

Lecture 5. Optimisation. Regularisation

COMP90051 Statistical Machine Learning

Semester 2, 2017
Lecturer: Andrey Kan



THE UNIVERSITY OF
MELBOURNE

This lecture

- Iterative optimisation
 - * Loss functions
 - * Coordinate descent
 - * Gradient descent
- Regularisation
 - * Model complexity
 - * Constrained modelling
 - * Bias-variance trade-off

Iterative Optimisation

A very brief summary of a few
basic optimisation methods

Supervised learning*

1. Assume a model (e.g., linear model)
 - * Denote parameters of the model as θ
 - * Model predictions are $\hat{f}(x, \theta)$
2. Choose a way to measure discrepancy between predictions and training data
 - * E.g., sum of squared residuals $\|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2$
3. Training = parameter estimation = optimisation

$$\hat{\theta} = \operatorname{argmin}_{\theta \in \Theta} L(\text{data}, \theta)$$

*This is the setup of what's called *frequentist supervised learning*. A different view on parameter estimation/training will be presented later in the subject.

Loss functions: Measuring discrepancy

- For a single training example the discrepancy between prediction and label is measured using a loss function

- Examples:

- * squared loss $l_{sq} = (y - \hat{f}(\mathbf{x}, \boldsymbol{\theta}))^2$

- * absolute loss $l_{abs} = |y - \hat{f}(\mathbf{x}, \boldsymbol{\theta})|$

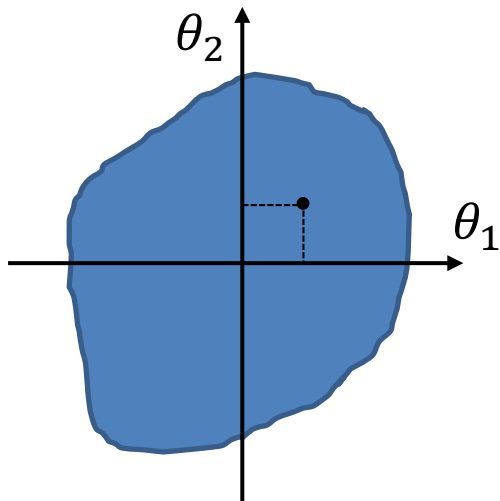
- * Perceptron loss (next lecture)

- * Hinge loss (later in the subject)

Solving optimisation problems

- Analytic (aka closed form) solution
 - * Known only in limited number of cases
 - * Use necessary condition: $\frac{\partial L}{\partial \theta_1} = \dots = \frac{\partial L}{\partial \theta_p} = 0$
- Approximate iterative solution
 1. Initialisation: choose starting guess $\boldsymbol{\theta}^{(1)}$, set $i = 1$
 2. Update: $\boldsymbol{\theta}^{(i+1)} \leftarrow \text{SomeRule}[\boldsymbol{\theta}^{(i)}]$, set $i \leftarrow i + 1$
 3. Termination: decide whether to **Stop**
 4. Go to **Step 2**
 5. **Stop**: return $\hat{\boldsymbol{\theta}} \approx \boldsymbol{\theta}^{(i)}$

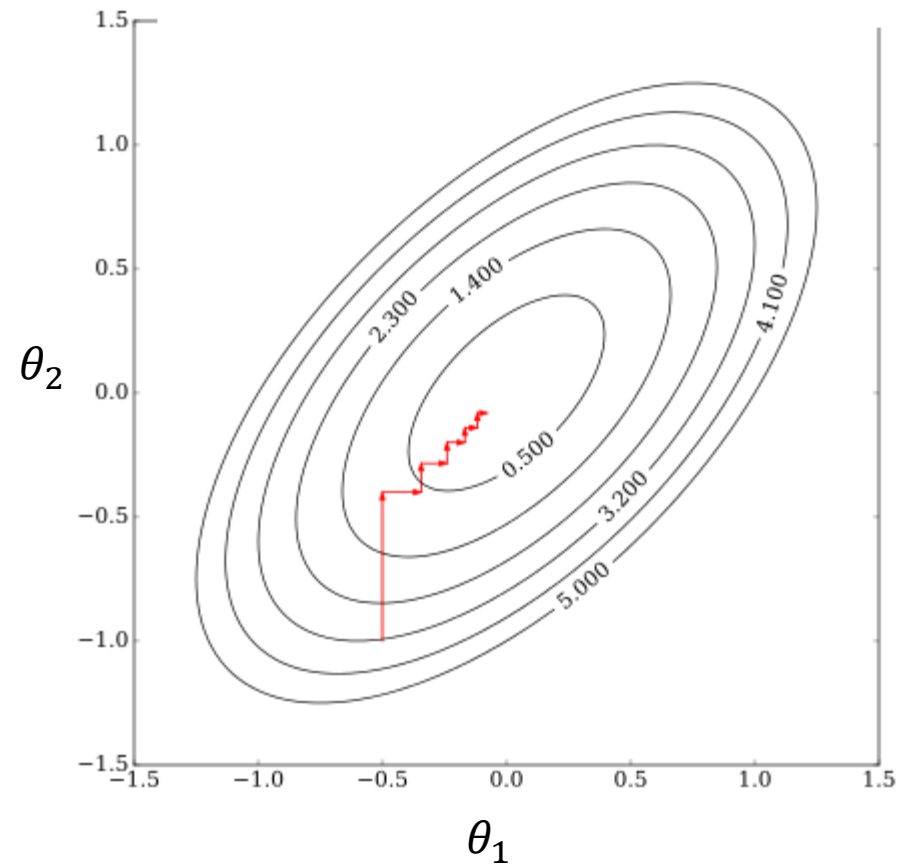
Finding the deepest point



- In this example, a model has 2 parameters
- Each location corresponds to a particular combination of parameter values
- Depth indicates discrepancy between model using those values and data

Coordinate descent

- Suppose $\boldsymbol{\theta} = [\theta_1, \dots, \theta_K]'$
- 1. Choose $\boldsymbol{\theta}^{(1)}$ and some T
- 2. For i from 1 to T^*
 1. $\boldsymbol{\theta}^{(i+1)} \leftarrow \boldsymbol{\theta}^{(i)}$
 2. For j from 1 to K
 1. Fix components of $\boldsymbol{\theta}^{(i+1)}$, except j -th component
 2. Find $\hat{\theta}_j^{(i+1)}$ that minimises $L(\theta_j^{(i+1)})$
 3. Update j -th component of $\boldsymbol{\theta}^{(i+1)}$
- 3. Return $\hat{\boldsymbol{\theta}} \approx \boldsymbol{\theta}^{(i)}$



*Other stopping criteria can be used

Gradient

- Gradient (at point θ) is defined as $\left[\frac{\partial L}{\partial \theta_1}, \dots, \frac{\partial L}{\partial \theta_p} \right]'$ computed at point θ
- One can show that gradient points to the direction of maximal change of $L(\theta)$ when departing from point θ
- Shorthand notation
 - * $\nabla L \stackrel{\text{def}}{=} \left[\frac{\partial L}{\partial \theta_1}, \dots, \frac{\partial L}{\partial \theta_p} \right]'$ computed at point θ
 - * Here ∇ is the nabla symbol

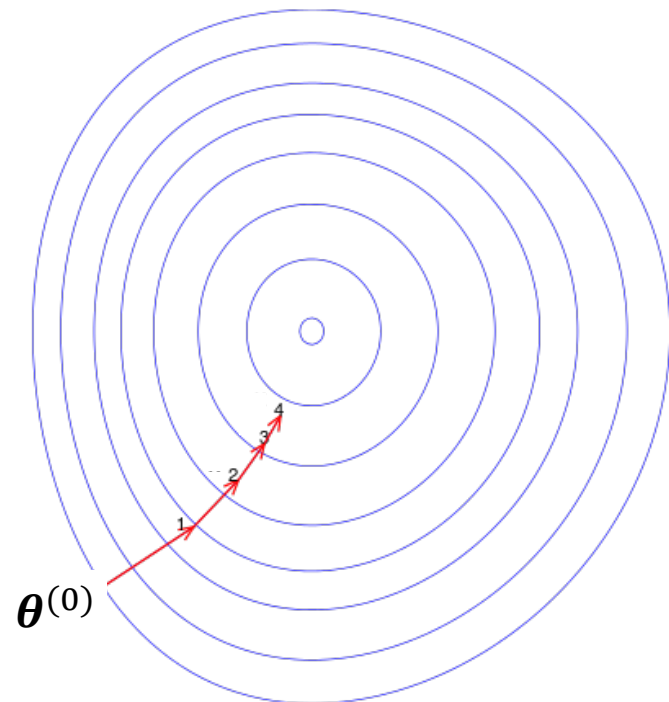


Harps, p. 984.

Gradient descent

1. Choose $\theta^{(1)}$ and some T
 2. For i from 1 to T^*
 1. $\theta^{(i+1)} = \theta^{(i)} - \eta \nabla L(\theta^{(i)})$
 3. Return $\hat{\theta} \approx \theta^{(i)}$
- Note: η is dynamically updated in each step

We assume L is differentiable



*Other stopping criteria can be used

Wikimedia Commons. Authors:
Olegalexandrov, Zerodamage

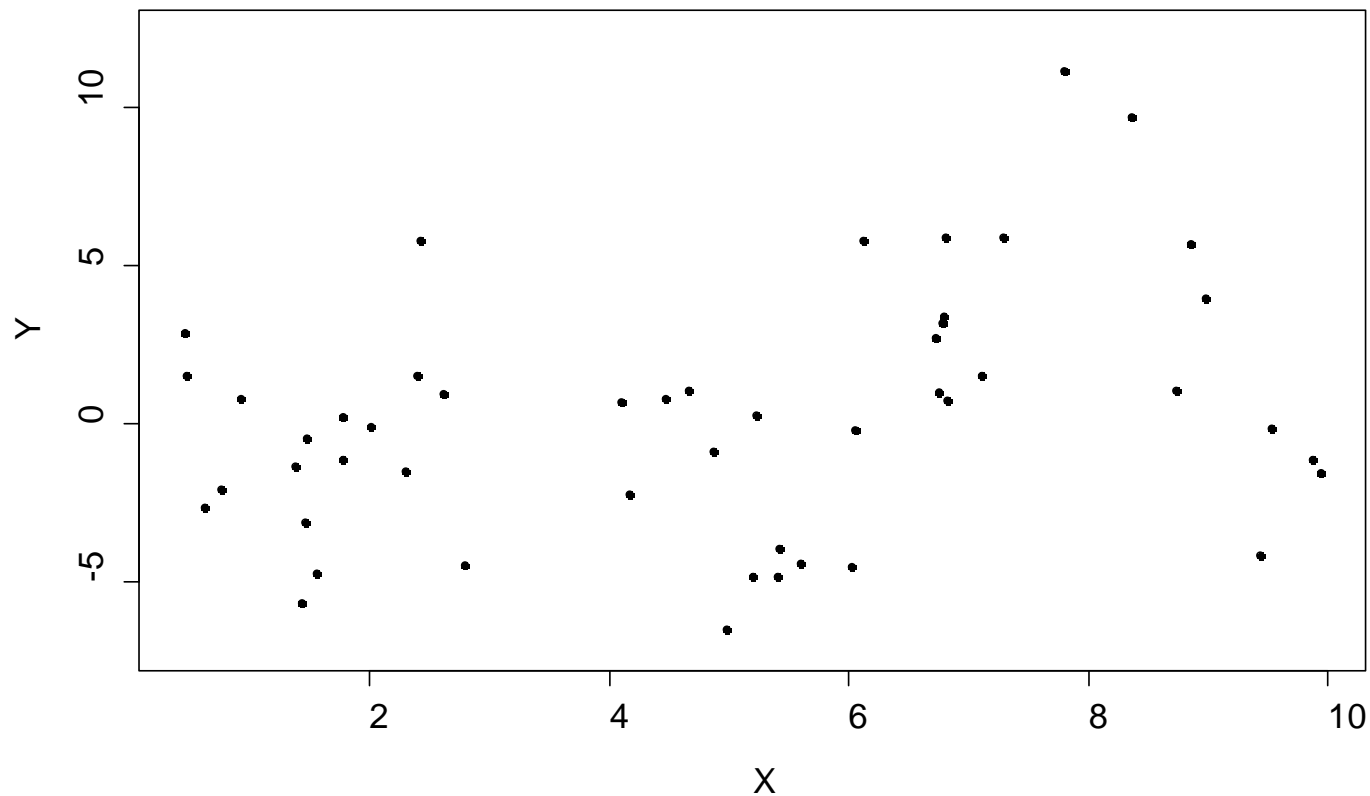
Regularisation

Process of introducing additional information in order to solve an ill-posed problem or to prevent overfitting (*Wikipedia*)

Previously: Regularisation

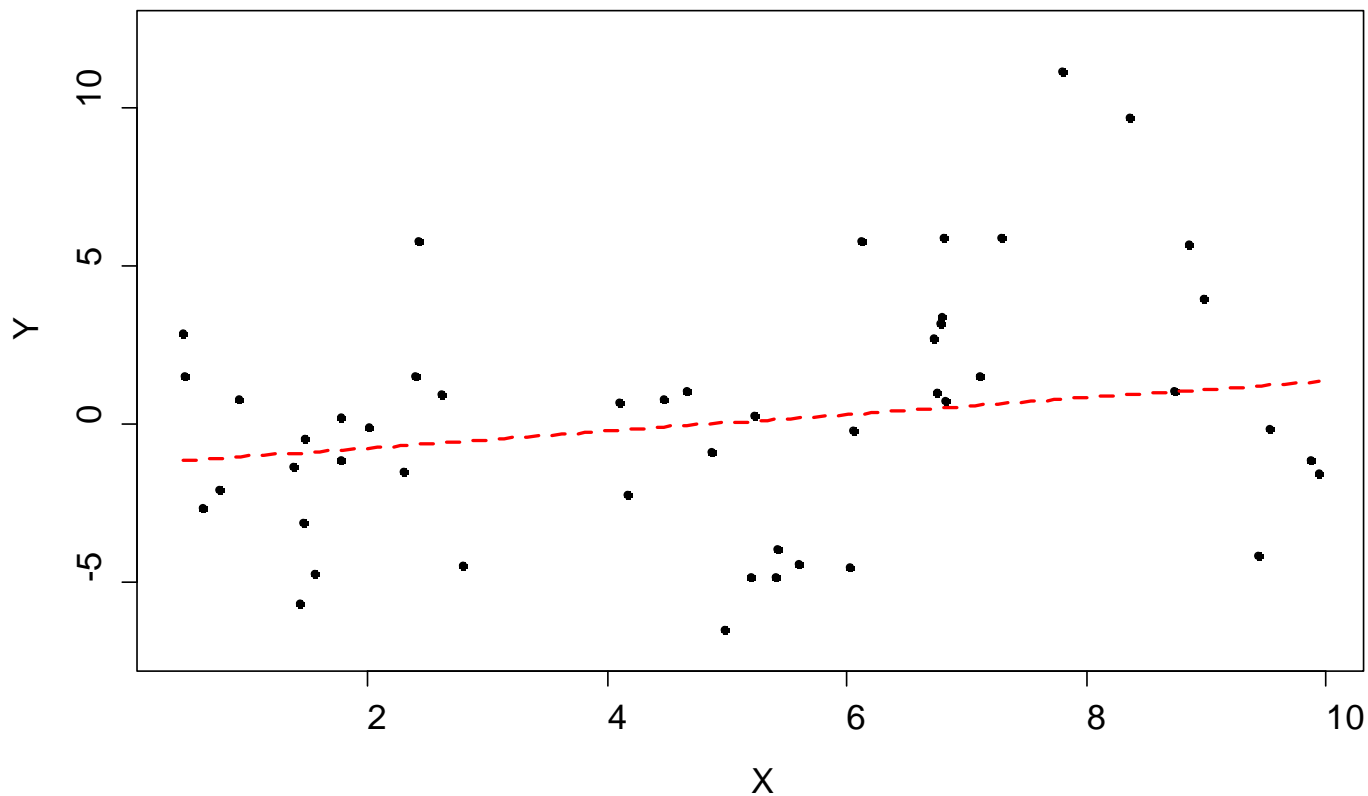
- Major technique, common in Machine Learning
- Addresses one or more of the following related problems
 - * Avoid ill-conditioning
 - * Introduce prior knowledge
 - * Constrain modelling
- This is achieved by augmenting the objective function
- Not just for linear methods.

Example regression problem



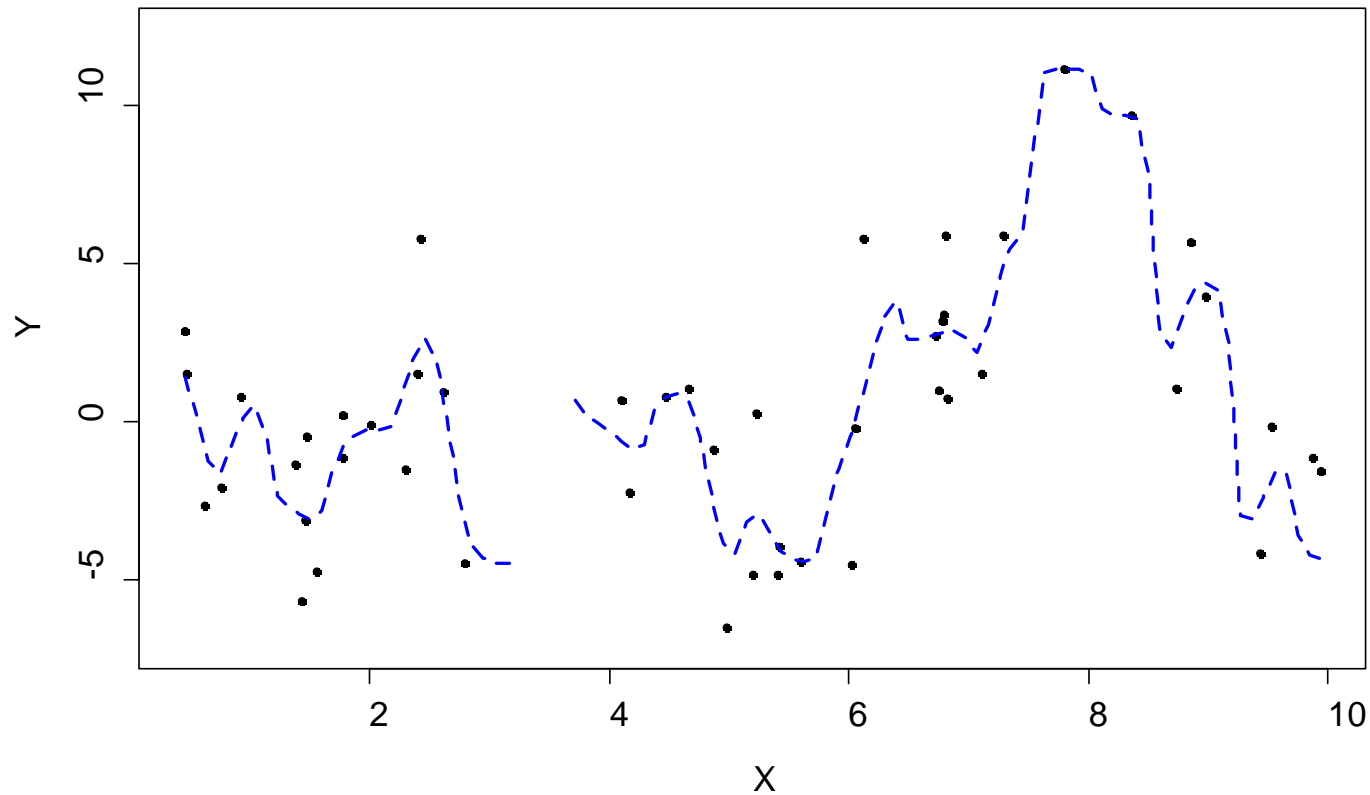
How complex a model should we use?

Underfitting (linear regression)



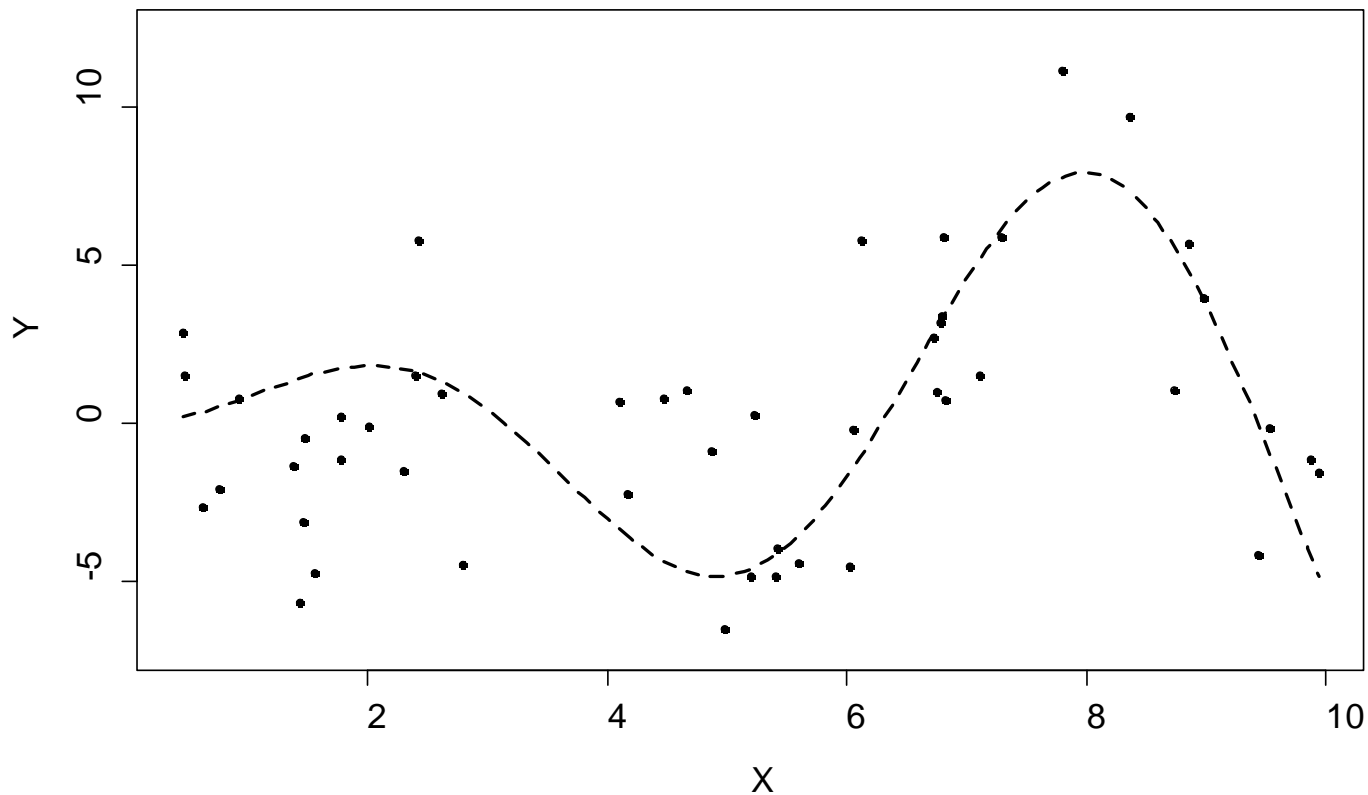
Model class Θ can be **too simple** to possibly fit true model.

Overfitting (non-parametric smoothing)



Model class Θ can be **so complex** it can fit true model + noise

Actual model ($x \sin x$)



The **right model class** Θ will sacrifice some training error, for test error.

How to “vary” model complexity

- Method 1: Explicit model selection
- Method 2: Regularisation
- Usually, method 1 can be viewed a special case of method 2

1. Explicit model selection

- Try different classes of models. Example, try polynomial models of various degree d (linear, quadratic, cubic, ...)
- Use held out validation (cross validation) to select the model
 1. Split training data into D_{train} and $D_{validate}$ sets
 2. For each degree d we have model f_d
 1. Train f_d on D_{train}
 2. Test f_d on $D_{validate}$
 3. Pick degree \hat{d} that gives the best test score
 4. Re-train model $f_{\hat{d}}$ using all data

2. Vary complexity by regularisation

- Augment the problem:

$$\hat{\boldsymbol{\theta}} = \operatorname{argmin}_{\boldsymbol{\theta} \in \Theta} (L(\text{data}, \boldsymbol{\theta}) + \lambda R(\boldsymbol{\theta}))$$

- E.g., ridge regression

$$\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w} \in W} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$$

- Note that $R(\boldsymbol{\theta})$ does not depend on data
- Use held out validation/cross validation to choose λ

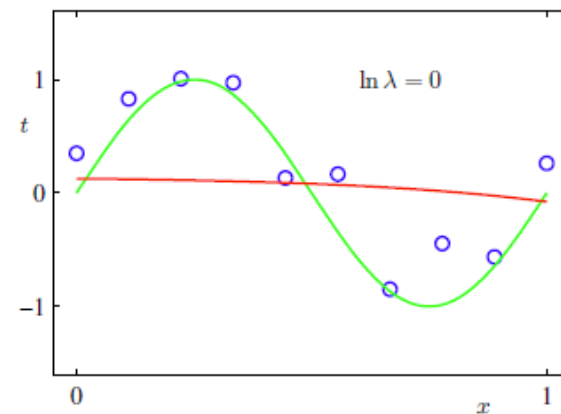
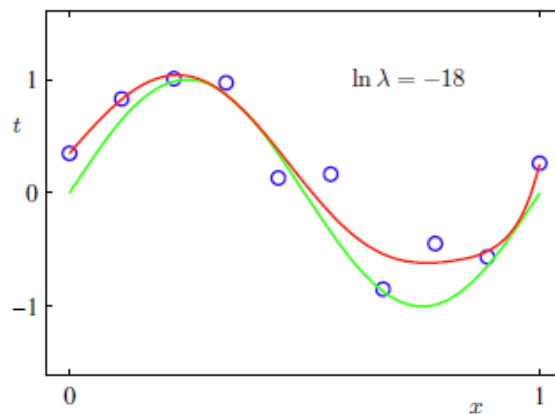
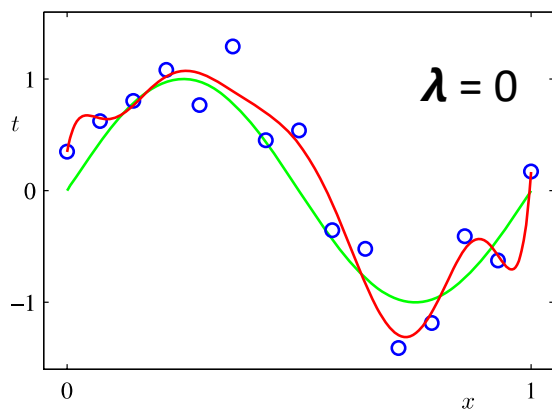
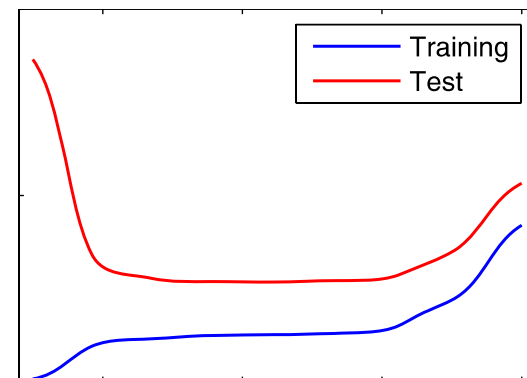
Example: polynomial regression

- 9th order polynomial regression

- * model of form

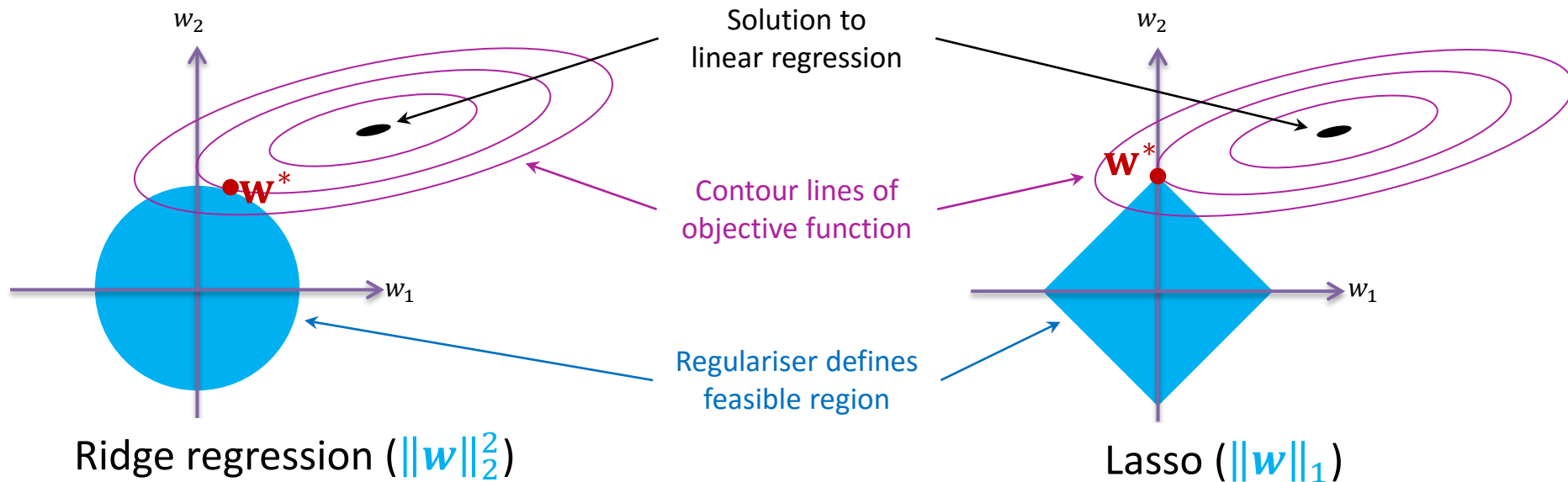
$$\hat{f} = w_0 + w_1 x + \dots + w_9 x^9$$

- * regularised with $\lambda \|\mathbf{w}\|_2^2$ term



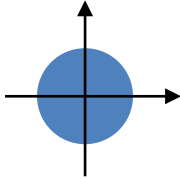
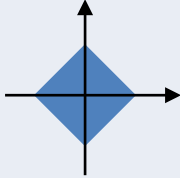
Regulariser as a constraint

- For illustrative purposes, consider a *modified problem*:
 minimise $\|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2$ subject to $\|\mathbf{w}\|_2^2 \leq \lambda$ for $\lambda > 0$



- Lasso (L1 regularisation) encourages solutions to sit on the axes
 → Some of the weights are set to zero → **Solution is sparse**

Regularised linear regression

Algorithm	Minimises	Regulariser	Solution
Linear regression	$\ \mathbf{y} - \mathbf{X}\mathbf{w}\ _2^2$	None	$(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$ (if inverse exists)
Ridge regression	$\ \mathbf{y} - \mathbf{X}\mathbf{w}\ _2^2 + \lambda\ \mathbf{w}\ _2^2$	L2 norm 	$(\mathbf{X}'\mathbf{X} + \lambda\mathbf{I})^{-1}\mathbf{X}'\mathbf{y}$
Lasso	$\ \mathbf{y} - \mathbf{X}\mathbf{w}\ _2^2 + \lambda\ \mathbf{w}\ _1$	L1 norm 	No closed-form, but solutions are sparse and suitable for high-dim data

Bias-variance trade-off

Analysis of relations between
train error, test error and
model complexity

Assessing generalisation capacity

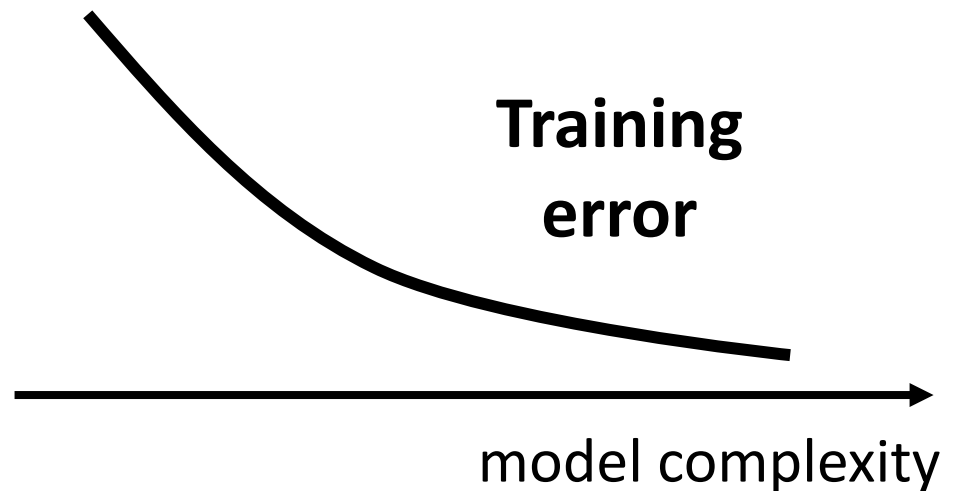
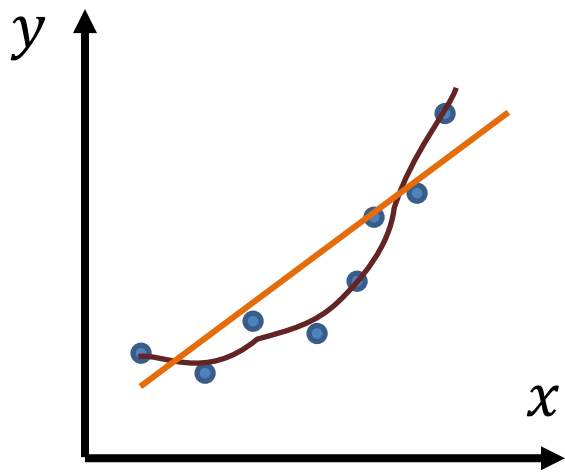
- Supervised learning: train the model on existing data, then make predictions on new data
 - * Generalisation capacity of the model is an important consideration
- Training the model: minimisation of training error
- Generalisation capacity is captured by the test error
- Model complexity is a major factor that influences the ability of the model to generalise
- In this section, our aim is to explore relations between training error, test error and model complexity

Training error

- Suppose training data is $D = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$
- Let $\hat{f}(\mathbf{x})$ denote predicted value for \mathbf{x} from the model trained on D
- Let $l(y_i, \hat{f}(\mathbf{x}_i))$ denote loss on the i -th training example
- Training error: $\frac{1}{n} \sum_{i=1}^n l(y_i, \hat{f}(\mathbf{x}_i))$

Training error and model complexity

- More complex model \rightarrow training error goes down
- Finite number of points \rightarrow usually can reduce training error to 0 (is it always possible?)



Test error

- Assume $\mathcal{Y} = h(\mathbf{x}_0) + \varepsilon$
 - * \mathbf{x}_0 is a fixed instance
 - * $h(\mathbf{x})$ is an unknown true function
 - * ε is noise, $\mathbb{E}[\varepsilon] = 0$ and $\text{Var}[\varepsilon] = \sigma^2$
- Treat training data as a random variable \mathcal{D}
 - * Draw $D \sim \mathcal{D}$, train model on D , make predictions
 - * Prediction $\hat{f}(\mathbf{x}_0)$ is a random variable, despite \mathbf{x}_0 fixed
- Test error for \mathbf{x}_0 : $\mathbb{E} \left[l \left(\mathcal{Y}, \hat{f}(\mathbf{x}_0) \right) \right]$
 - * The expectation is taken with respect to \mathcal{D} and ε
 - * \mathcal{D} and ε are assumed to be independent

Bias-variance decomposition

- For the following analysis, we consider squared loss as an important special case

$$l(y, \hat{f}(x_0)) = (y - \hat{f}(x_0))^2$$

- Lemma: Bias-Variance Decomposition

$$\mathbb{E} \left[l(y, \hat{f}(x_0)) \right] = (\mathbb{E}[y] - \mathbb{E}[\hat{f}])^2 + \text{Var}[\hat{f}] + \text{Var}[y]$$

test error
for x_0

(bias)²

variance

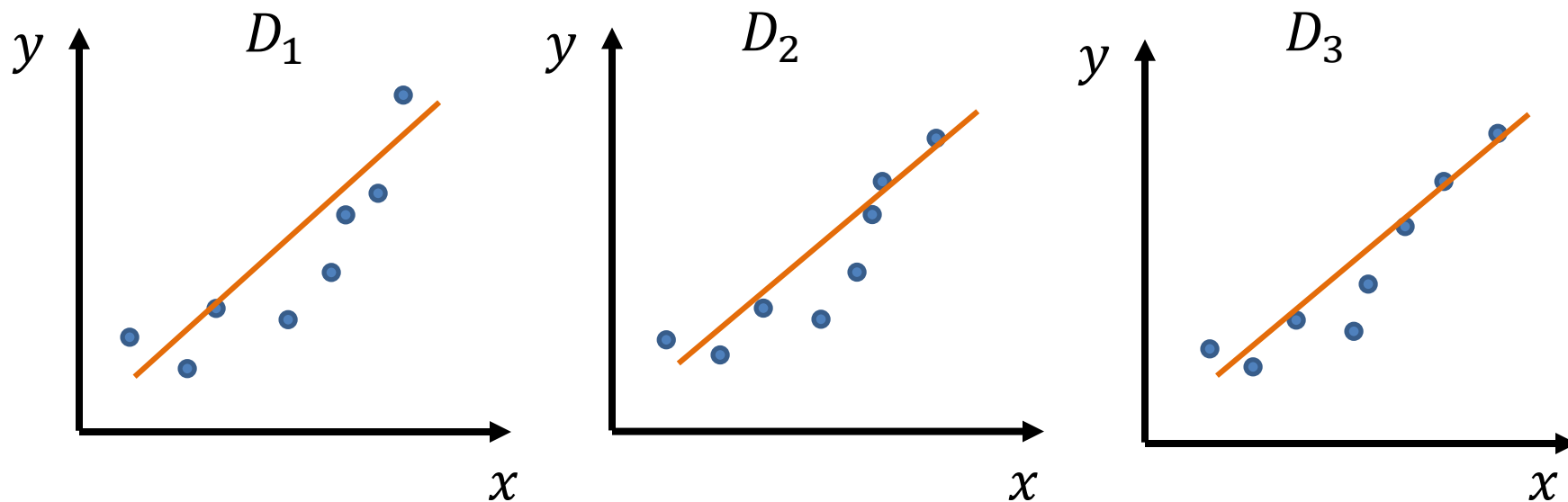
irreducible
error

Decomposition proof sketch

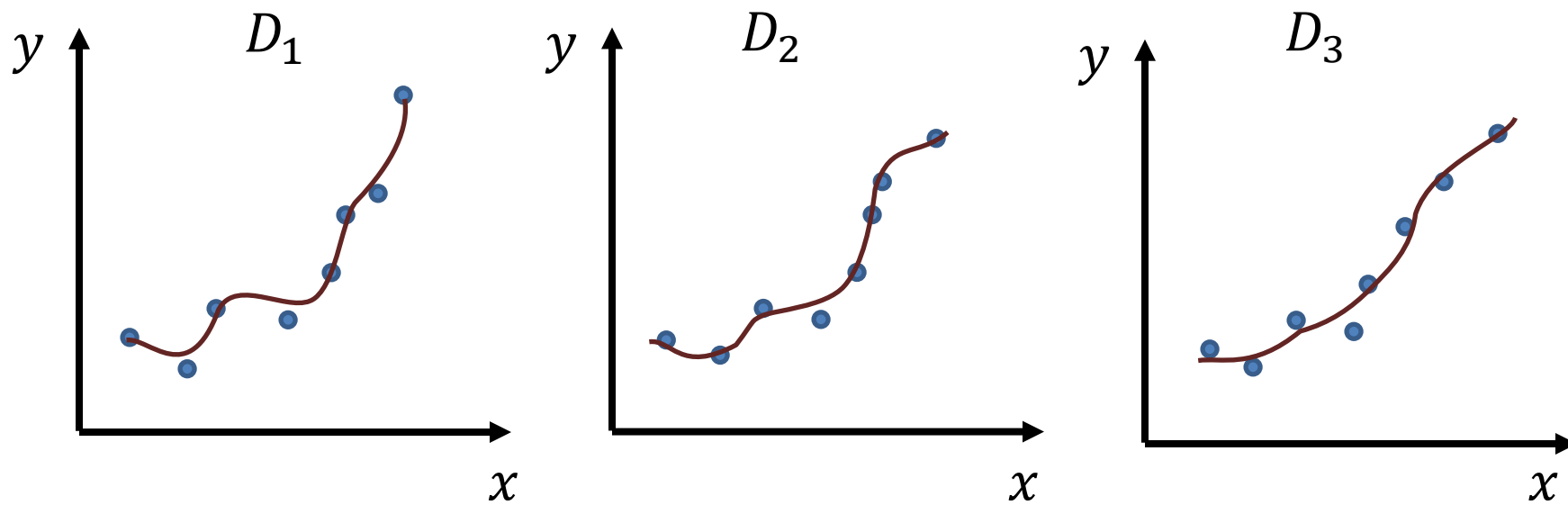
- Here (\mathbf{x}) is omitted to de-clutter notation
- $\mathbb{E}[(y - \hat{f})^2] = \mathbb{E}[y^2 + \hat{f}^2 - 2y\hat{f}]$
- $= \mathbb{E}[y^2] + \mathbb{E}[\hat{f}^2] - \mathbb{E}[2y\hat{f}]$
- $= \text{Var}[y] + \mathbb{E}[y]^2 + \text{Var}[\hat{f}] + \mathbb{E}[\hat{f}]^2 - 2\mathbb{E}[y]\mathbb{E}[\hat{f}]$
- $= \text{Var}[y] + \text{Var}[\hat{f}] + (\mathbb{E}[y]^2 - 2\mathbb{E}[y]\mathbb{E}[\hat{f}] + \mathbb{E}[\hat{f}]^2)$
- $= \text{Var}[y] + \text{Var}[\hat{f}] + (\mathbb{E}[y] - \mathbb{E}[\hat{f}])^2$

* Green slides are non-examinable

Training data as a random variable

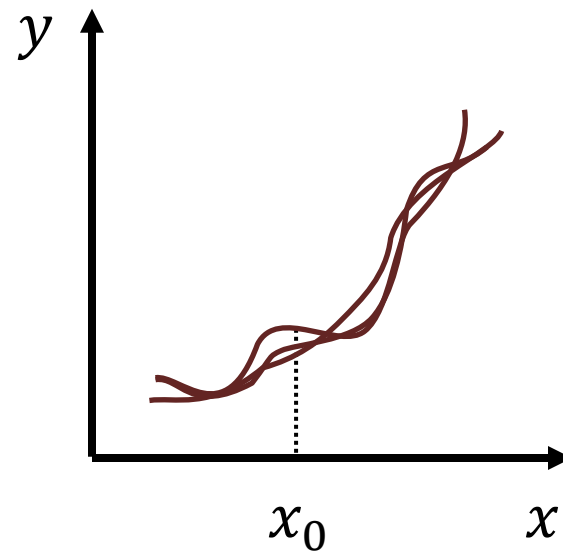
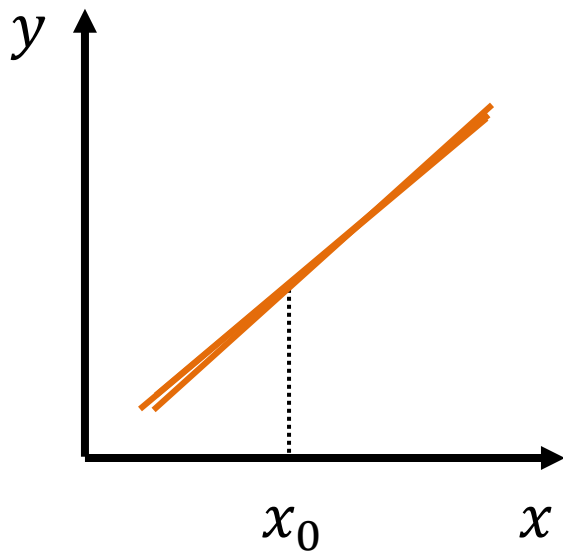


Training data as a random variable



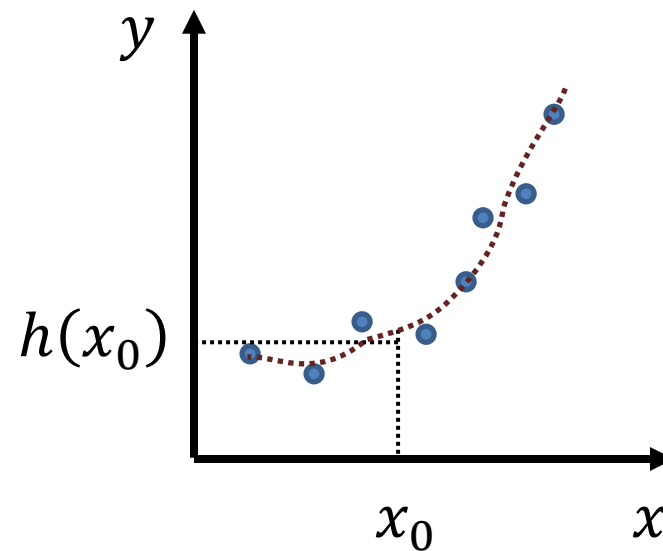
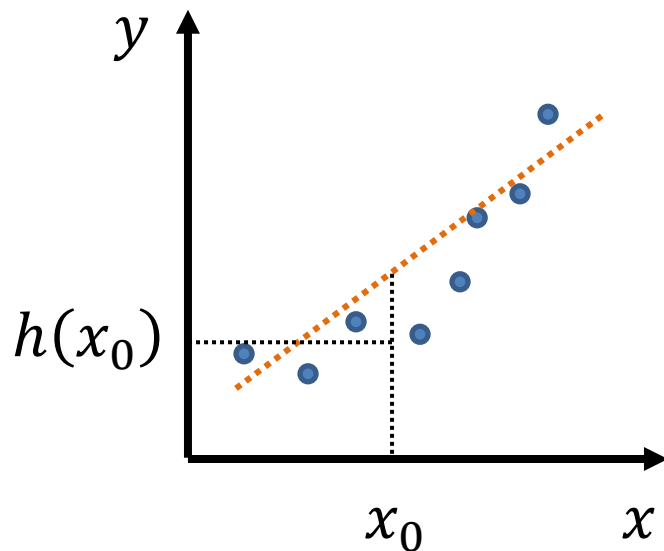
Model complexity and variance

- simple model \rightarrow low variance
- complex model \rightarrow high variance



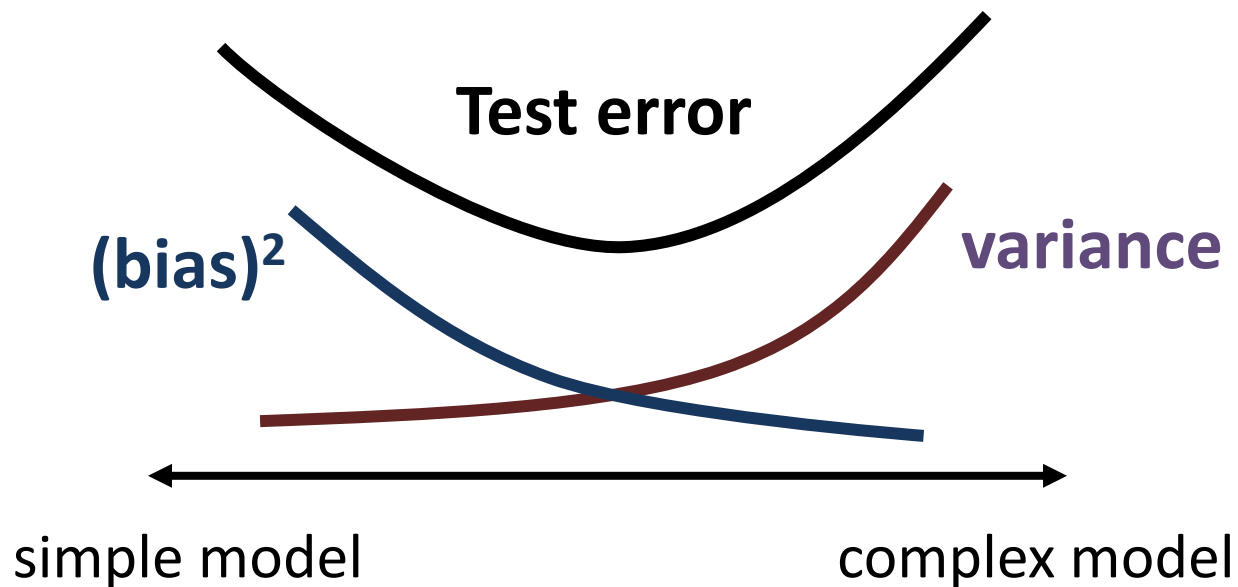
Model complexity and bias

- simple model \rightarrow high bias
- complex model \rightarrow low bias

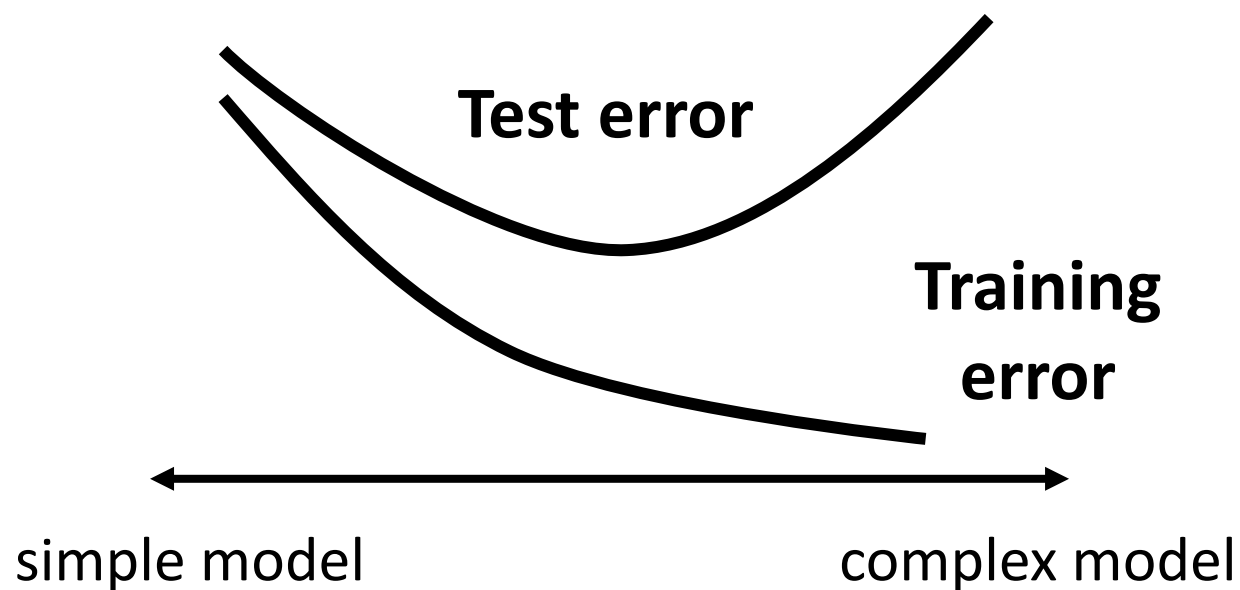


Bias-variance trade-off

- simple model \rightarrow high bias, low variance
- complex model \rightarrow low bias, high variance



Test error and training error



This lecture

- Iterative optimisation
 - * Loss functions
 - * Coordinate descent
 - * Gradient descent
- Regularisation
 - * Model complexity
 - * Constrained modelling
 - * Bias-variance trade-off