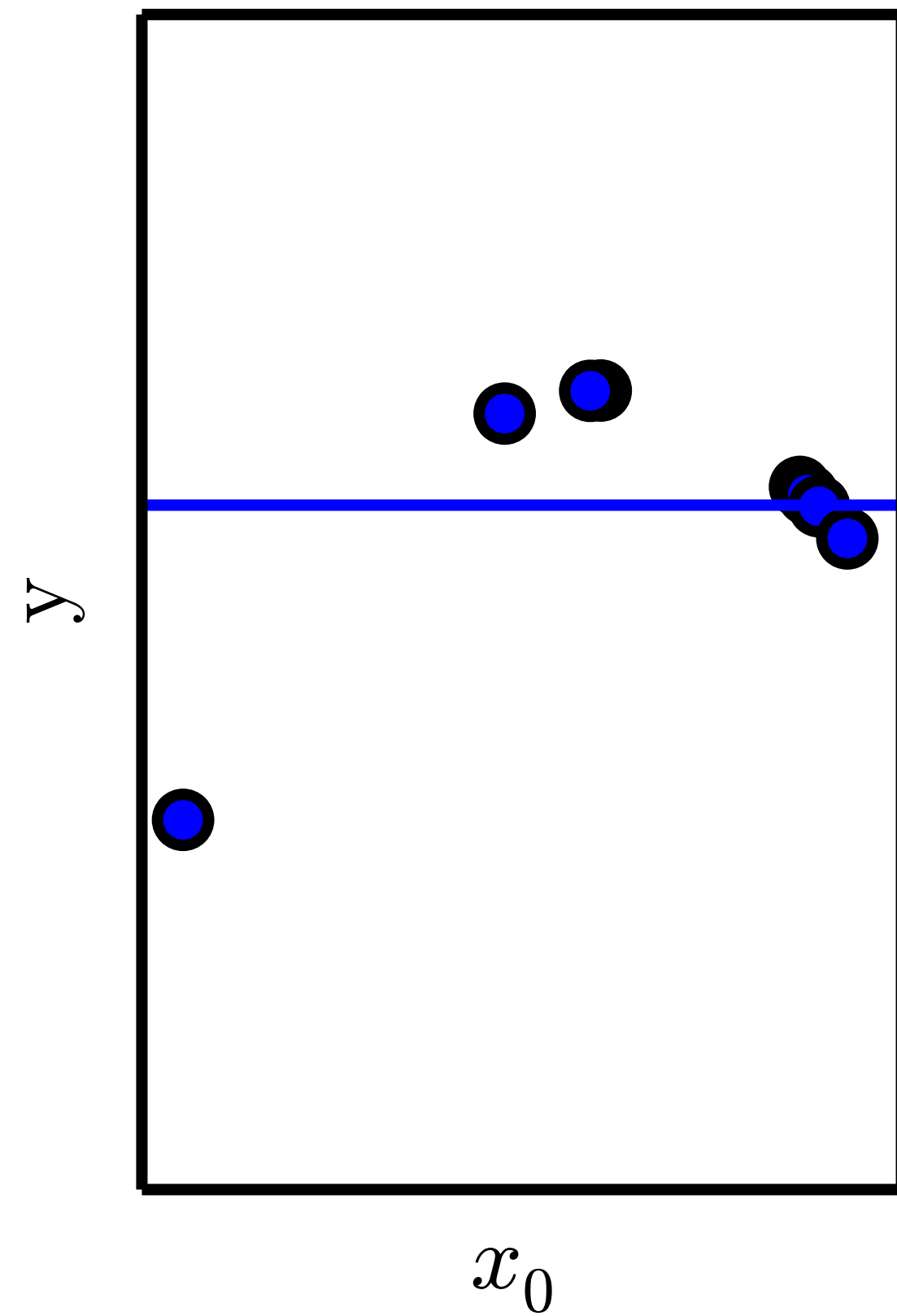
# L2-regularization

- A popular form of regularization

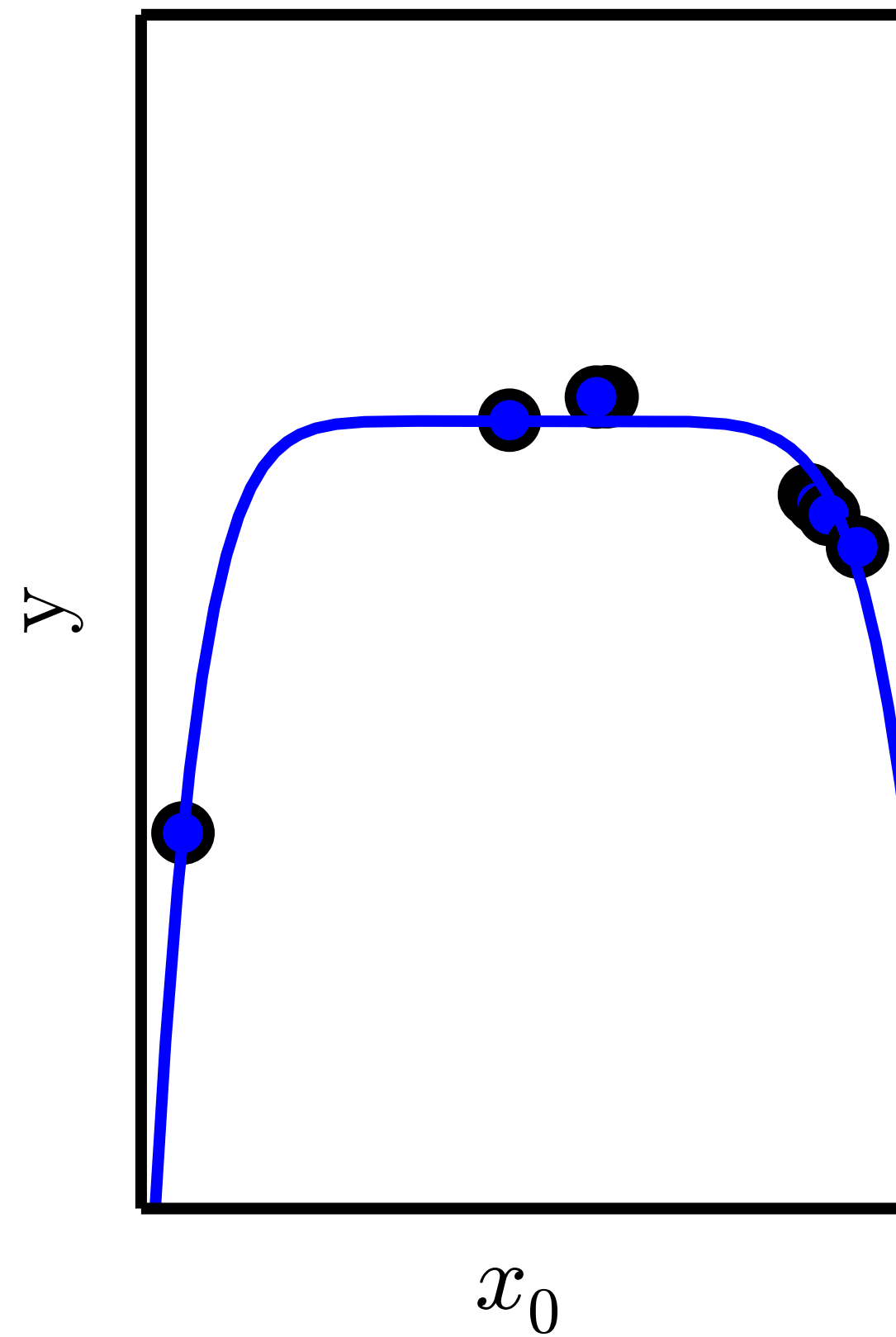$$\text{Loss} := \text{MSE}^{\text{train}} + \lambda \underbrace{\mathbf{w}^\top \mathbf{w}}_{\|\mathbf{w}\|_2}$$

- Penalizes the size of the weights

  - Smaller weights means simpler models (next slide)

$$\text{Loss} := \text{MSE}^{\text{train}} + \lambda\mathbf{w}^{\top}\mathbf{w}$$
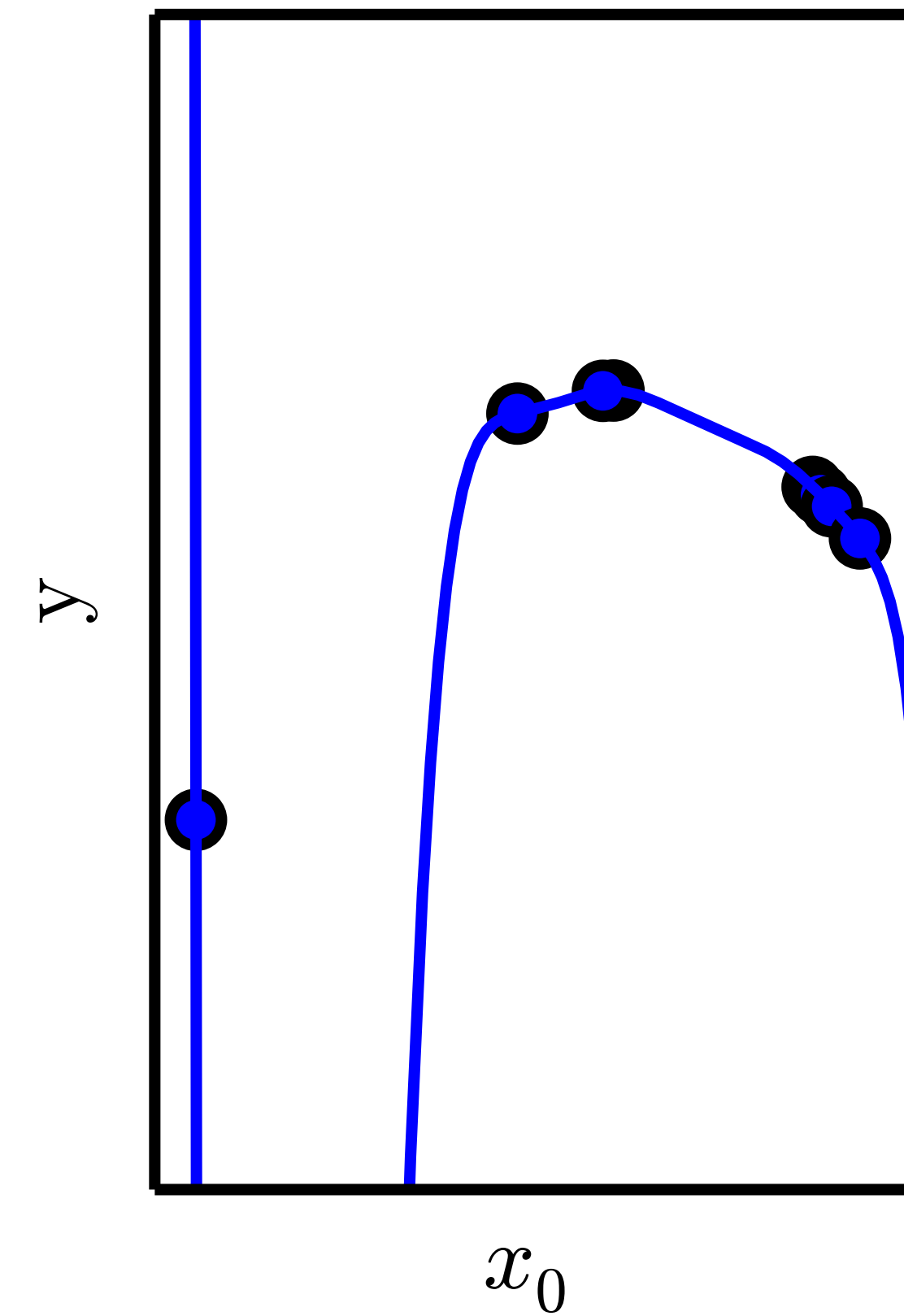


Underfitting
(Excessive $\lambda$)

Appropriate weight decay
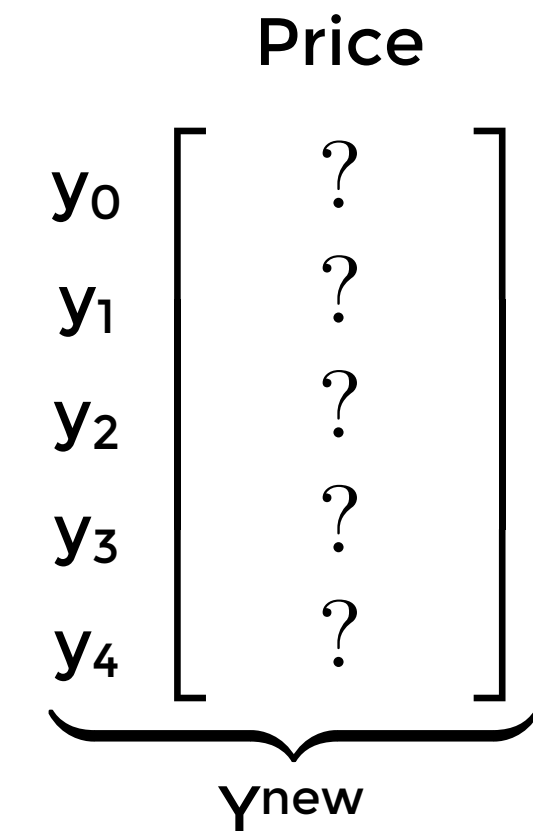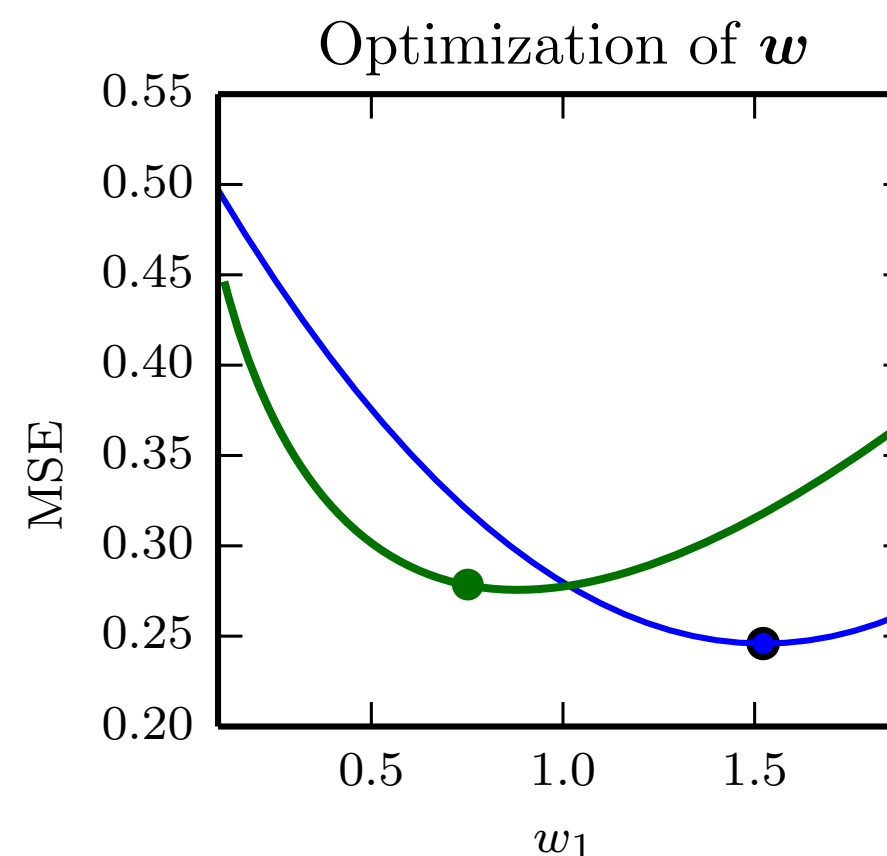(Medium $\lambda$)

Overfitting
($\lambda \rightarrow 0$)

[Figure 5.5, Chapter 5, Deep Learning]

# Validating a model

# Generalization

Price

$$y_0 \begin{bmatrix} ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix}$$
$y_1$
$y_2$
$y_3$
$y_4$

$Y^{new}$

- $\text{Loss}^{(X,Y)}$ : The one you can evaluate

- $\text{Loss}^{(X^{new},Y^{new})}$ : The one that you care about

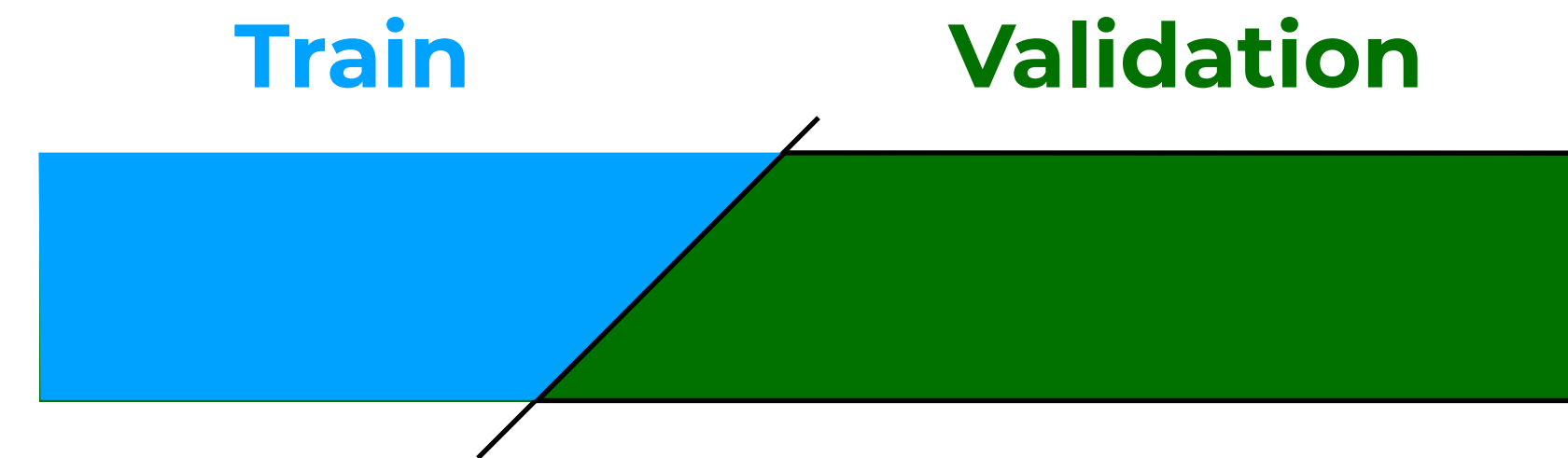- In general minimizing the former will not yield the best loss on the latter:

$$\arg\min_{w} \textbf{Loss}^{\textbf{(X,Y)}} \neq \arg\min_{w'} \textbf{Loss}^{\textbf{(X}^{\textbf{new}}\textbf{,Y}^{\textbf{new}}\textbf{)}}$$
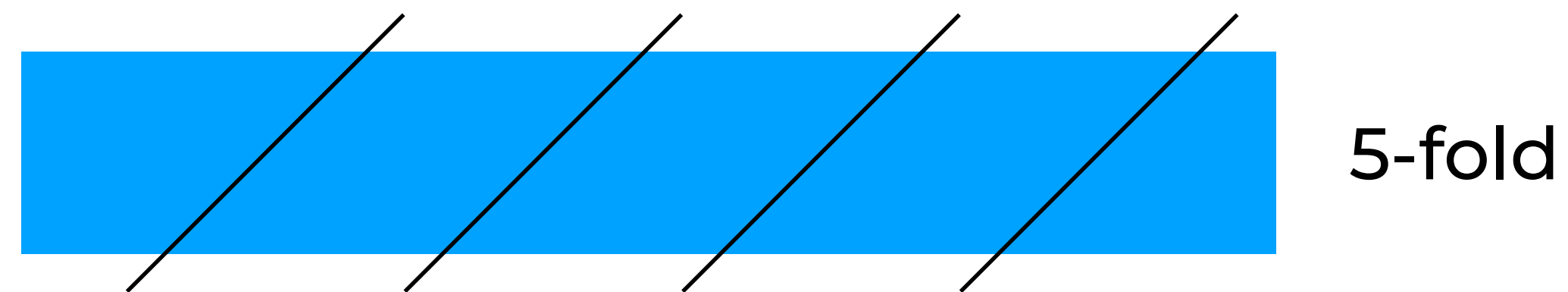

Optimization of $w$

# Validation set

- How do we choose the right model and set its hyper parameters (e.g. $\lambda$)?

  - Use a validation set

    - Split the original data into two:

      1. Train set

      2. Validation set

        - Proxy to the test set

  - Train different models/hyperparameter settings on the train set

  - Pick the best according to their performance on the validation set
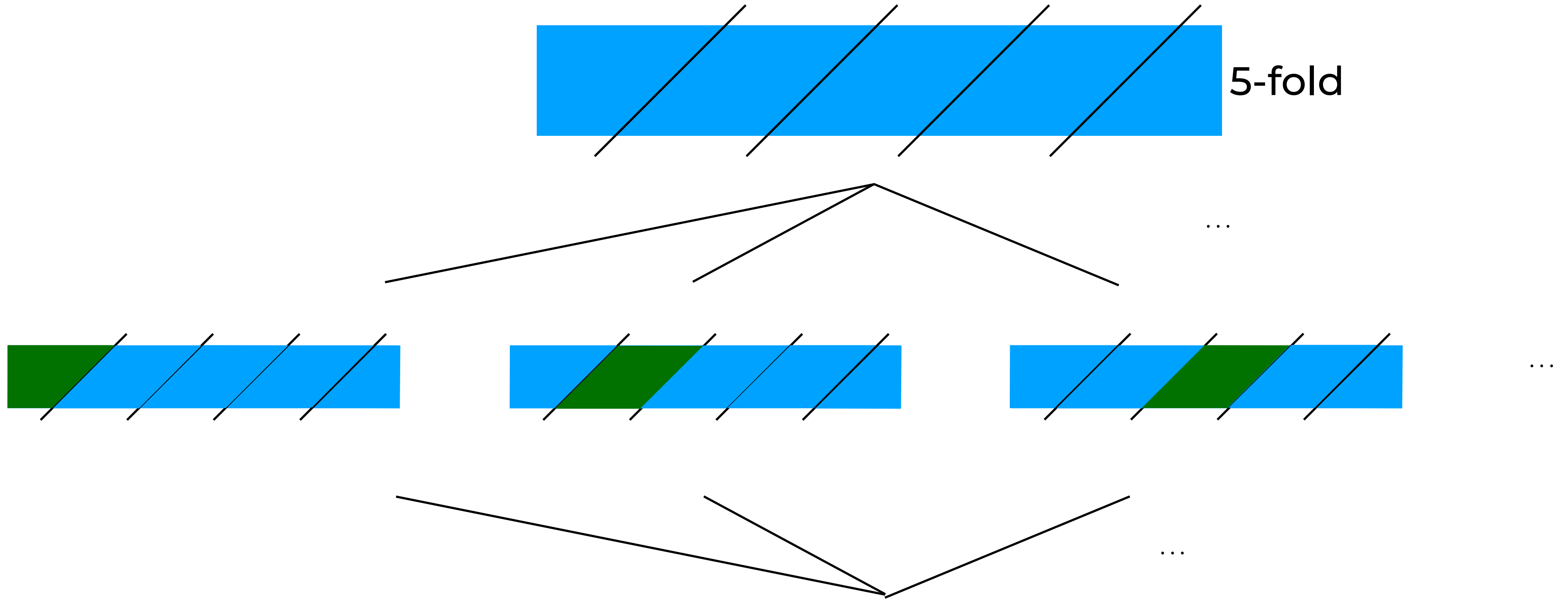
**Train**   **Validation**

# Cross-validation (CV)

- Splitting the data into train/validation can be detrimental

  - e.g., if data is small to begin with (small train and validation sets)

- K-fold CV: Split the data into k-folds
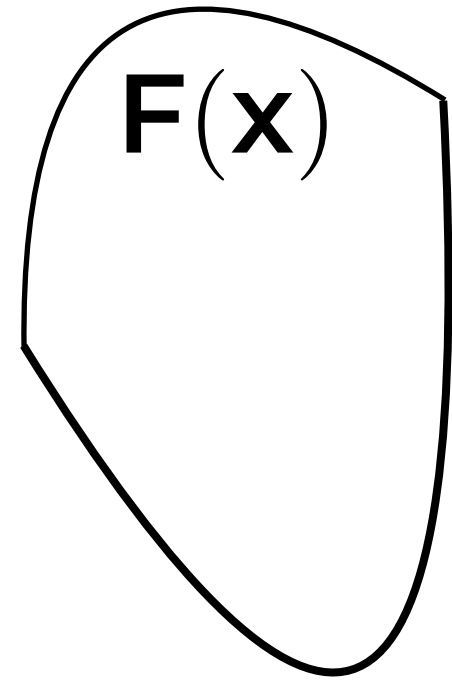


5-fold

5-fold

...

**Train**

**Validation**

...

...

Pick the model/hyperparameters that
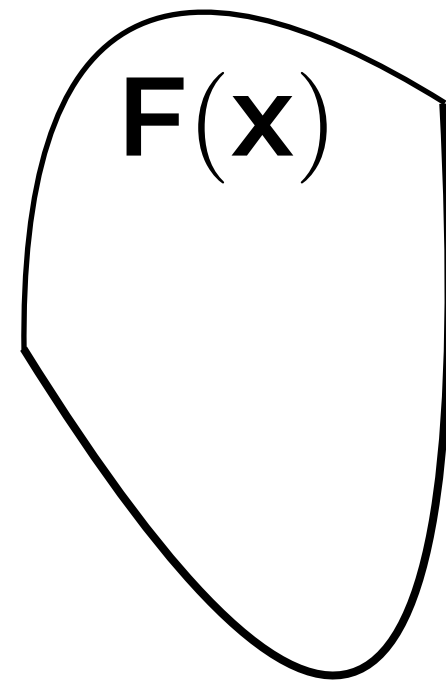does best (e.g., smallest loss) according to
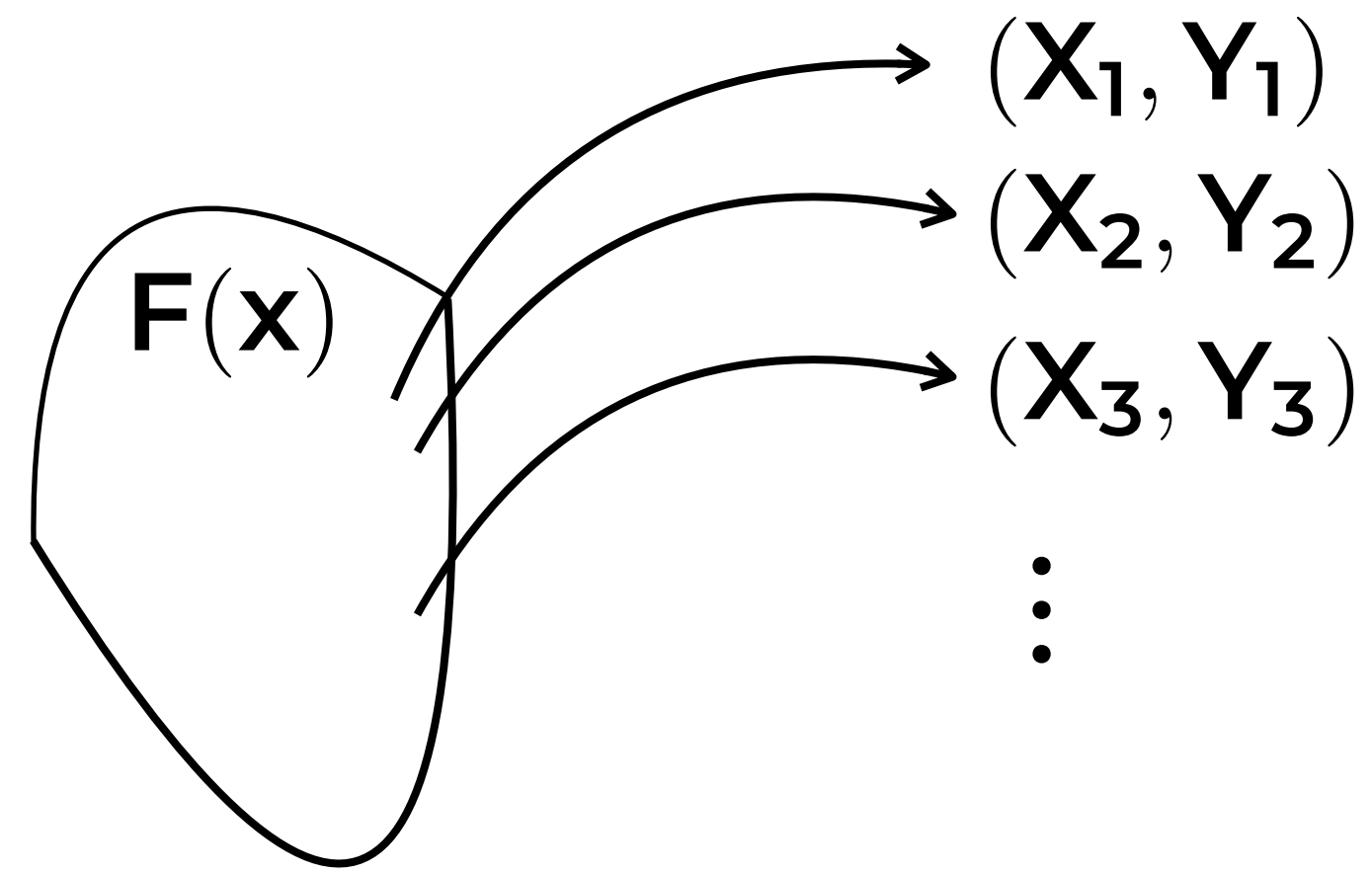the average of the validation sets
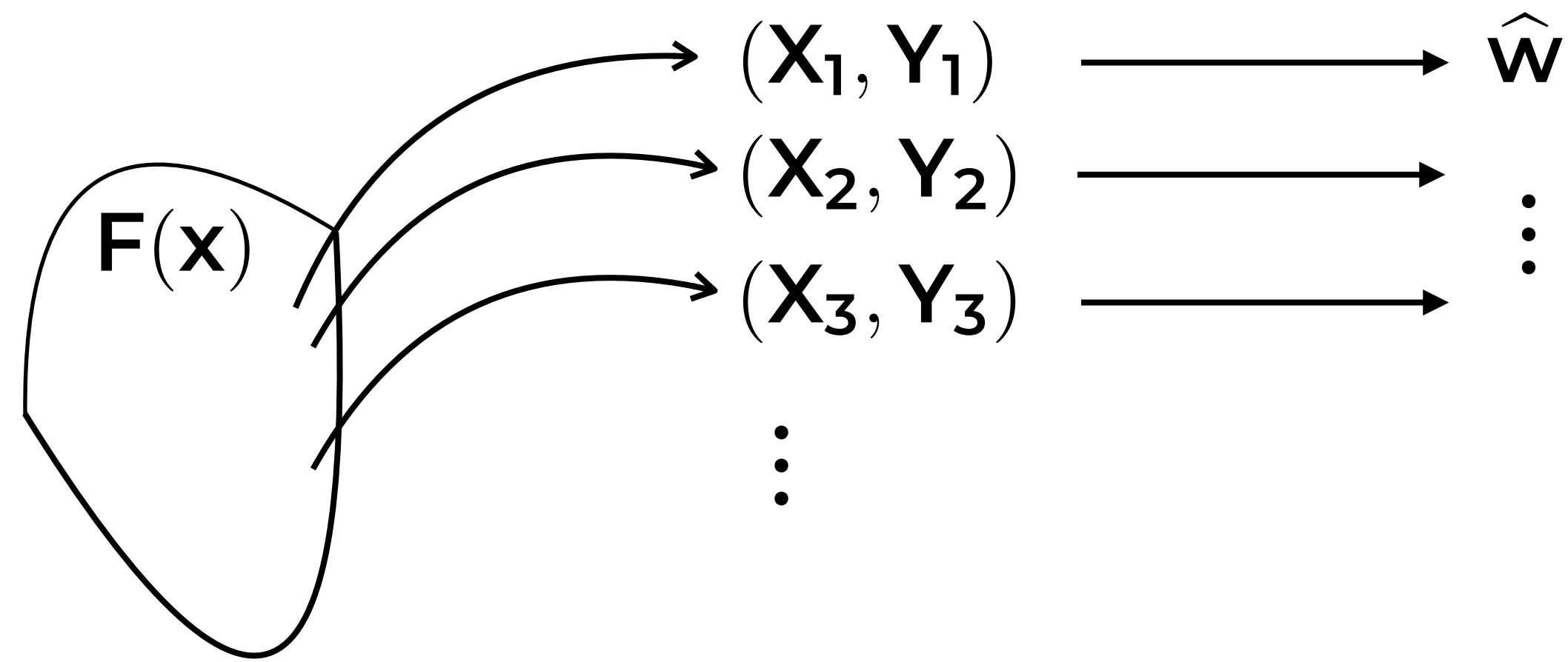
# Bias/Variance:
## A second perspective on generalization

F(x)

**F(x)**

$F(x) =$ **linear regression with parameters** $w^*$

$F(x) = w_0^* + w_1^* x + w_2^* x^2 + w_3^* x^3$

$F(x) = $ **linear regression with parameters $w^*$**

$$F(x) = w_0^* + w_1^* x + w_2^* x^2 + w_3^* x^3$$

$(X_1, Y_1) \longrightarrow \widehat{w}$

$(X_2, Y_2) \longrightarrow$

$(X_3, Y_3) \longrightarrow$

$F(x)$

$F(x) =$ linear regression with parameters $w^*$

$F(x) = w_0^* + w_1^* x + w_2^* x^2 + w_3^* x^3$

# Bias/Variance tradeoff (in 4 slides)

- An alternate framework

# Bias/Variance tradeoff (in 4 slides)

- **An alternate framework**

$$w^* := \text{the value of parameters that generated the data}$$
$$\widehat{w} := \text{estimator of the true } w^*$$

# Bias/Variance tradeoff (in 4 slides)

- An alternate framework

$$w^* := \text{the value of parameters that generated the data}$$

$$\widehat{w} := \text{estimator of the true } w^*$$

$$\textbf{MSE} := \mathbb{E}[(\widehat{w} - w^*)^2]$$

$$= \textbf{Bias}(\widehat{w})^2 + \textbf{Var}(\widehat{w})$$

# Bias/Variance tradeoff (in 4 slides)

- An alternate framework
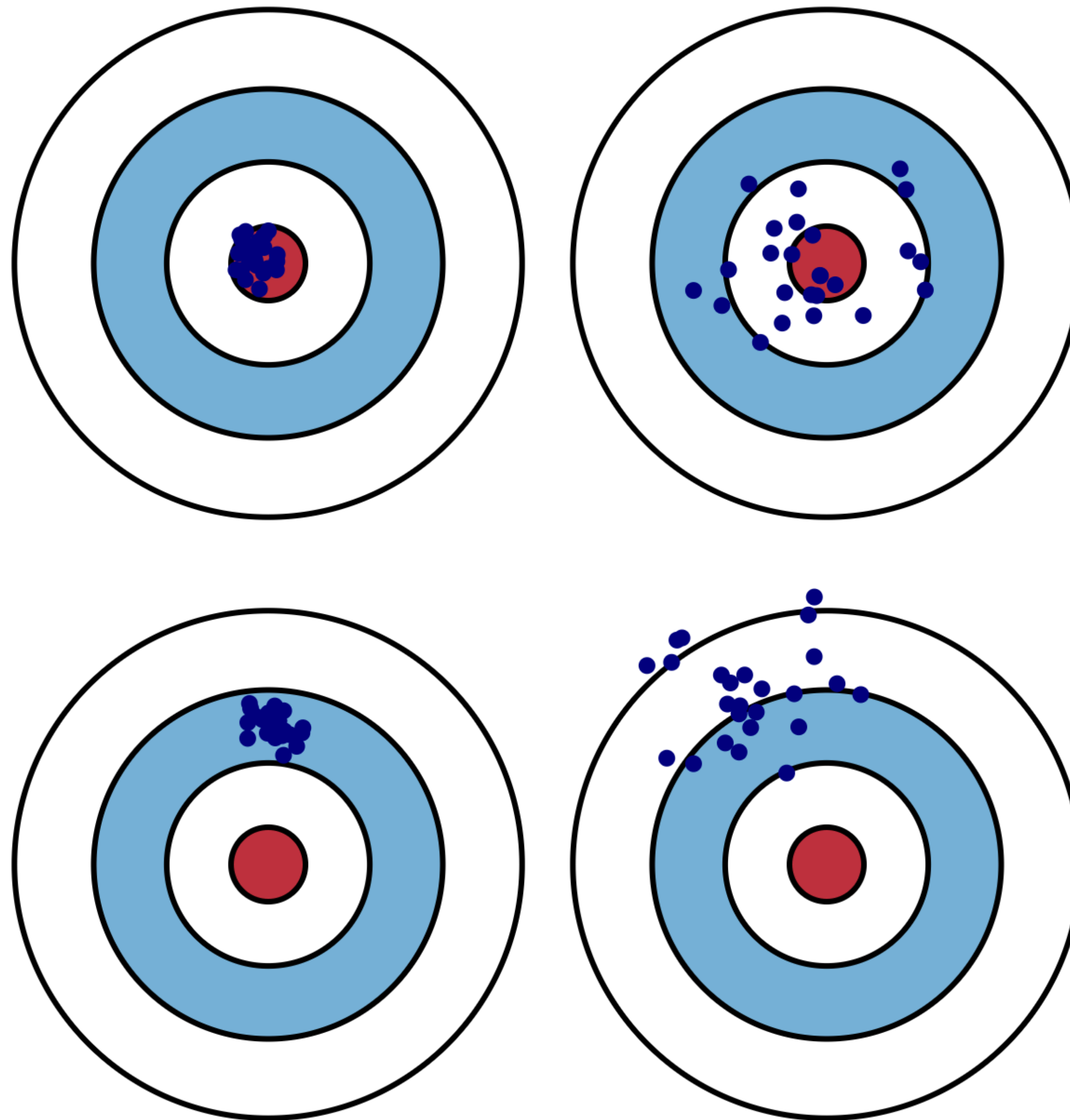
$w^* :=$ **the value of parameters that generated the data**

$\widehat{w} :=$ **estimator of the true** $w^*$

$$\mathbf{MSE} := \mathbb{E}[(\widehat{w} - w^*)^2]$$

$$= \mathbf{Bias}(\widehat{w})^2 + \mathbf{Var}(\widehat{w})$$

$$\mathbf{Bias}(\widehat{w}) = \mathbb{E}[\widehat{w}] - w^*$$

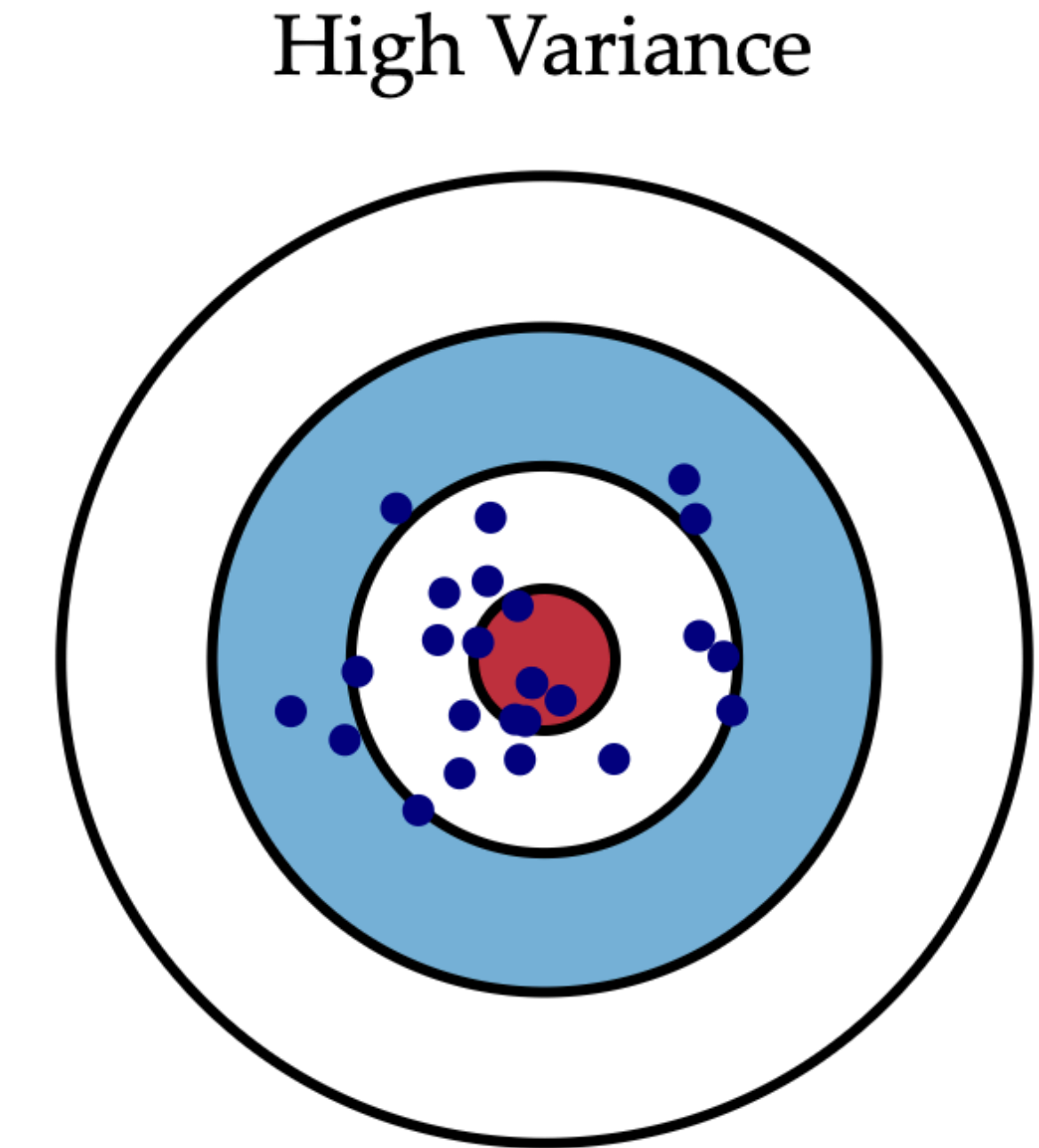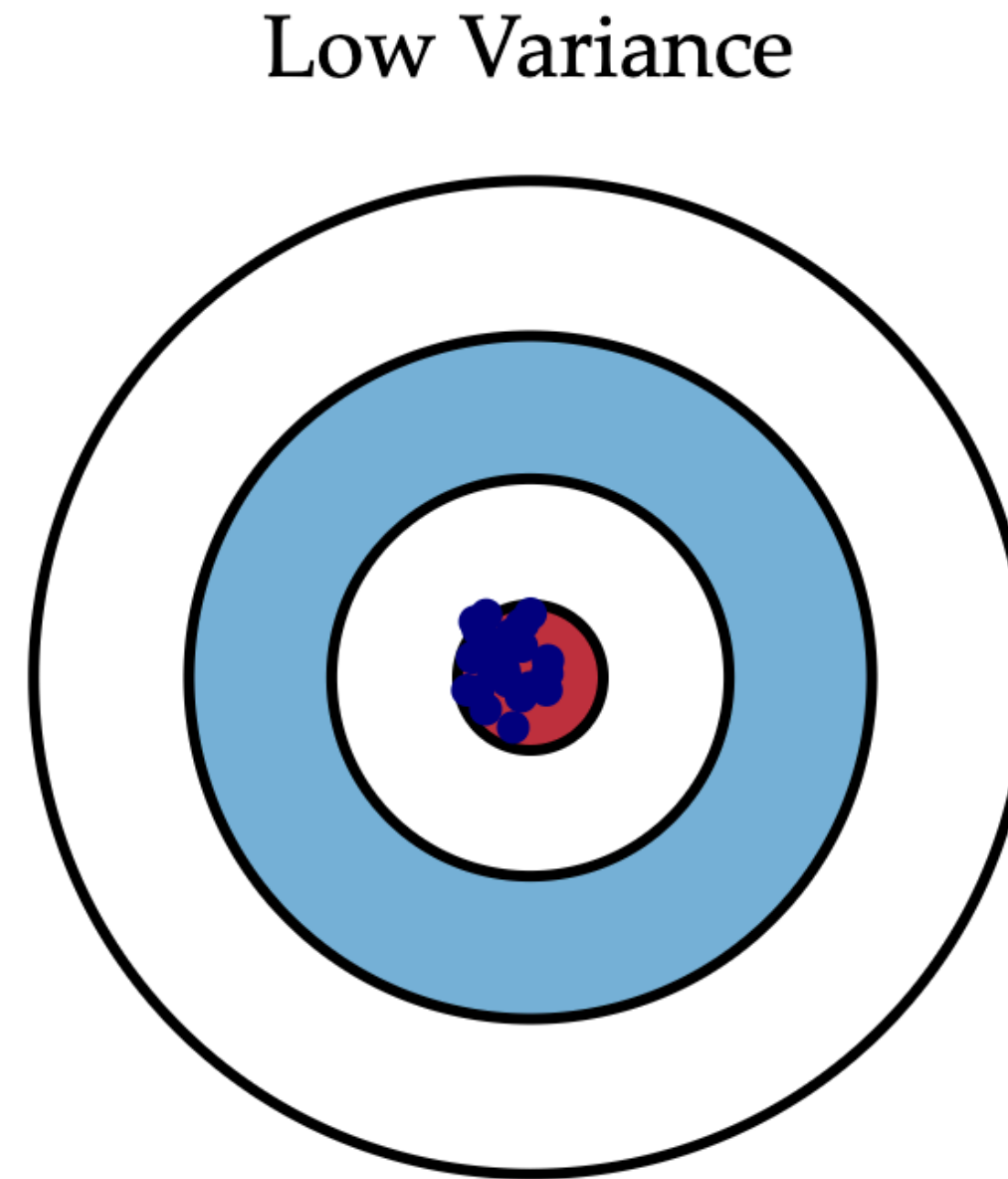$$\mathbf{Var}(\widehat{w}) = \mathbb{E}[(\widehat{w} - \mathbb{E}[\widehat{w}])^2]$$

- The goal is to hit the bull's eye
- Each blue dot represents the "performance" of a fixed model on different data from the same distribution

http://scott.fortmann-roe.com/docs/BiasVariance.html
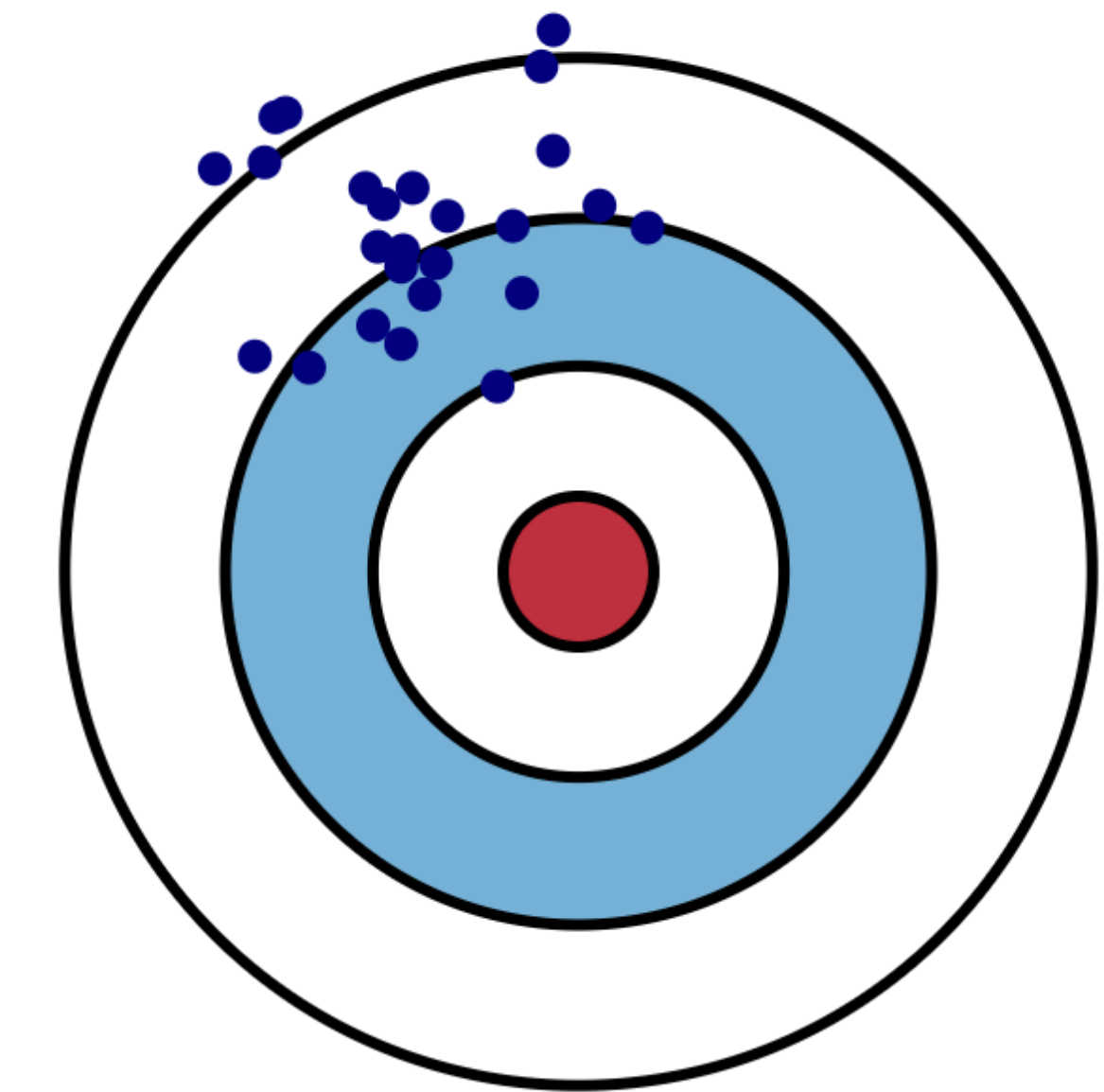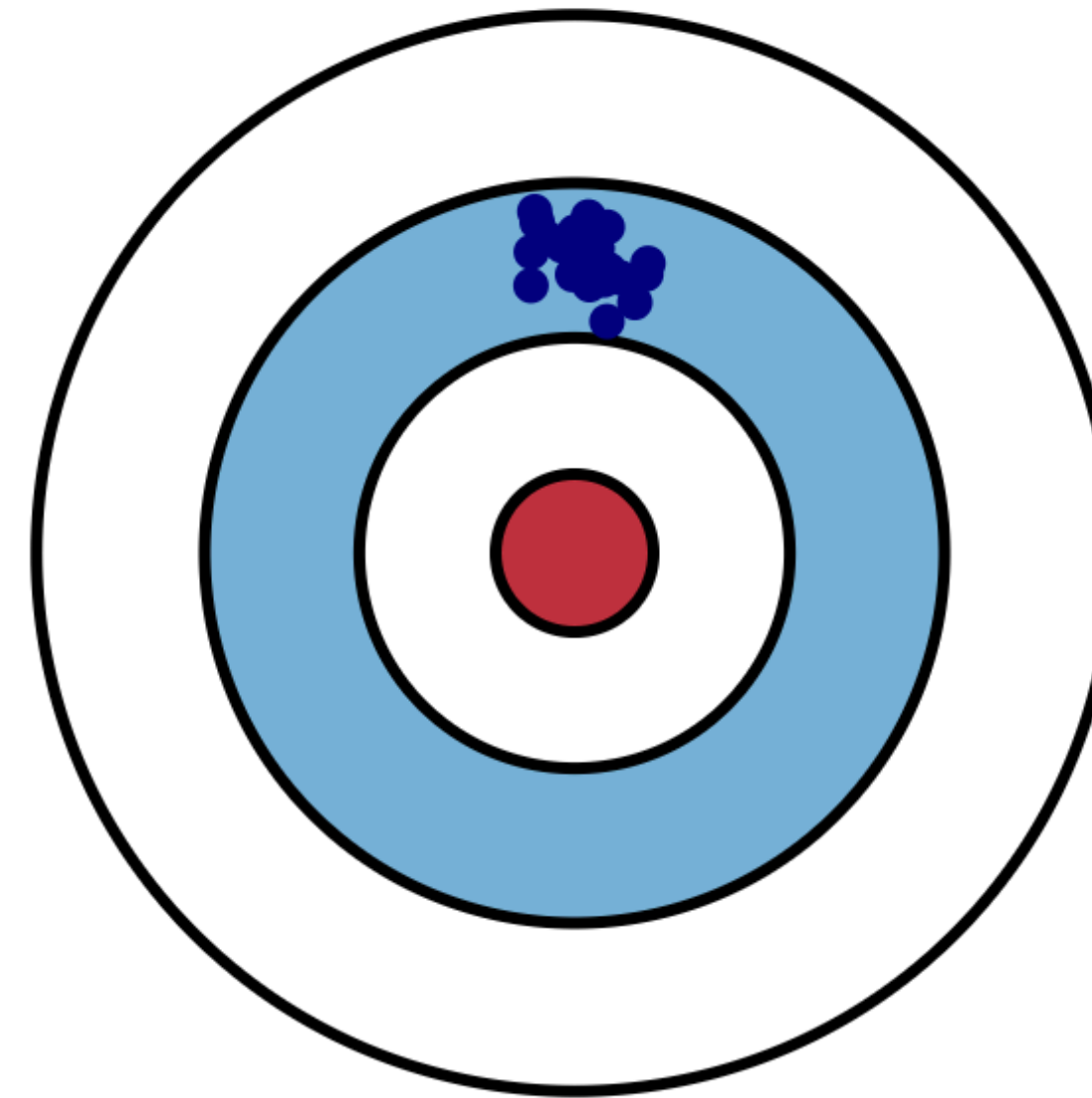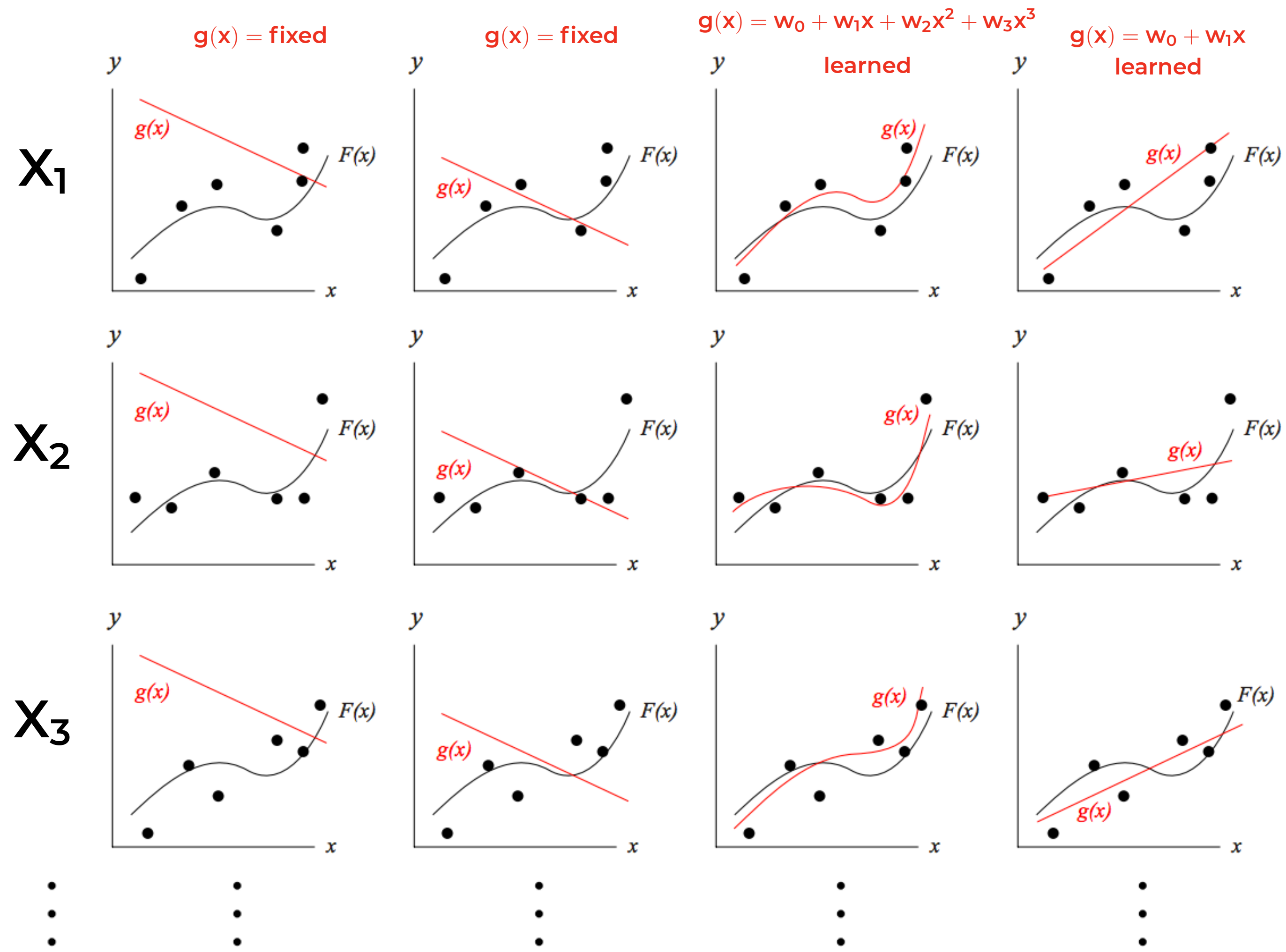
- The goal is to hit the bull's eye
- Each blue dot represents the "performance" of a fixed model on different data from the same distribution

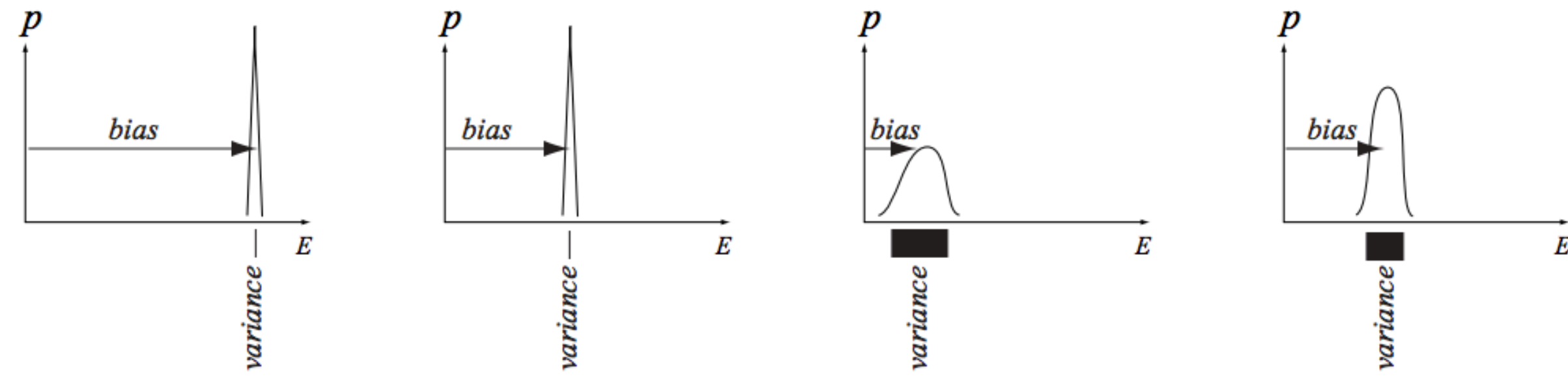http://scott.fortmann-roe.com/docs/BiasVariance.html

Top row column headers:
$g(x) =$ fixed    $g(x) =$ fixed    $g(x) = w_0 + w_1x + w_2x^2 + w_3x^3$ learned    $g(x) = w_0 + w_1x$ learned

Row labels: $X_1$, $X_2$, $X_3$

Histogram of the error over training sets

[Figure 9.2. Pattern Classification. Duda, Hart, Stork. 2001]

Underfitting zone

Overfitting zone

Bias

Generalization error

Variance

Optimal capacity

Capacity

# Other frameworks

- **Bayesian**

  - **Uncertainty indicates you degree of belief**

  - **Unknown quantities are random variables**

# Other evaluations

- Is test set evaluation enough?

  - The test error may be a proxy for what you are really trying to evaluate

  - You model may be used inside a larger system

  - How can you convince that an X % improvement in test error is meaningful?

# Other evaluations

- Model exploration

  - Are the parameter values it has learned sensible?

  - Plot the residuals

  - Dive into your model's predictions

    - Where does it do better/worse than others?

- Model criticism

  - How do generated data from your fitted model look like?