

CMPT 318
Data Analytics

Learning Algorithms

Adapted in part from
Chapter 11 of Data Science from Scratch (Grus)
Chapter 5 of Deep Learning (Goodfellow)

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 , as measured by P , improves with experience E .

Tasks

- Classification
- Regression
- NLP: Transcription, Translation
- Structured output
- Anomaly detection
- Synthesis and sampling
- Density estimation

Computer Vision Tasks

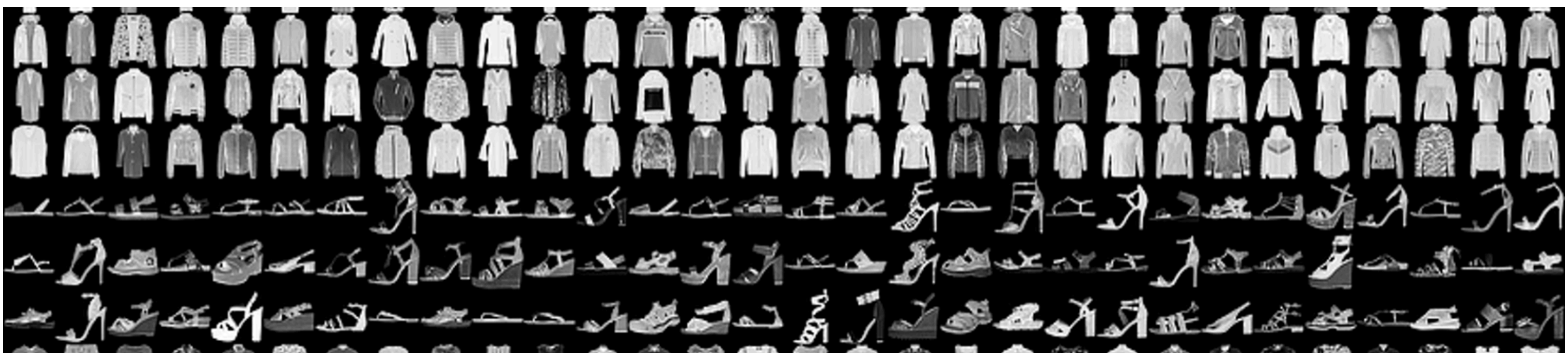
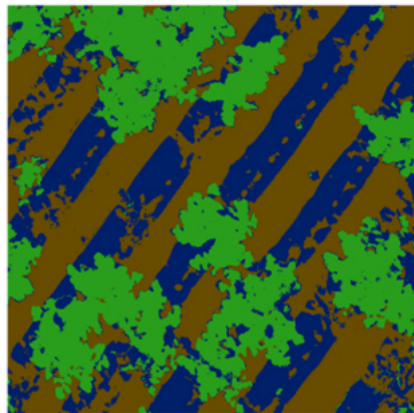


Image Classification

Computer Vision Tasks



Regression



Semantic
Segmentation



Object Detection



Instance
Segmentation

Performance

- Measure the “accuracy” of the model
- Easy to measure for some tasks, e.g. classification
- Hard to measure for others, e.g. density estimation
- May be intractable

Experience

- aka, the *Data* (\mathbf{x}) & *Labels* (y)
- Unsupervised: learn $p(\mathbf{x})$
- Supervised: learn $p(\mathbf{y} \mid \mathbf{x})$
- *Form design matrix:* $\mathbf{X} \in \mathbb{R}^{150 \times 4}$
 - *150 examples (rows); 4 features (columns)*

Generalization

- In practice: minimize training error
- Real goal: minimize generalization error
 - “expected value of the error on a new input”
- Assume: 1) examples independent
2) test/training are identically distributed

Model Capacity

- ↑ Increase number of features/parameters
- ↑ Represent more functions (families of functions)
- ↓ Imperfect optimization procedure

Capacity should fit the problem/data

Generalization and Capacity

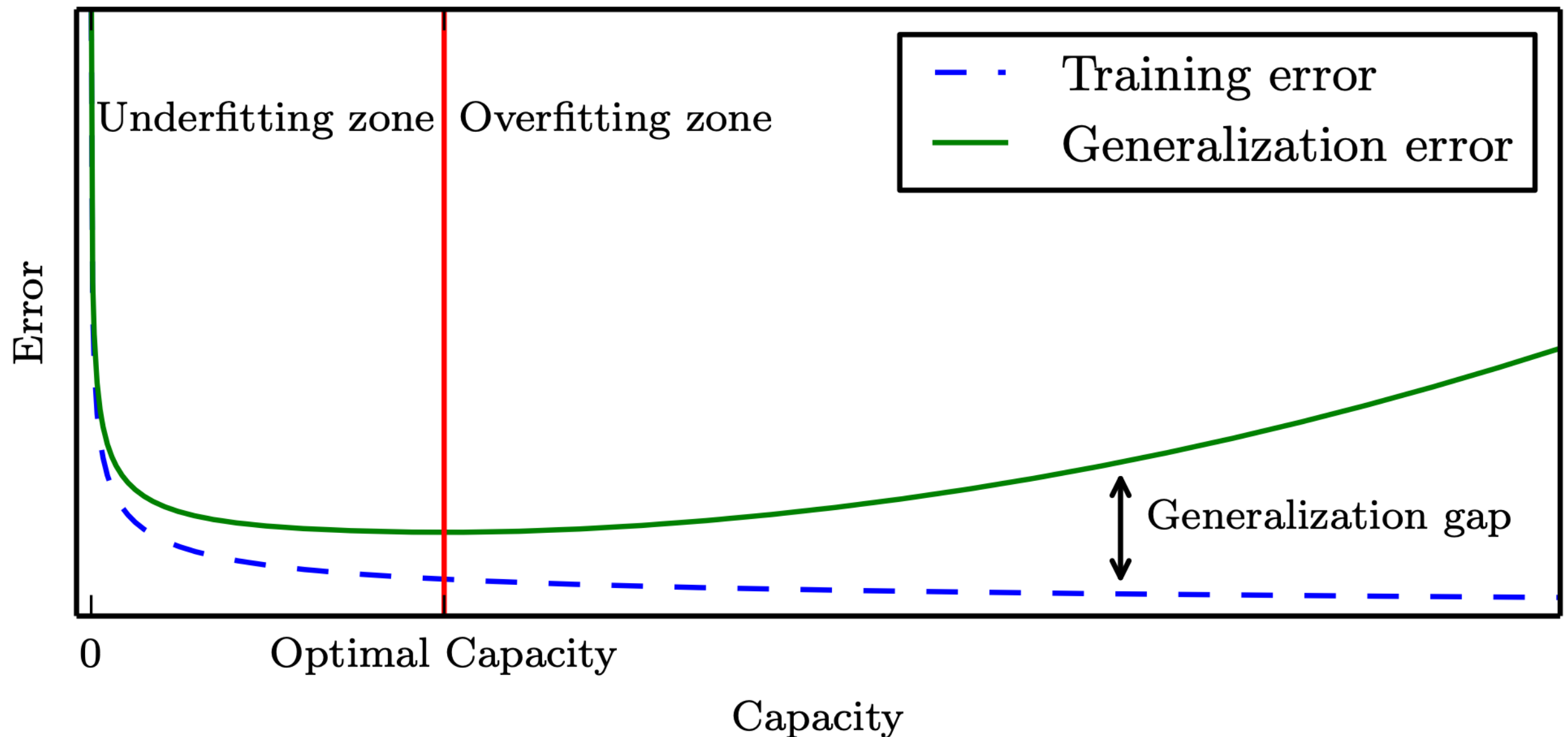


Figure 5.3

Increasing Model Capacity

Recall

Polynomial Model

$$\hat{y} = \sum_j \beta_j x^j$$

Underfitting and Overfitting in Polynomial Estimation

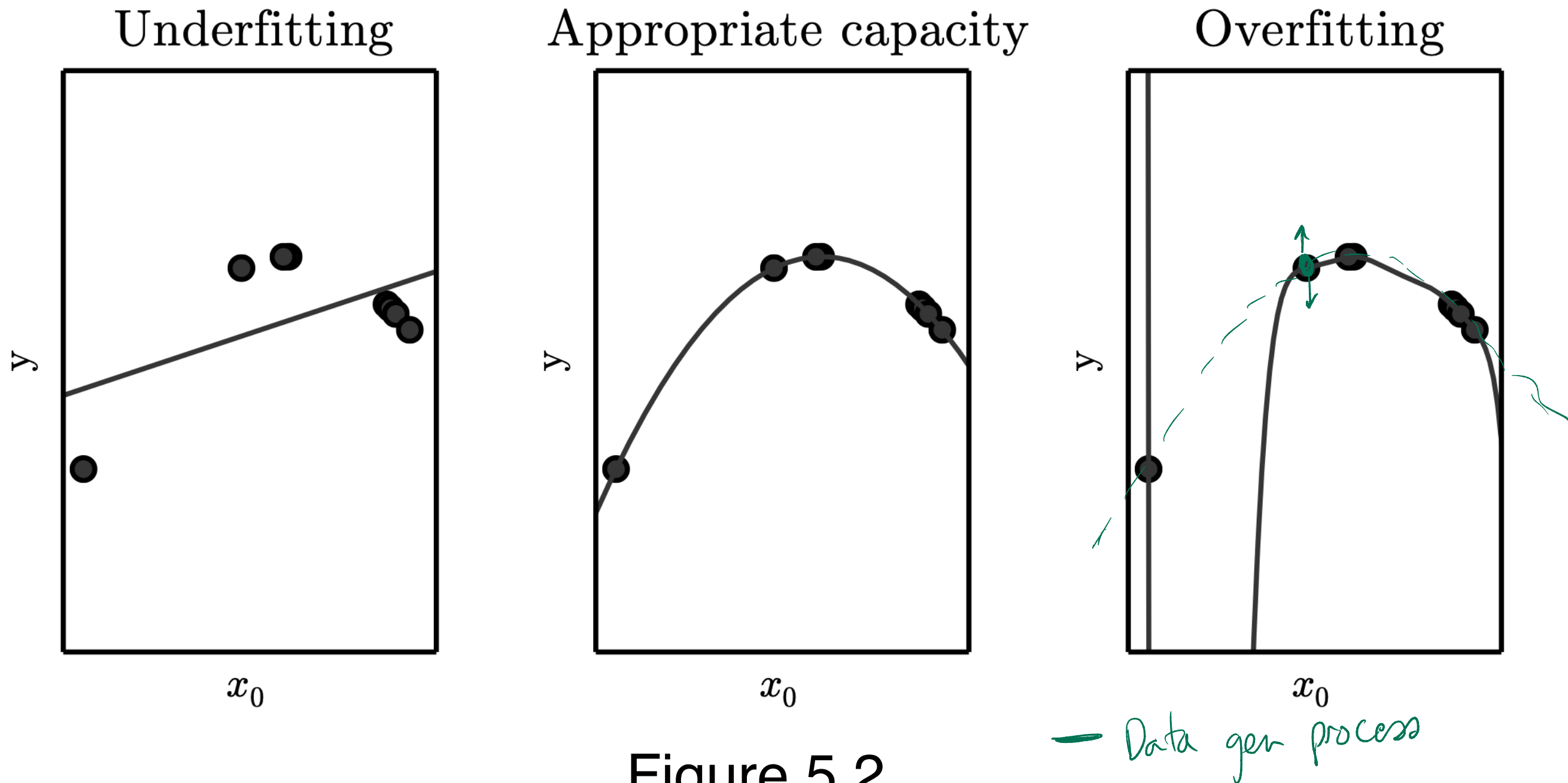


Figure 5.2

Training Set Size Example

synthetic regression problem for noisy degree-5 polynomial

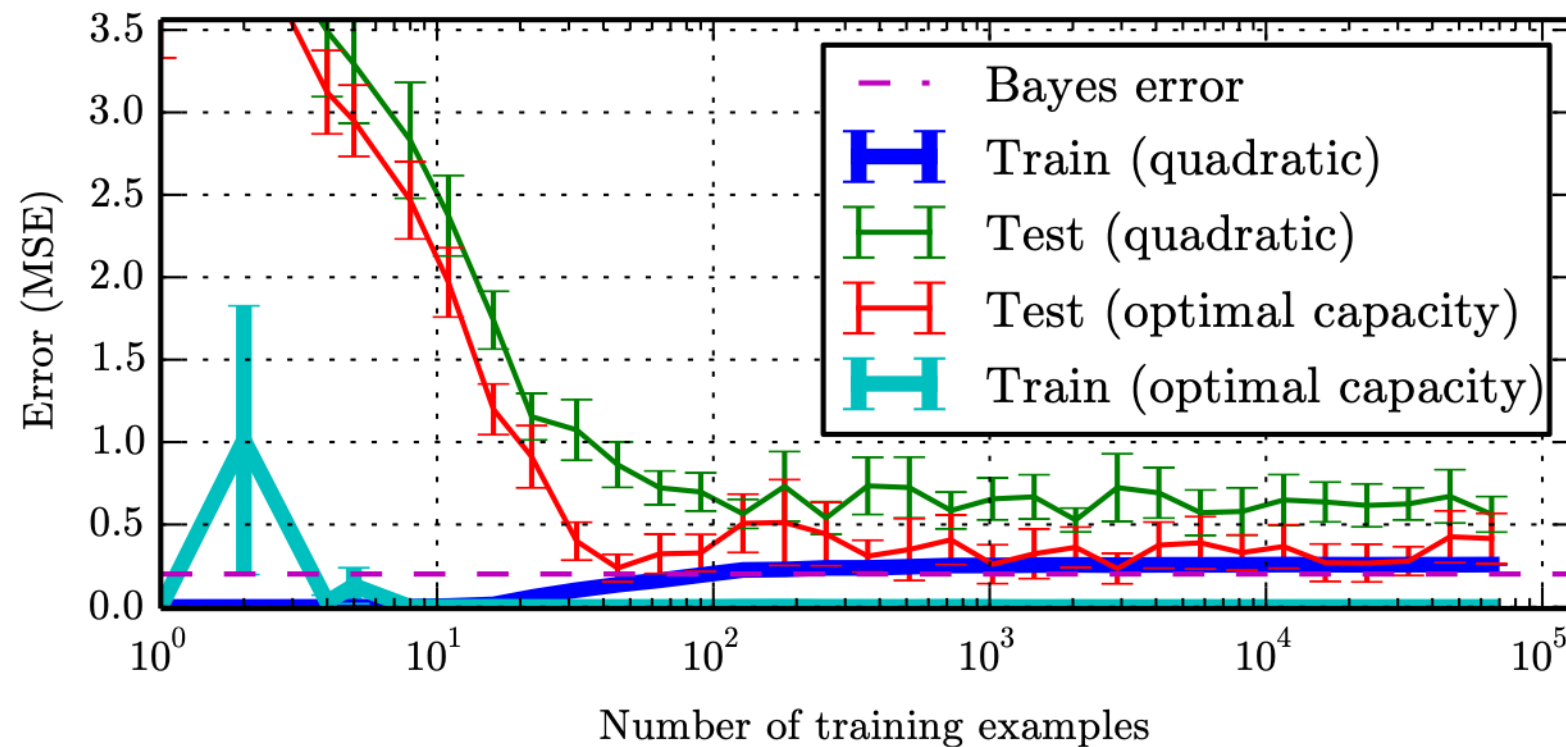
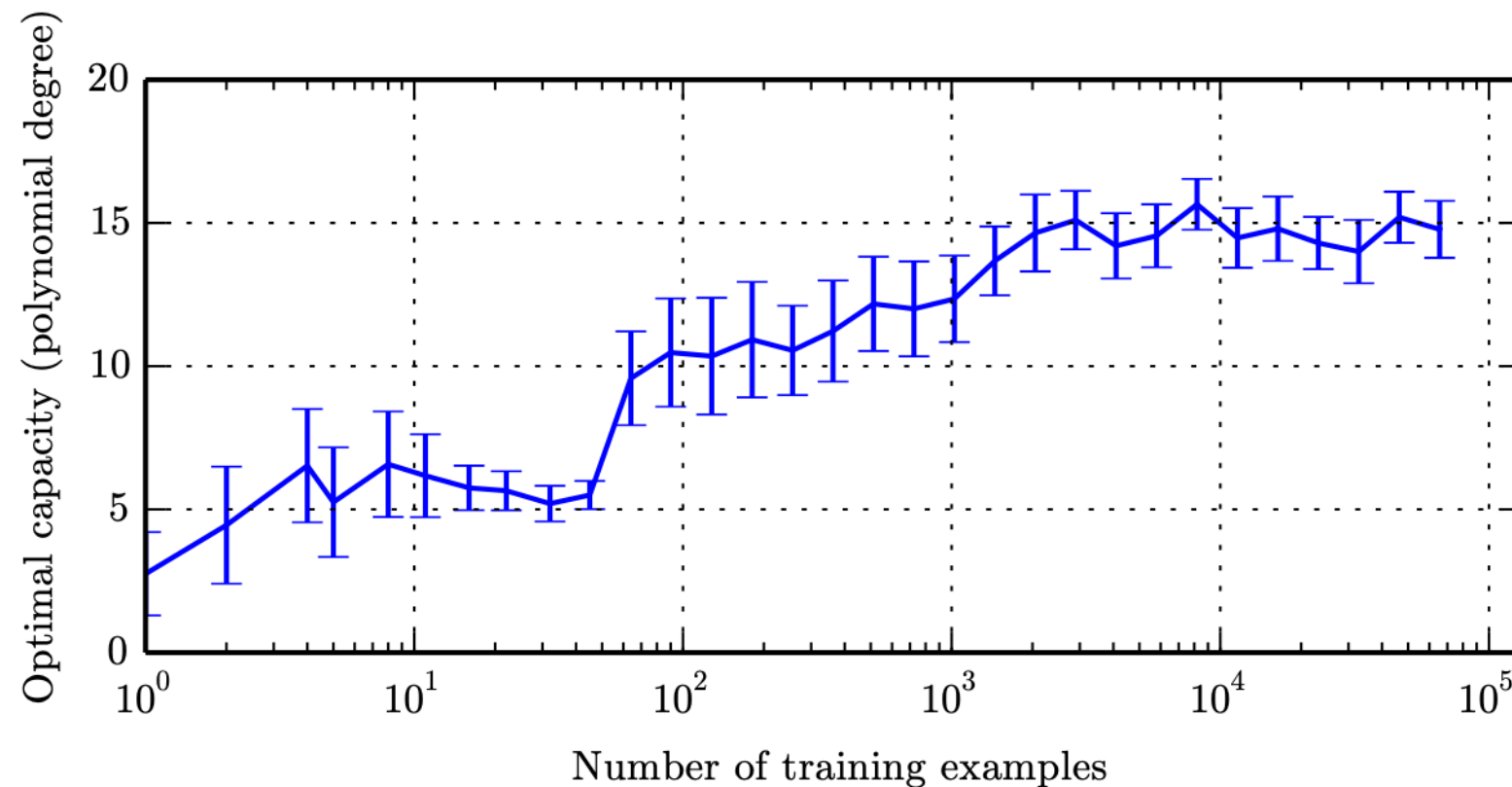


Figure 5.4



Regularization

- We can reduce model capacity by:
 1. Restricting the type and amount of functions
 2. Allowing lots of functions, but providing a preference for one solution over another, e.g. prefer solutions with smaller norm

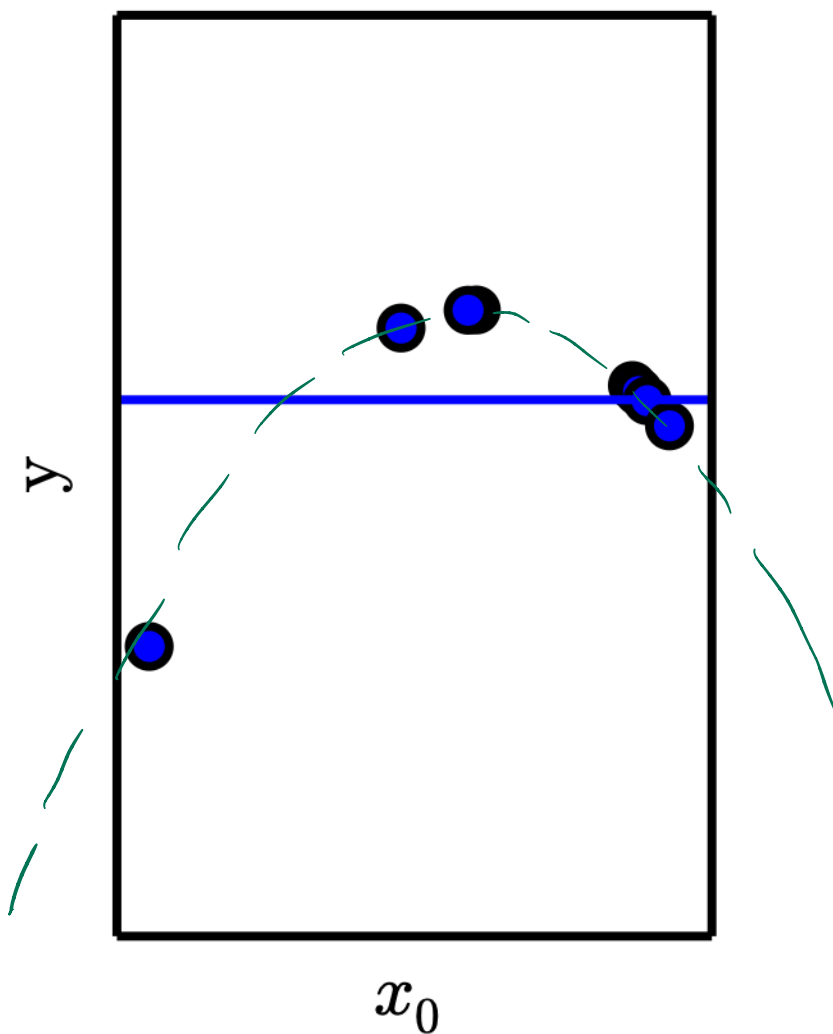
$$J(\mathbf{w}) = \text{MSE}_{\text{train}} + \lambda \mathbf{w}^\top \mathbf{w}, \quad (5.18)$$

Regularization: any modification we make to a learning algorithm that is intended to reduce its generalization error but not its training error.

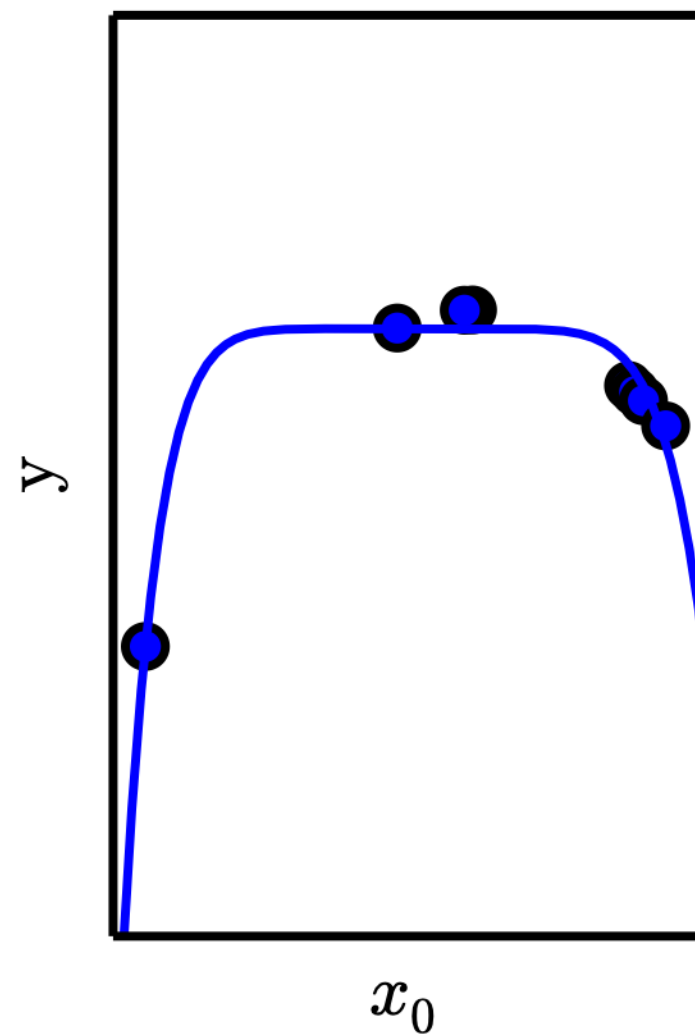
9-order polynomial

Weight Decay

Underfitting
(Excessive λ)



Appropriate weight decay
(Medium λ)



Overfitting
($\lambda \rightarrow 0$)

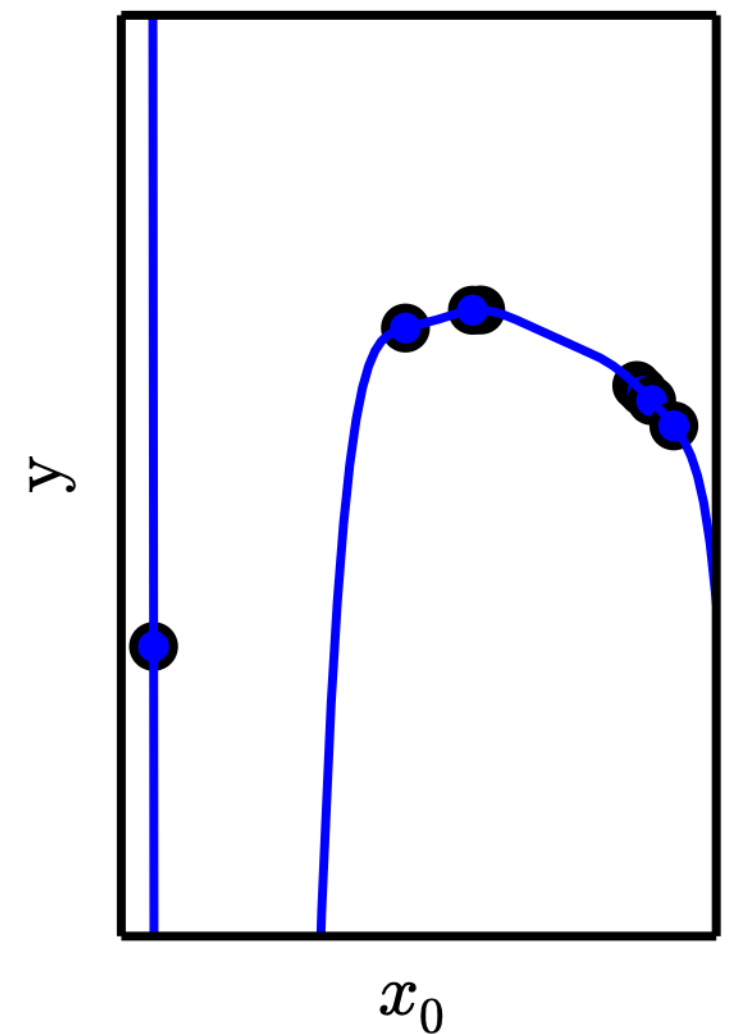


Figure 5.5

Ridge Regression

l_2 -norm penalty on "weights"
 $\text{loss} + \|w\|_2^2 \rightarrow \sum_i w_i^2$

w_i will be small

Lasso Regression

\mathcal{L}_1 - norm penalty on weights

$$\text{loss} + \|w\|_1 \quad \text{---} \quad \sum_i |w_i|$$

w will be sparse

many w_i will be zero

Hyperparameters & Cross-Validation

- e.g., λ is a hyperparameter for regularization
- Usually many other hyperparameters
- Use ***validation set*** that the training algo doesn't see
- For small test sets: k-fold cross-validation