# Lecture 5. Regularisation

COMP90051 Statistical Machine Learning

Semester 2, 2019
Lecturer:  Ben Rubinstein

THE UNIVERSITY OF
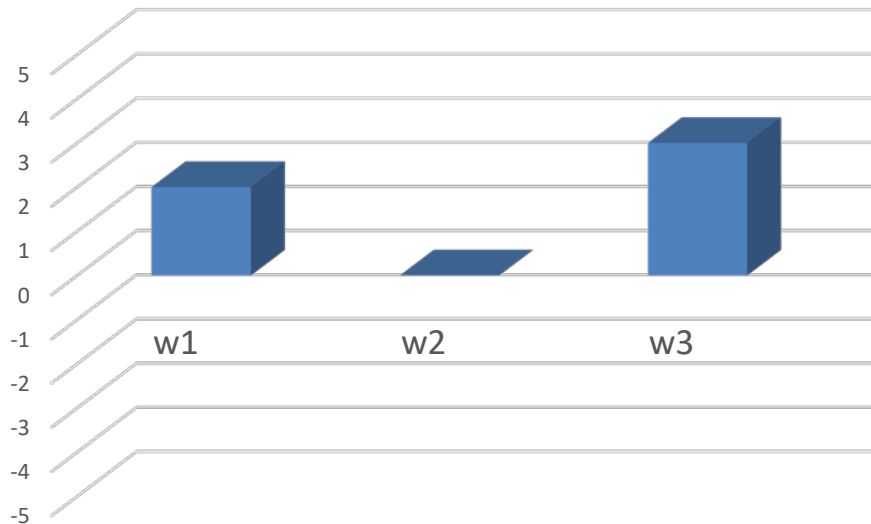MELBOURNE

# This lecture: Regularisation

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

• Major technique & theme, throughout ML

• Addresses one or more of the following related problems
  * Avoids ill-conditioning (a computational problem)
  * Avoids overfitting (a statistical problem)
  * Introduce prior knowledge into modelling

• This is achieved by augmenting the objective function

• In this lecture: we cover the first two aspects. We will cover more of regularisation throughout the subject

# Example 1: Feature importance

- Linear model on three features
  - $X$ is matrix on $n = 4$ instances (rows)
  - Model: $y = w_1 x_1 + w_2 x_2 + w_3 x_3 + w_0$
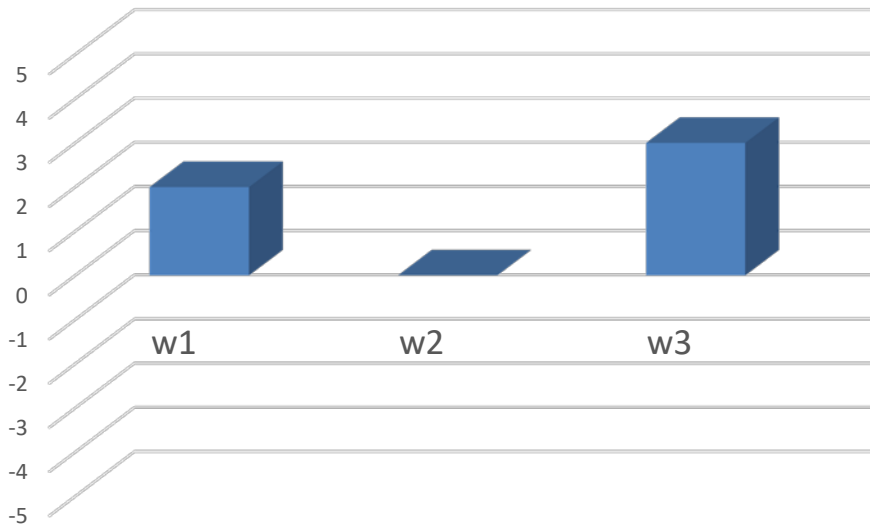


**Question: Which feature is more important?**

# Question: Which feature is more important?

1

2

3

I don't
know

# Example 1: Feature importance

- Linear model on three features
  - $X$ is matrix on $n = 4$ instances (rows)
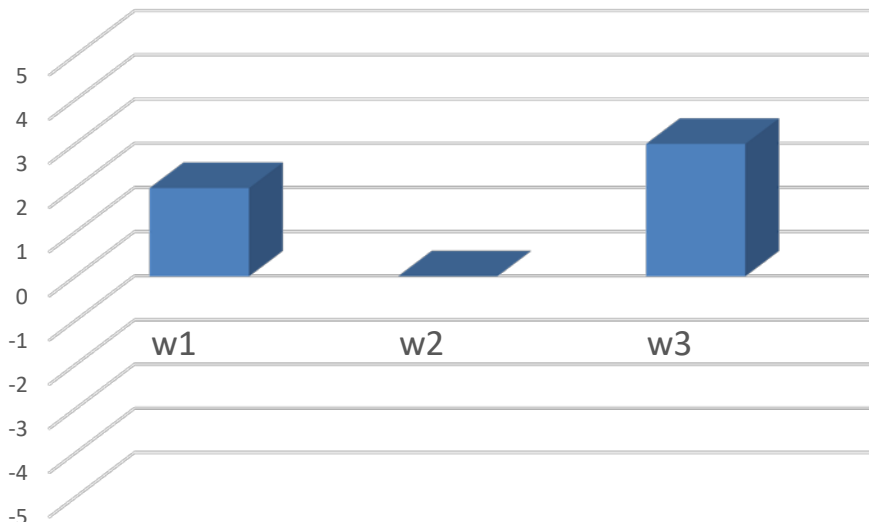  - Model: $y = w_1 x_1 + w_2 x_2 + w_3 x_3 + w_0$

# Example 1: Irrelevant features

- Linear model on three features, first two same
  - $X$ is matrix on $n = 4$ instances (rows)
  - Model: $y = w_1 x_1 + w_2 x_2 + w_3 x_3 + w_0$
  - First two columns of $X$ identical
  - Feature 2 (or 1) is **irrelevant**

| 3 | 3 | 7 |
|---|---|---|
| 6 | 6 | 9 |
| 21 | 21 | 79 |
| 34 | 34 | 2 |



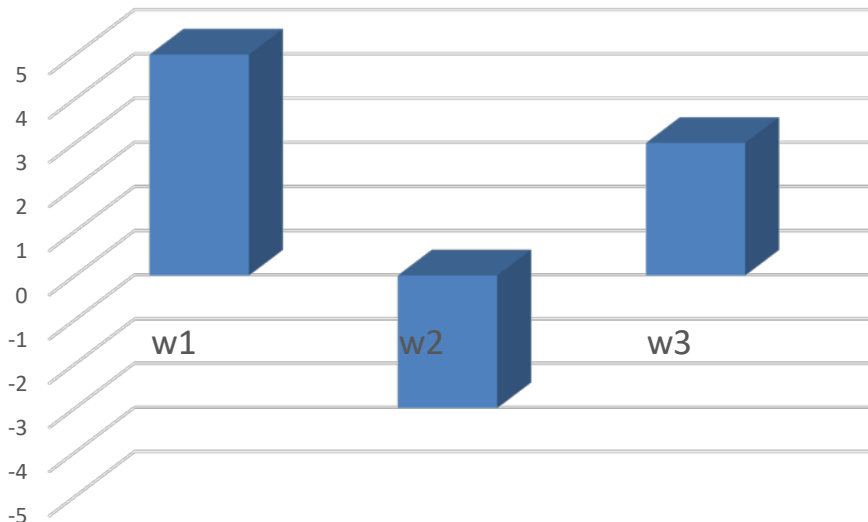- Effect of perturbations on model predictions?
  - Add $\Delta$ to $w_1$
  - Subtract $\Delta$ from $w_2$

6

# Example 1: Irrelevant features

- Linear model on three features, first two same
  - $X$ is matrix on $n = 4$ instances (rows)
  - Model: $y = w_1 x_1 + w_2 x_2 + w_3 x_3 + w_0$
  - First two columns of $X$ identical
  - Feature 2 (or 1) is **irrelevant**

| | | |
|---|---|---|
| 3 | 3 | 7 |
| 6 | 6 | 9 |
| 21 | 21 | 79 |
| 34 | 34 | 2 |

- Effect of perturbations on model predictions?
  - Add $\Delta$ to $w_1$
  - Subtract $\Delta$ from $w_2$

7

# Question: Which feature is more important?

1
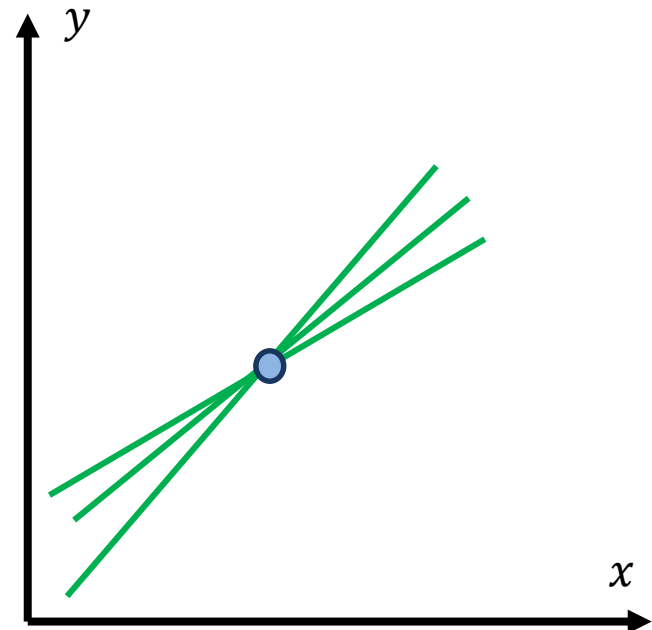
2

3

I don't
know

# Problems with irrelevant features

- In example, suppose $[\widehat{w}_0, \widehat{w}_1, \widehat{w}_2, \widehat{w}_3]'$ is "optimal"

- For any $\delta$ new $[\widehat{w}_0, \widehat{w}_1 + \delta, \widehat{w}_2 - \delta, \widehat{w}_3]'$ get

  * *Same* predictions!
  * *Same* sum of squared errors!

- Problems this highlights

  * The solution is not unique
  * Lack of interpretability
  * Optimising to learn parameters is ill-posed problem

# Irrelevant (co-linear) features in general

- Extreme case: features complete clones

- For linear models, more generally
  * Feature $X_{\cdot j}$ is irrelevant if
  * $X_{\cdot j}$ is a linear combination of other columns

  $$X_{\cdot j} = \sum_{l \neq j} \alpha_l X_{\cdot l}$$

  … for some scalars $\alpha_l$. Also called multicollinearity
  * Equivalently: Some eigenvalue of $X'X$ is zero

- Even *near*-irrelevance/colinearity can be problematic
  * V small eigenvalues of $X'X$

- Not just a pathological extreme; *easy to happen!*

$X_{\cdot j}$ denotes the $j$-th column of $X$
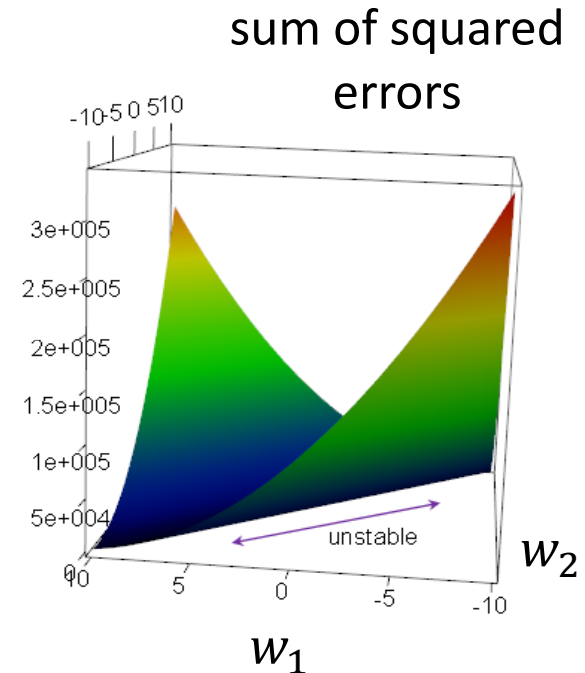
# Example 2: Lack of data

- Extreme example:
  - ∗ Model has two parameters (slope and intercept)
  - ∗ Only one data point

- Underdetermined system

# Ill-posed problems

- In both examples, finding the best parameters becomes an ill-posed problem

- This means that the problem solution is not defined
  * In our case $w_1$ and $w_2$ cannot be uniquely identified

- Remember normal equations solution of linear regression:
$$\widehat{w} = (X'X)^{-1}X'y$$

- With irrelevant/multicolinear features, matrix $X'X$ has no inverse

sum of squared errors

convex, but not strictly convex

12

# Re-conditioning the problem

- Regularisation: introduce an additional condition into the system

- The original problem is to minimise $\|\boldsymbol{y} - \boldsymbol{Xw}\|_2^2$

- The regularised problem is to minimise

$$\|\boldsymbol{y} - \boldsymbol{Xw}\|_2^2 + \lambda\|\boldsymbol{w}\|_2^2 \text{ for } \lambda > 0$$

- The solution is now
$$\widehat{\boldsymbol{w}} = (\boldsymbol{X}'\boldsymbol{X} + \lambda\boldsymbol{I})^{-1}\boldsymbol{X}'\boldsymbol{y}$$

- This formation is called ridge regression
  * Turns the ridge into a peak
  * Adds $\lambda$ to eigenvalues of $\boldsymbol{X}'\boldsymbol{X}$: makes invertible

sum of squared errors



stable

$w_2$

$w_1$

strictly convex

# Regulariser as a prior

- Without regularisation, parameters found based entirely on the information contained in the training set $X$
  * Regularisation introduces additional information

- Recall our probabilistic model $Y = x'w + \varepsilon$
  * Here $Y$ and $\varepsilon$ are random variables, where $\varepsilon$ denotes noise

- Now suppose that $w$ is also a random variable (denoted as $W$) with a Normal prior distribution
$$W \sim \mathcal{N}(0, 1/\lambda)$$
  * I.e. we expect small weights and that no one feature dominates
  * Is this always appropriate? E.g. data centring and scaling
  * We could encode much more elaborate problem knowledge

# Computing posterior using Bayes rule

- The prior is then used to compute the posterior



$$p(\boldsymbol{w}|\boldsymbol{X}, \boldsymbol{y}) = \frac{p(\boldsymbol{y}|\boldsymbol{X}, \boldsymbol{w})p(\boldsymbol{w})}{p(\boldsymbol{y}|\boldsymbol{X})}$$

posterior, likelihood, prior, marginal likelihood

- Instead of maximum likelihood (MLE), take *maximum a posteriori* estimate (MAP)

- Apply log trick, so that
$$\log(posterior) = \log(likelihood) + \log(prior) - \log(marg)$$

- Arrive at the problem of minimising
$$\|\boldsymbol{y} - \boldsymbol{X}\boldsymbol{w}\|_2^2 + \lambda\|\boldsymbol{w}\|_2^2$$

this term doesn't affect optimisation
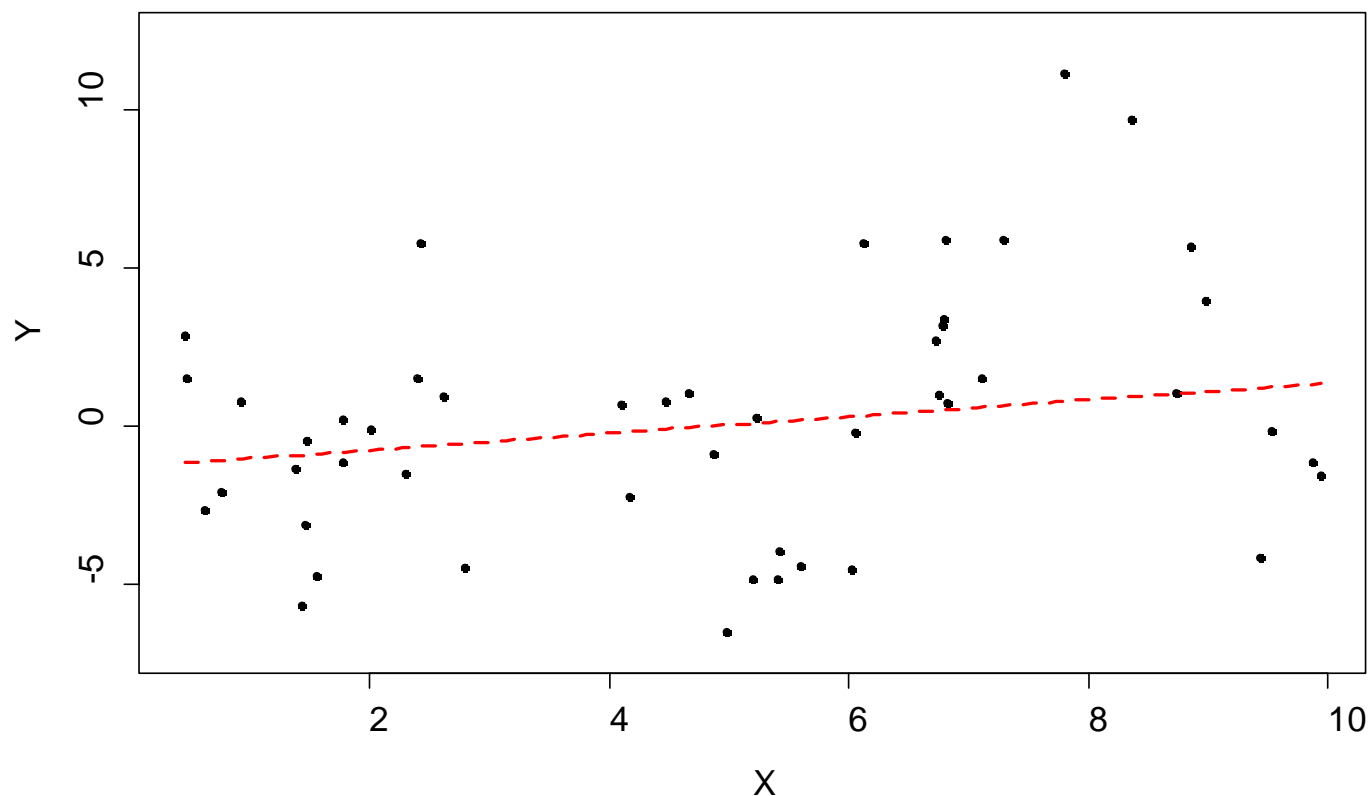
15

# Regularisation in Non-Linear Models

*Model selection in ML*
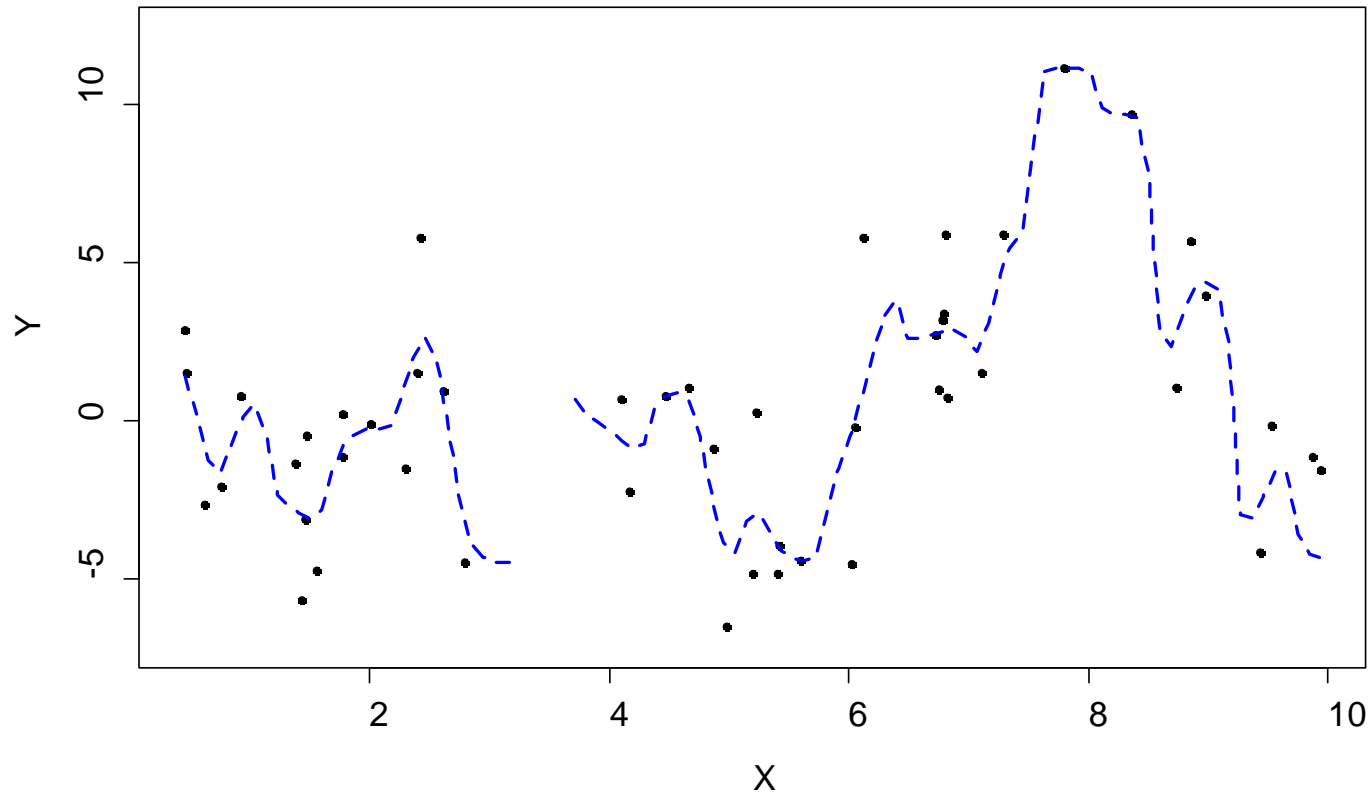
# Example regression problem



**How complex** a model should we use?

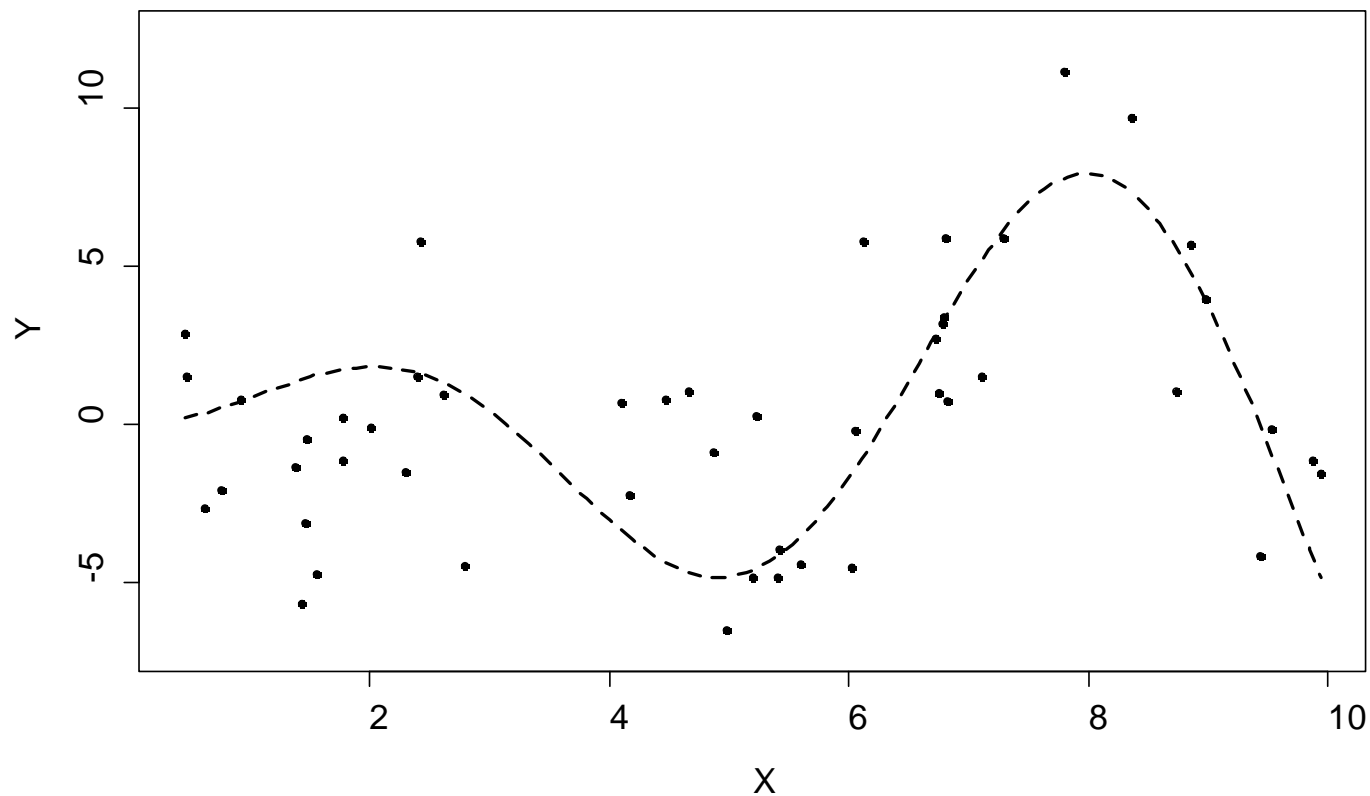# Underfitting (linear regression)



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**

19

# Actual model ($x \sin x$)



The **right model class** $\Theta$ 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 <u>held out validation</u> (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:

$$\widehat{\boldsymbol{\theta}} \in \underset{\boldsymbol{\theta} \in \Theta}{\text{argmin}} \big( L(data, \boldsymbol{\theta}) + \lambda R(\boldsymbol{\theta}) \big)$$

- E.g., ridge regression

$$\widehat{\boldsymbol{w}} \in \underset{\boldsymbol{w} \in W}{\text{argmin}} \|\boldsymbol{y} - \boldsymbol{Xw}\|_2^2 + \lambda \|\boldsymbol{w}\|_2^2$$

- Note that regulariser $R(\boldsymbol{\theta})$ does not depend on data

- Use held out validation/cross validation to choose $\lambda$

# Example: Polynomial regression
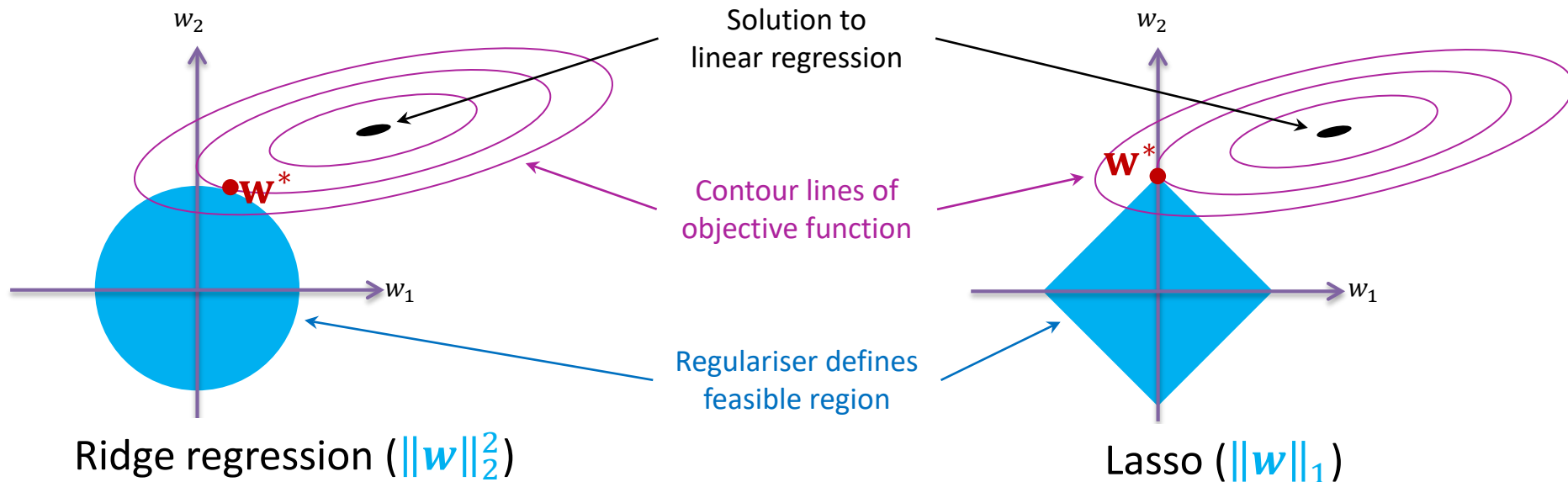
- 9$^{th}$-order polynomial regression
  * model of form
    $$\hat{f} = w_0 + w_1 x + \ldots + w_9 x^9$$
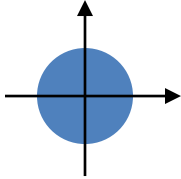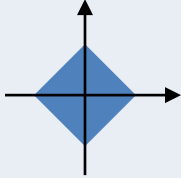  * regularised with $\lambda \|\boldsymbol{w}\|_2^2$ term





$\boldsymbol{\lambda} = 0$

$\ln \lambda = -18$

$\ln \lambda = 0$

# Regulariser as a constraint

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



Ridge regression ($\|\boldsymbol{w}\|_2^2$)

Lasso ($\|\boldsymbol{w}\|_1$)

- Lasso (L$_1$ 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 | $\|y - Xw\|_2^2$ | None | $(X'X)^{-1}X'y$ (if inverse exists) |
| Ridge regression | $\|y - Xw\|_2^2 + \lambda\|w\|_2^2$ | $L_2$ norm | $(X'X + \lambda I)^{-1}X'y$ |
| Lasso | $\|y - Xw\|_2^2 + \lambda\|w\|_1$ | $L_1$ norm | No closed-form, but solutions are sparse and suitable for high-dim data |

*Gaussian error* (handwritten annotation)

*Laplace prior* (handwritten annotation)
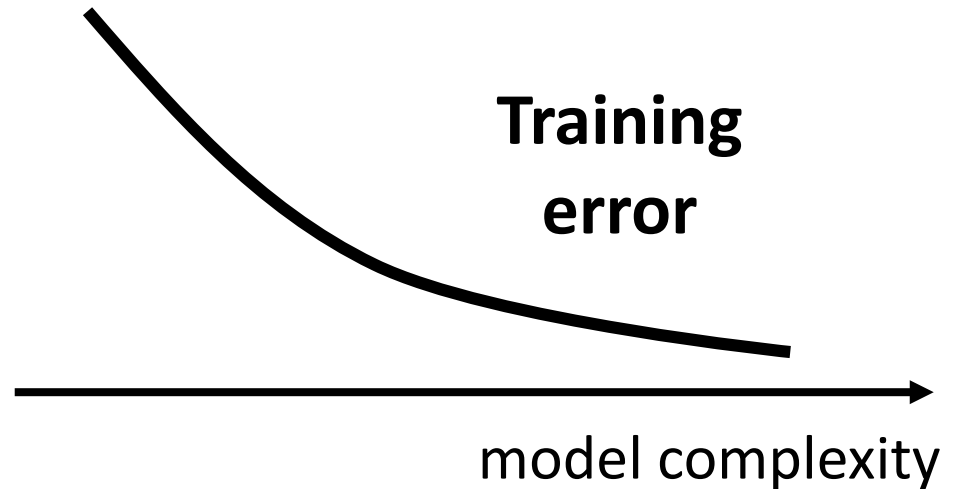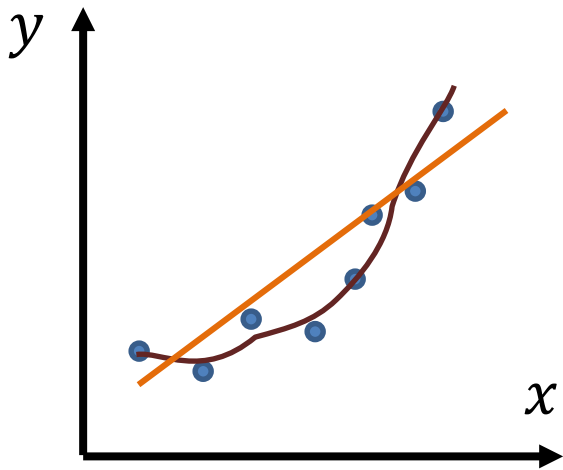
# Bias-variance trade-off

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 <u>new data</u>

- Training the model: ERM / minimisation of <u>training error</u>

- Generalisation capacity is captured by risk / <u>test error</u>

- <u>Model complexity</u> 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 and model complexity

- More complex model → training error goes down

- Finite number of points → usually can reduce training error to 0 (is it always possible?)

# (Another) Bias-variance decomposition

- Squared loss for supervised-regression predictions

$$l\left(Y, \hat{f}(\boldsymbol{X}_0)\right) = \left(Y - \hat{f}(\boldsymbol{X}_0)\right)^2$$

- Lemma: Bias-variance decomposition

$$\mathbb{E}\left[l\left(Y, \hat{f}(\boldsymbol{X}_0)\right)\right] = \left(\mathbb{E}[Y] - \mathbb{E}[\hat{f}]\right)^2 + Var[\hat{f}] + Var[Y]$$

**Risk /
test error
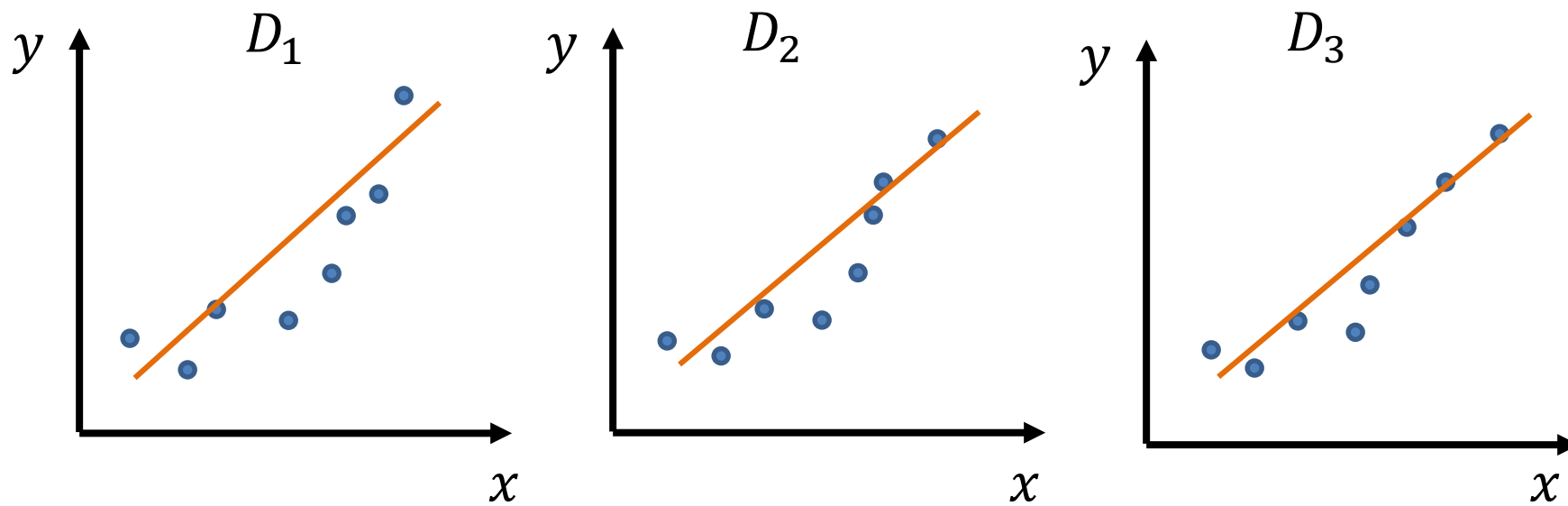for $x_0$**

**(bias)²**

**variance**

**irreducible
error**

\* Prediction randomness comes from randomness in test features AND training data
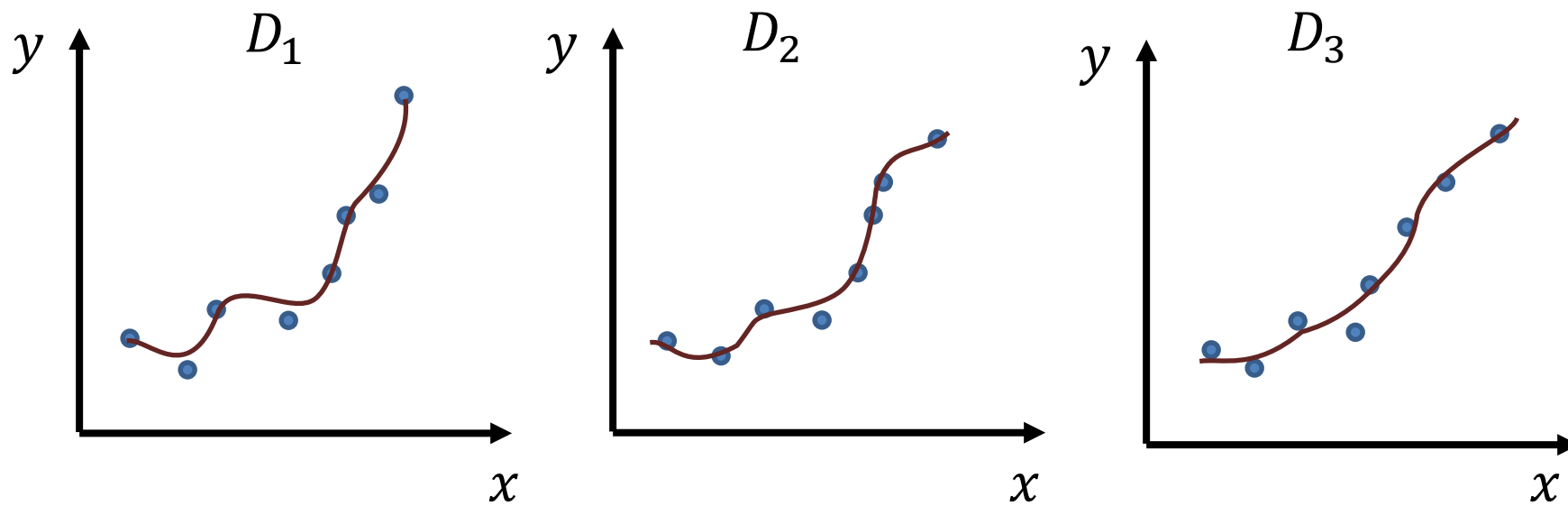
# Decomposition proof sketch

- Here $(\boldsymbol{x})$ is omitted to de-clutter notation

- $\mathbb{E}\left[\left(Y - \hat{f}\right)^2\right] = \mathbb{E}\left[Y^2 + \hat{f}^2 - 2Y\hat{f}\right]$

- $= \mathbb{E}[Y^2] + \mathbb{E}[\hat{f}^2] - \mathbb{E}[2Y\hat{f}]$

- $= Var[Y] + \mathbb{E}[Y]^2 + Var[\hat{f}] + \mathbb{E}[\hat{f}]^2 - 2\mathbb{E}[Y]\mathbb{E}[\hat{f}]$

- $= Var[Y] + Var[\hat{f}] + \left(\mathbb{E}[Y]^2 - 2\mathbb{E}[Y]\mathbb{E}[\hat{f}] + \mathbb{E}[\hat{f}]^2\right)$

- $= Var[Y] + Var[\hat{f}] + \left(\mathbb{E}[Y] - \mathbb{E}[\hat{f}]\right)^2$

* Green slides are non-examinable

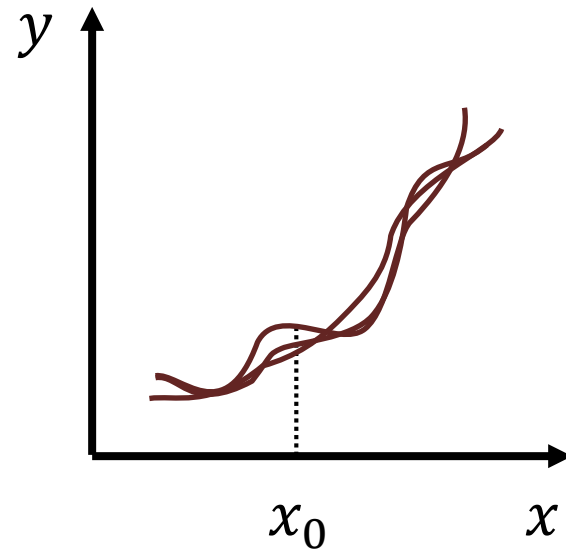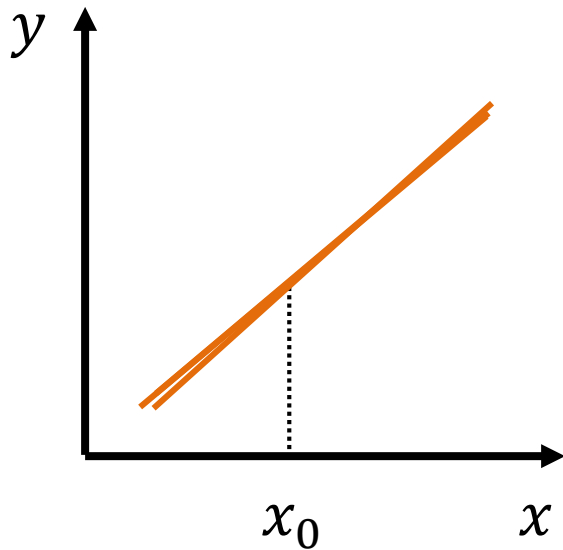# Training data as a random variable
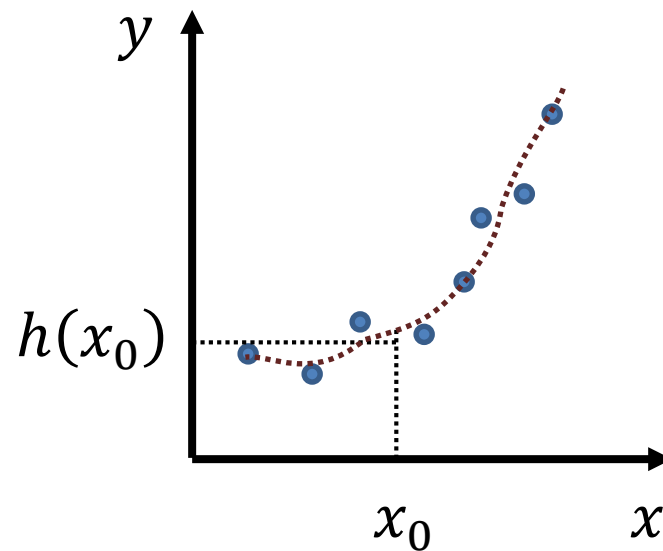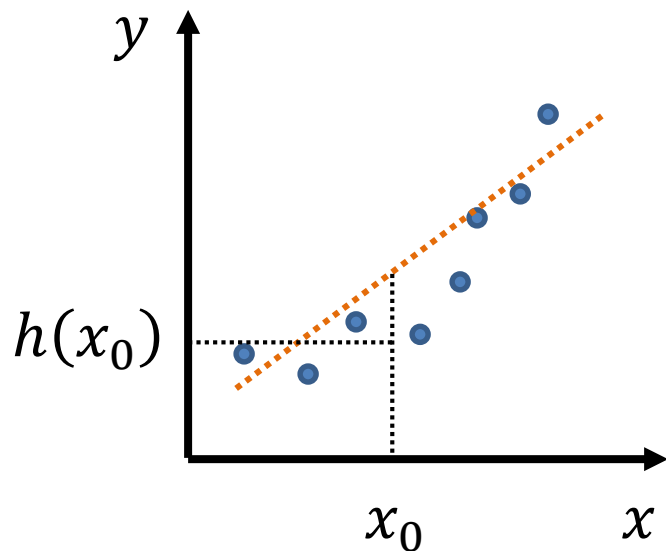
# Training data as a random variable

# Model complexity and variance

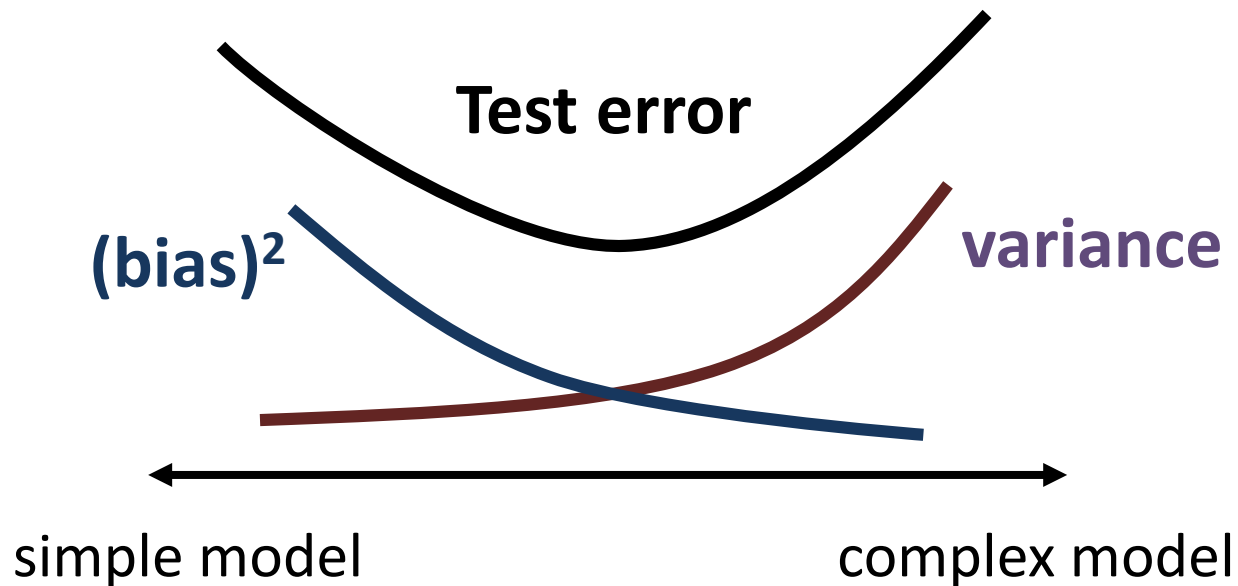- simple model ➔ low variance

- complex model ➔ high variance

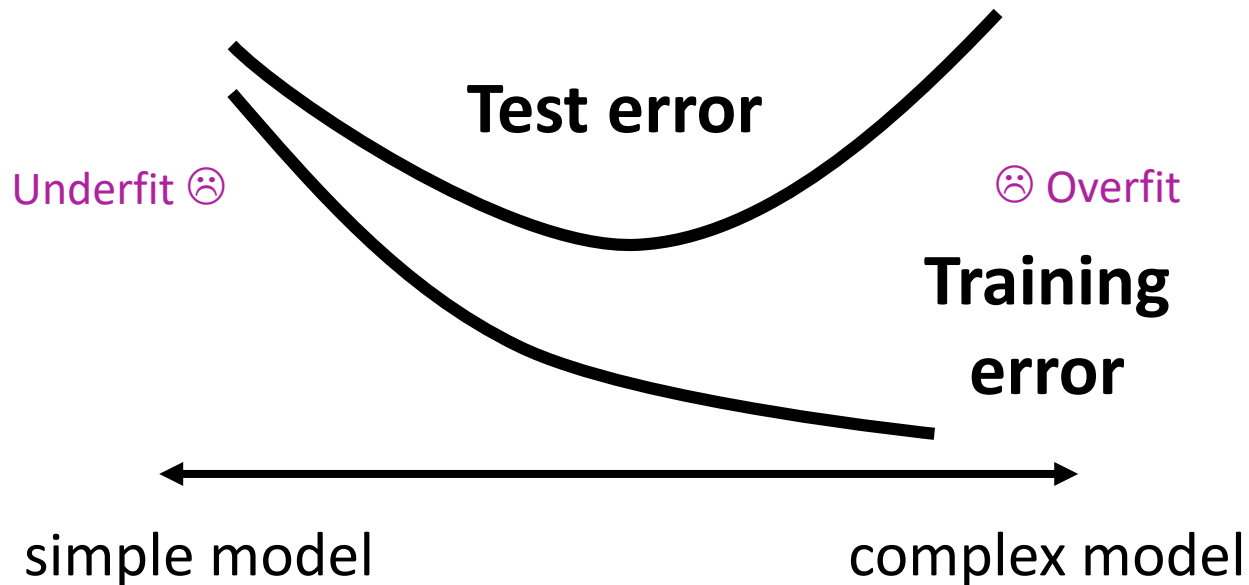# Model complexity and bias

- simple model ➔ high bias

- complex model ➔ low bias

# Bias-variance trade-off

- simple model ➔ high bias, low variance

- complex model ➔ low bias, high variance

# Test error and training error

# Summary

- Regularisation
  - * Irrelevant/multicolinear features → ill-posed problems
  - * Model complexity
  - * Bias-variance trade-off

- Next lecture: Towards neural nets with perceptron