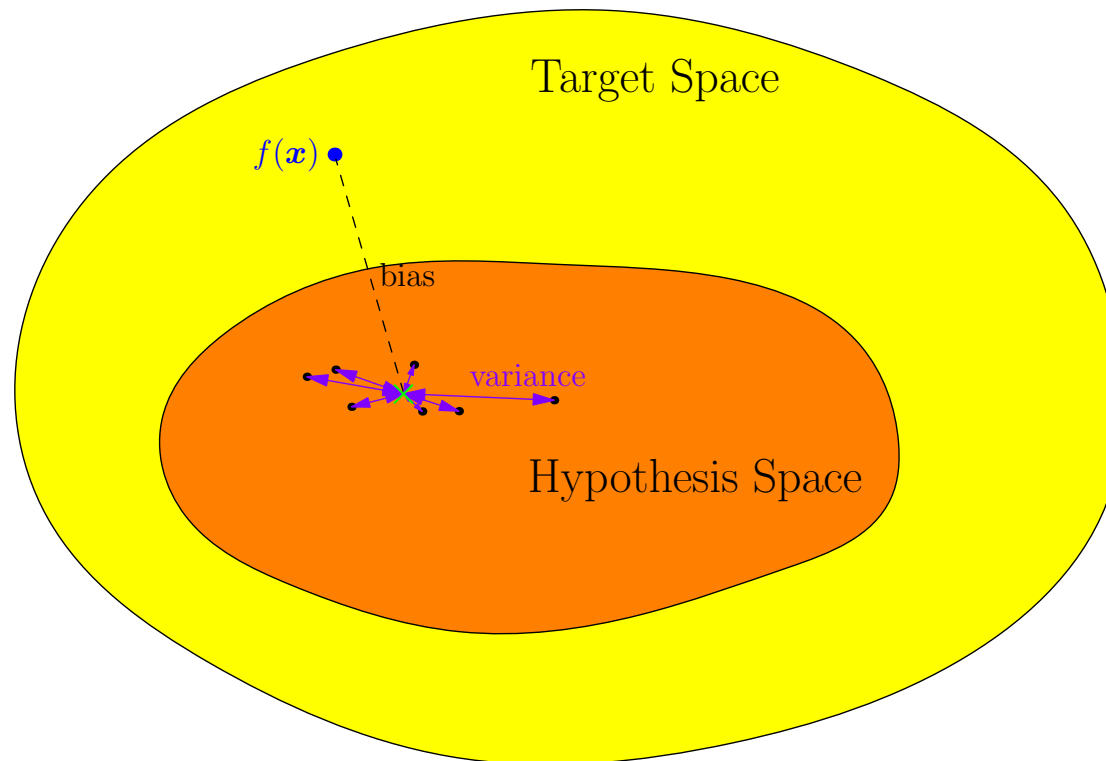


# DISC-NET ML Workshop

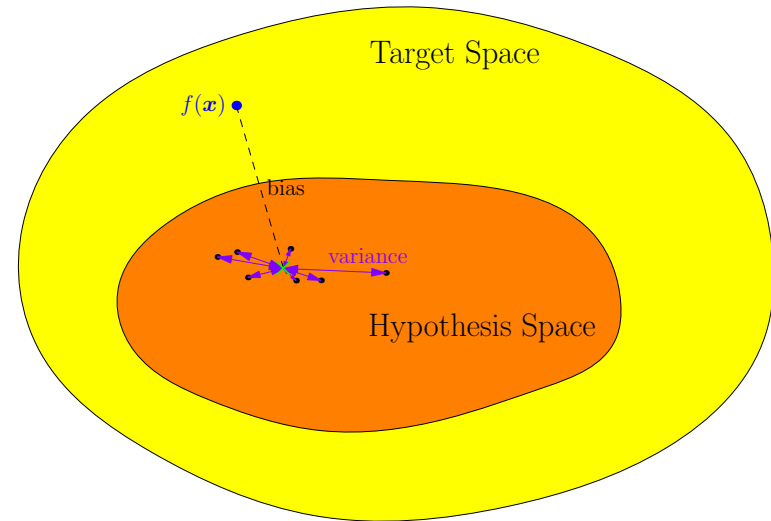
## *Advanced Machine Learning*



*When ML Works, SVMs, Decision Trees, Ensemble Methods,  
Bayesian Inference*

# Outline

1. **What Makes a Good Learning Machine?**
2. SVMs
3. Ensemble Methods
4. Bayesian Inference



# What Makes a Good Learning Machine?

- We want to understand why some machine learning techniques work well and other don't
- To understand why these works we need to understand what makes a good learning machine
- For this we have to get conceptual and think about **generalisation** performance

*generalisation: how well do we do on unseen data as opposed to the training data*

# What Makes a Good Learning Machine?

- We want to understand why some machine learning techniques work well and other don't
- To understand why these works we need to understand what makes a good learning machine
- For this we have to get conceptual and think about **generalisation** performance

*generalisation: how well do we do on unseen data as opposed to the training data*

# What Makes a Good Learning Machine?

- We want to understand why some machine learning techniques work well and other don't
- To understand why these works we need to understand what makes a good learning machine
- For this we have to get conceptual and think about **generalisation** performance

*generalisation: how well do we do on unseen data as opposed to the training data*

# What Makes a Good Learning Machine?

- We want to understand why some machine learning techniques work well and other don't
- To understand why these works we need to understand what makes a good learning machine
- For this we have to get conceptual and think about **generalisation** performance

***generalisation:** how well do we do on unseen data as opposed to the training data*

# What Makes Machine Learning Hard?

- Typically work in high dimensions (i.e. have many features)
- The problem can be over-constrained (i.e. we have conflicting data to deal with)
- The problem can be under-constrained (i.e. there are many possible solutions that are consistent with the data)
- Typically in machine learning the data will be over-constrained in some dimensions and under-constrained in others
- We can't visualise the data to know what is going on

# What Makes Machine Learning Hard?

- Typically work in high dimensions (i.e. have many features)
- The problem can be over-constrained (i.e. we have conflicting data to deal with)
- The problem can be under-constrained (i.e. there are many possible solutions that are consistent with the data)
- Typically in machine learning the data will be over-constrained in some dimensions and under-constrained in others
- We can't visualise the data to know what is going on



# What Makes Machine Learning Hard?

- Typically work in high dimensions (i.e. have many features)
- The problem can be over-constrained (i.e. we have conflicting data to deal with)—solve by minimising an error function
- The problem can be under-constrained (i.e. there are many possible solutions that are consistent with the data)
- Typically in machine learning the data will be over-constrained in some dimensions and under-constrained in others
- We can't visualise the data to know what is going on

# What Makes Machine Learning Hard?

- Typically work in high dimensions (i.e. have many features)
- The problem can be over-constrained (i.e. we have conflicting data to deal with)—solve by minimising an error function
- The problem can be under-constrained (i.e. there are many possible solutions that are consistent with the data)
- Typically in machine learning the data will be over-constrained in some dimensions and under-constrained in others
- We can't visualise the data to know what is going on

# What Makes Machine Learning Hard?

- Typically work in high dimensions (i.e. have many features)
- The problem can be over-constrained (i.e. we have conflicting data to deal with)—solve by minimising an error function
- The problem can be under-constrained (i.e. there are many possible solutions that are consistent with the data)—need to choose a plausible solution
- Typically in machine learning the data will be over-constrained in some dimensions and under-constrained in others
- We can't visualise the data to know what is going on

# What Makes Machine Learning Hard?

- Typically work in high dimensions (i.e. have many features)
- The problem can be over-constrained (i.e. we have conflicting data to deal with)—solve by minimising an error function
- The problem can be under-constrained (i.e. there are many possible solutions that are consistent with the data)—need to choose a plausible solution
- Typically in machine learning the data will be over-constrained in some dimensions and under-constrained in others
- We can't visualise the data to know what is going on

# What Makes Machine Learning Hard?

- Typically work in high dimensions (i.e. have many features)
- The problem can be over-constrained (i.e. we have conflicting data to deal with)—solve by minimising an error function
- The problem can be under-constrained (i.e. there are many possible solutions that are consistent with the data)—need to choose a plausible solution
- Typically in machine learning the data will be over-constrained in some dimensions and under-constrained in others
- We can't visualise the data to know what is going on

# Least Squared Errors

- Suppose we want to learn some function  $f(\mathbf{x})$
- We construct a learning machine that makes a prediction  $\hat{f}(\mathbf{x})$
- We typically choose the machine that minimise a *training error*

$$E_T(\mathcal{D}) = \sum_{\mathbf{x} \in \mathcal{D}} \left( \hat{f}(\mathbf{x}) - f(\mathbf{x}) \right)^2$$

# Least Squared Errors

- Suppose we want to learn some function  $f(\mathbf{x})$
- We construct a learning machine that makes a prediction  $\hat{f}(\mathbf{x})$
- We typically choose the machine that minimise a *training error*

$$E_T(\mathcal{D}) = \sum_{\mathbf{x} \in \mathcal{D}} \left( \hat{f}(\mathbf{x}) - f(\mathbf{x}) \right)^2$$

# Least Squared Errors

- Suppose we want to learn some function  $f(\mathbf{x})$
- We construct a learning machine that makes a prediction  $\hat{f}(\mathbf{x})$
- We typically choose the machine that minimise a *training error*

$$E_T(\mathcal{D}) = \sum_{\mathbf{x} \in \mathcal{D}} \left( \hat{f}(\mathbf{x}) - f(\mathbf{x}) \right)^2$$



# Least Squared Errors

- Suppose we want to learn some function  $f(\mathbf{x})$
- We construct a learning machine that makes a prediction  $\hat{f}(\mathbf{x})$
- We typically choose the machine that minimise a *training error*

$$E_T(\mathcal{D}) = \sum_{\mathbf{x} \in \mathcal{D}} \left( \hat{f}(\mathbf{x}) - f(\mathbf{x}) \right)^2 = \sum_{i=1}^m \left( \hat{f}(\mathbf{x}_i) - y_i \right)^2$$

where  $\mathcal{D} = \{(\mathbf{x}_i, y_i = f(\mathbf{x}_i))\}_{i=1}^m$  is a set of size  $m$ , sampled from the set of all inputs,  $\mathcal{X}$ , according to a probability distribution  $p(\mathbf{x})$  describing where our data is

# Least Squared Errors

- Suppose we want to learn some function  $f(\mathbf{x})$
- We construct a learning machine that makes a prediction  $\hat{f}(\mathbf{x})$
- We typically choose the machine that minimise a *training error*

$$E_T(\mathcal{D}) = \sum_{\mathbf{x} \in \mathcal{D}} \left( \hat{f}(\mathbf{x}) - f(\mathbf{x}) \right)^2 = \sum_{i=1}^m \left( \hat{f}(\mathbf{x}_i) - y_i \right)^2$$

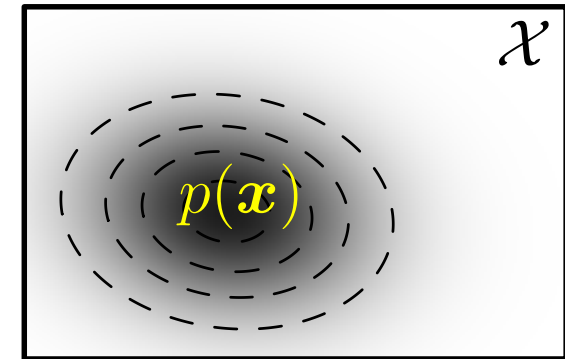
where  $\mathcal{D} = \{(\mathbf{x}_i, y_i = f(\mathbf{x}_i))\}_{i=1}^m$  is a set of size  $m$ , sampled from the set of all inputs,  $\mathcal{X}$ , according to a probability distribution  $p(\mathbf{x})$  describing where our data is

- Choose machine with smallest training error  $\hat{f}(\mathbf{x}|\mathcal{D})$

# Generalisation Error

- We want to minimise the *generalisation error* which in this case is

$$E_G(\mathcal{D}) = \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \left( \hat{f}(\mathbf{x}|\mathcal{D}) - f(\mathbf{x}) \right)^2$$



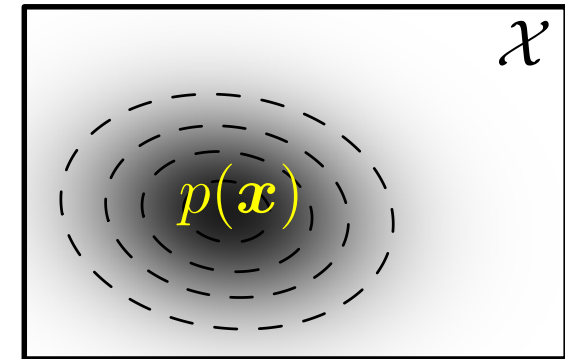
(we can estimate this if we have some examples with known labels  $y_i = f(\mathbf{x}_i)$  which we have not trained on)

- We want to minimise  $E_G(\mathcal{D})$  but in practice we are minimising  $E_T(\mathcal{D})$ , *what could possibly go wrong?*

# Generalisation Error

- We want to minimise the *generalisation error* which in this case is

$$E_G(\mathcal{D}) = \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \left( \hat{f}(\mathbf{x}|\mathcal{D}) - f(\mathbf{x}) \right)^2$$



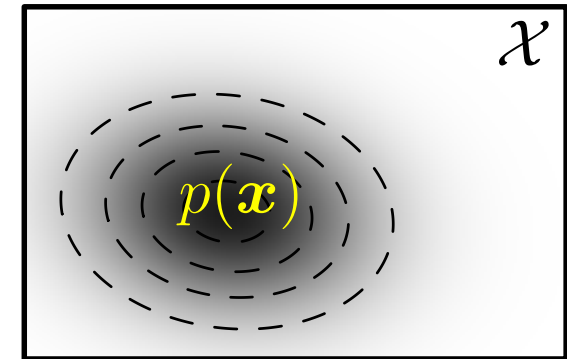
(we can estimate this if we have some examples with known labels  $y_i = f(\mathbf{x}_i)$  which we have not trained on)

- We want to minimise  $E_G(\mathcal{D})$  but in practice we are minimising  $E_T(\mathcal{D})$ , what could possibly go wrong?

# Generalisation Error

- We want to minimise the *generalisation error* which in this case is

$$E_G(\mathcal{D}) = \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \left( \hat{f}(\mathbf{x}|\mathcal{D}) - f(\mathbf{x}) \right)^2$$



(we can estimate this if we have some examples with known labels  $y_i = f(\mathbf{x}_i)$  which we have not trained on)

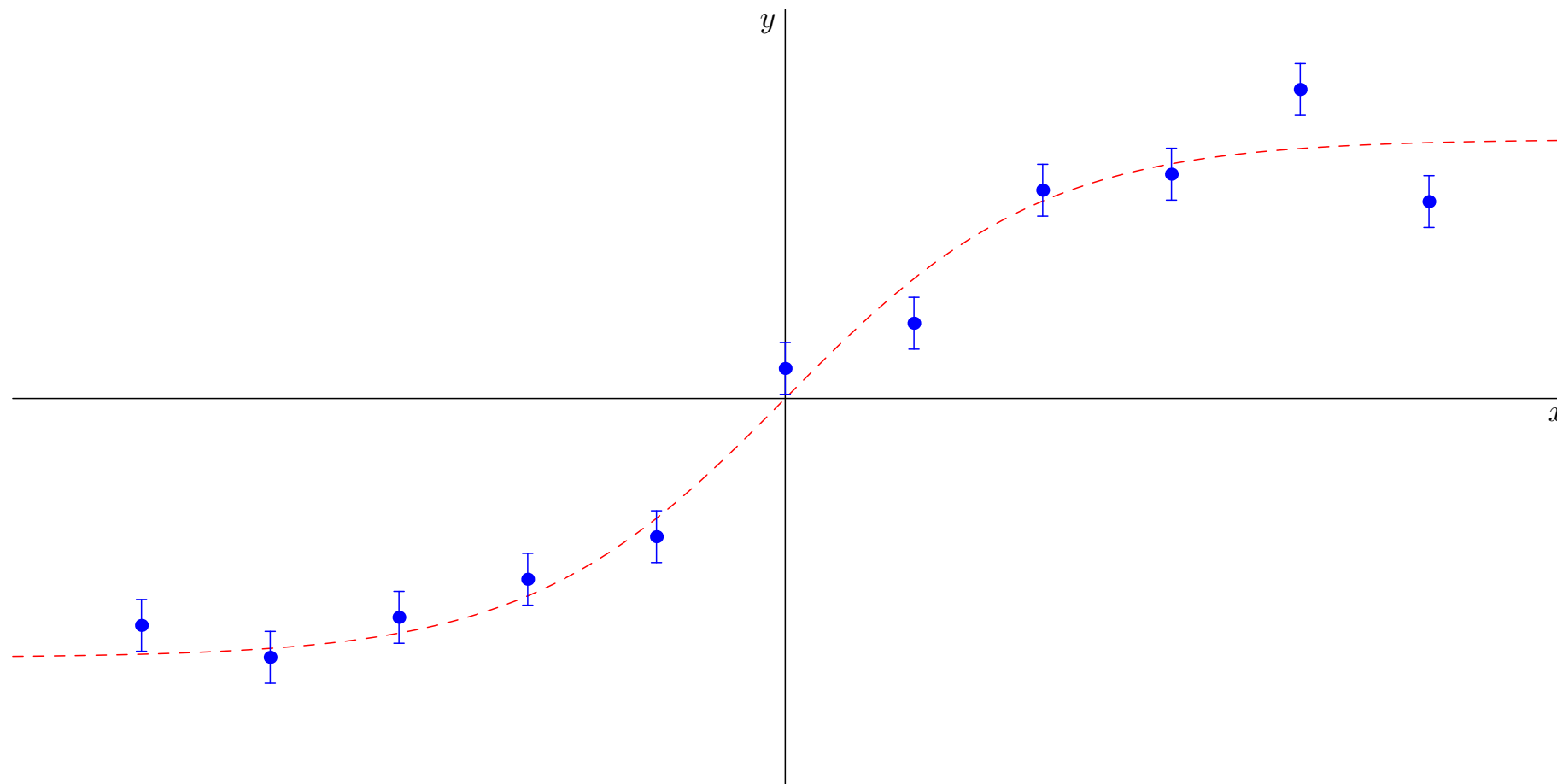
- We want to minimise  $E_G(\mathcal{D})$  but in practice we are minimising  $E_T(\mathcal{D})$ , what could possibly go wrong?

# Too Simple or Too Complex?

- Fit  $\hat{f}(x|\mathcal{D})$  to data

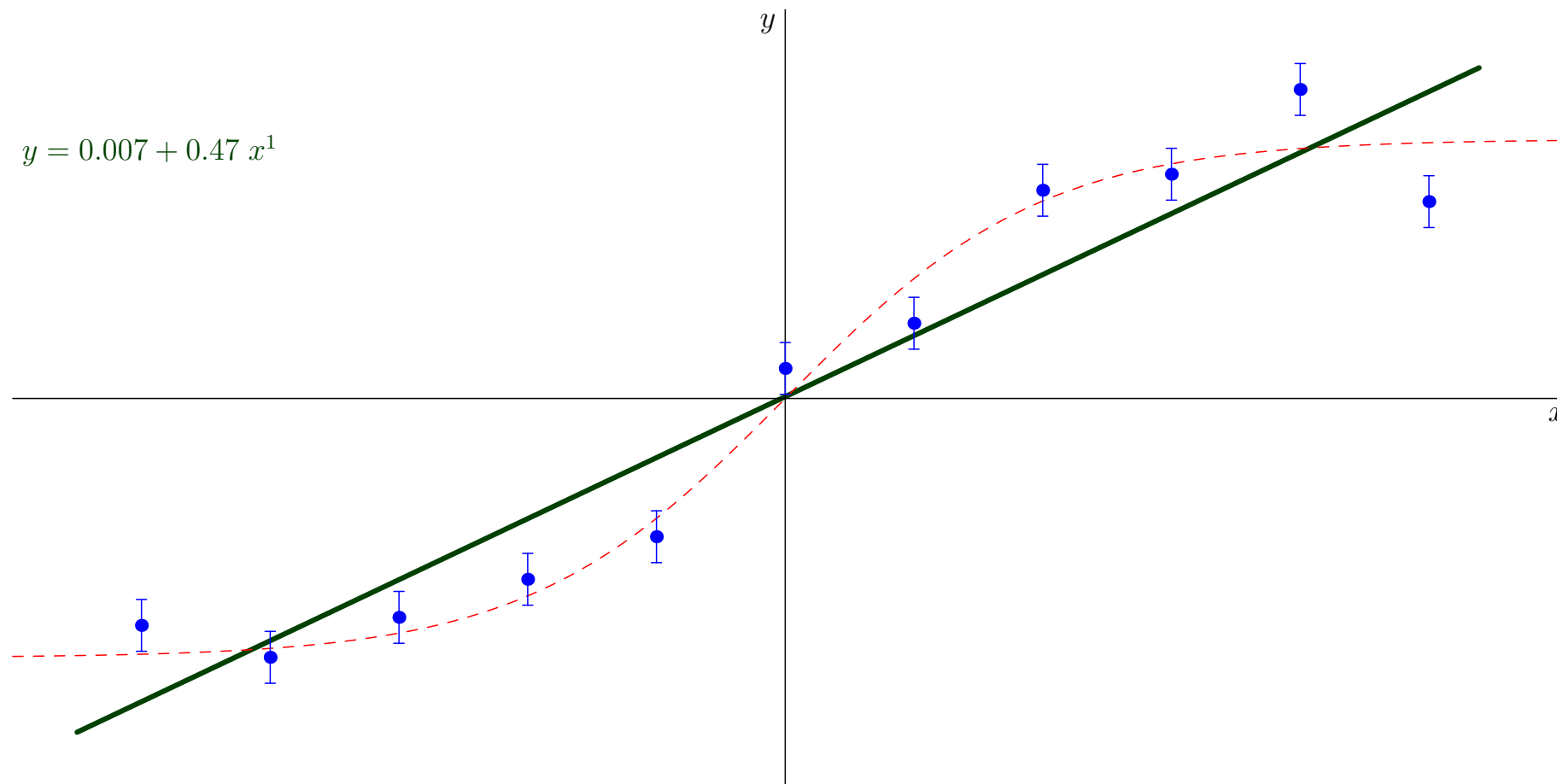
# Too Simple or Too Complex?

- Fit  $\hat{f}(x|\mathcal{D})$  to data



# Too Simple or Too Complex?

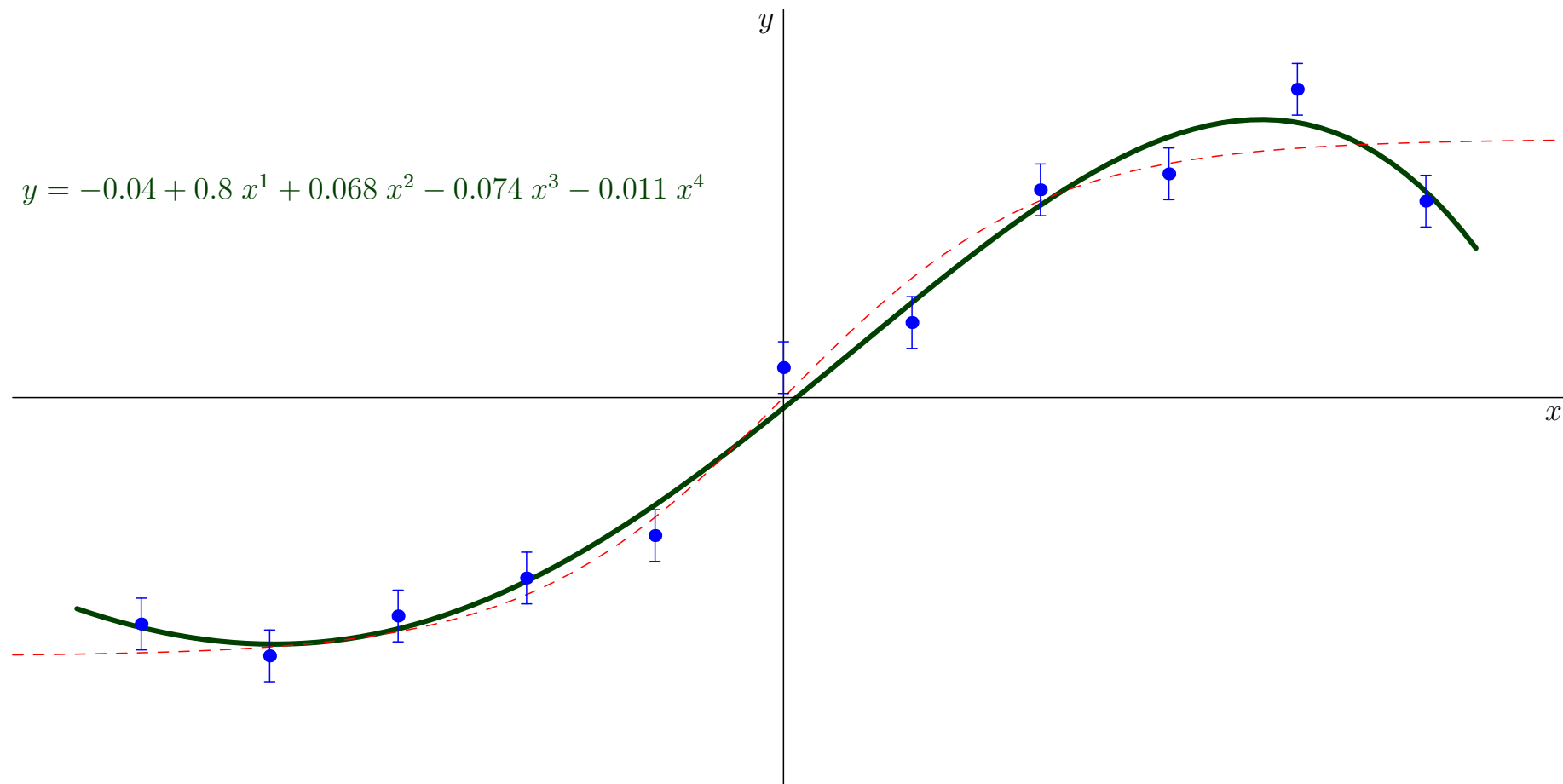
- Fit  $\hat{f}(x|\mathcal{D})$  to data





# Too Simple or Too Complex?

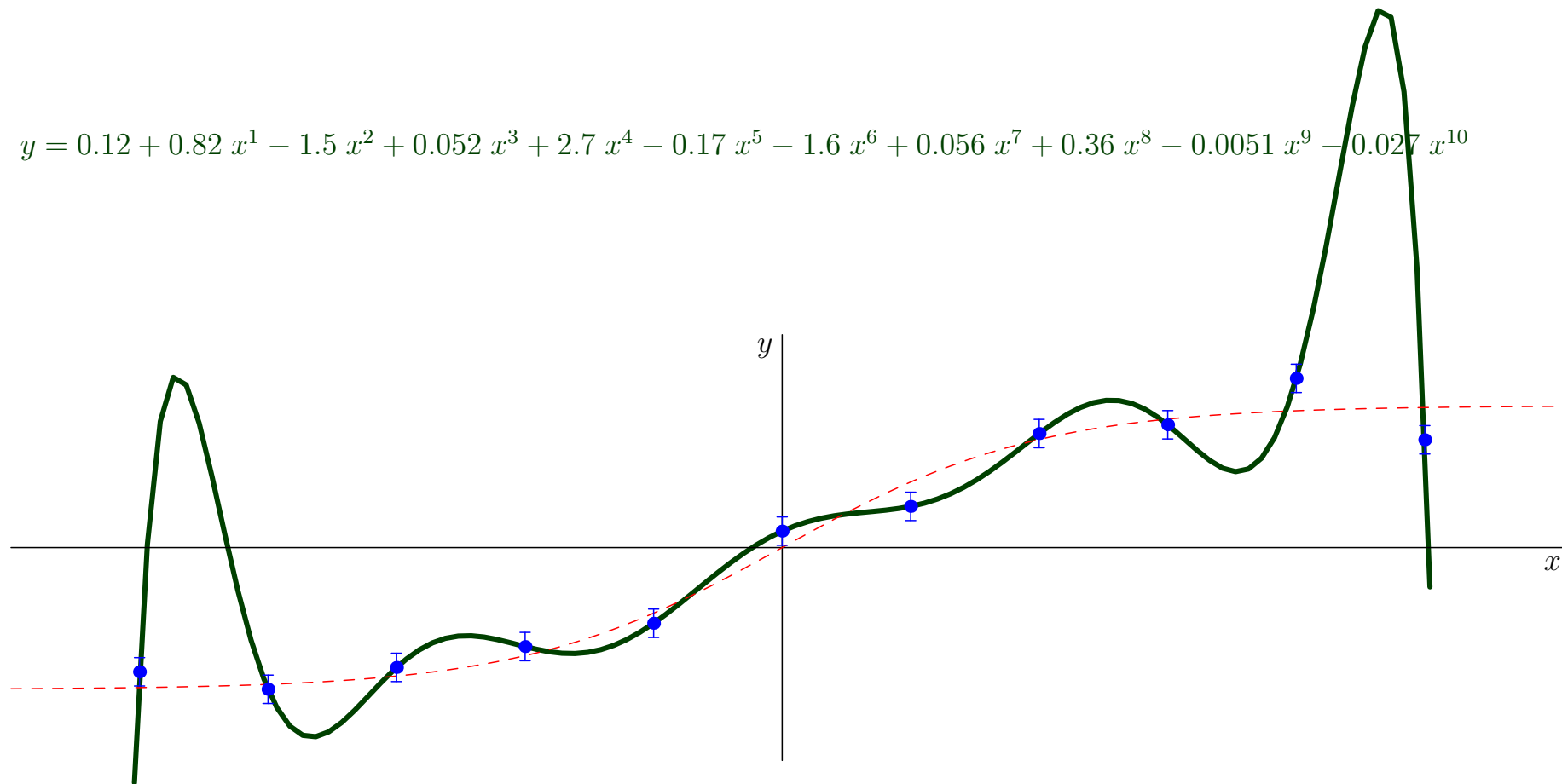
- Fit  $\hat{f}(x|\mathcal{D})$  to data



# Too Simple or Too Complex?

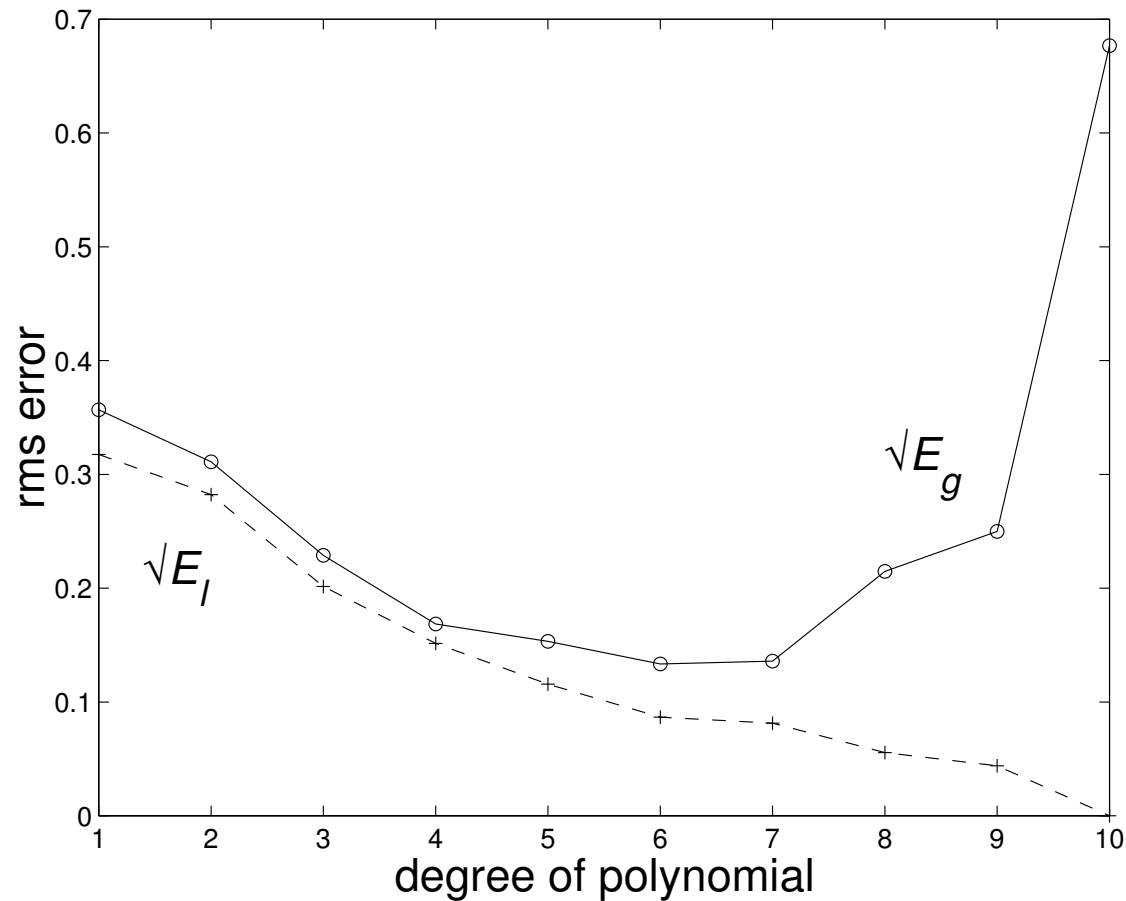
- Fit  $\hat{f}(x|\mathcal{D})$  to data

$$y = 0.12 + 0.82 x^1 - 1.5 x^2 + 0.052 x^3 + 2.7 x^4 - 0.17 x^5 - 1.6 x^6 + 0.056 x^7 + 0.36 x^8 - 0.0051 x^9 - 0.027 x^{10}$$



# Measuring Generalisation Error for Regression

- Consider the regression example. The root mean squared error is



# Expected Generalisation Performance

- Our generalisation performance will depend on our training set,  $\mathcal{D}$
- To reason about generalisation we can ask what is the *expected generalisation*, that is, when we average over all different data sets of size  $m$  drawn independently from  $p(\mathbf{x})$
- For each data set,  $\mathcal{D}$ , we would learn a different approximator  $\hat{f}(\mathbf{x}|\mathcal{D})$  (usually through weights  $\mathbf{w}_{\mathcal{D}}$ )
- Note that in practice we only get one data set. We might be lucky and do better than the expected generalisation or we might be unlucky and do worse

# Expected Generalisation Performance

- Our generalisation performance will depend on our training set,  $\mathcal{D}$
- To reason about generalisation we can ask what is the *expected generalisation*, that is, when we average over all different data sets of size  $m$  drawn independently from  $p(\mathbf{x})$
- For each data set,  $\mathcal{D}$ , we would learn a different approximator  $\hat{f}(\mathbf{x}|\mathcal{D})$  (usually through weights  $\mathbf{w}_{\mathcal{D}}$ )
- Note that in practice we only get one data set. We might be lucky and do better than the expected generalisation or we might be unlucky and do worse

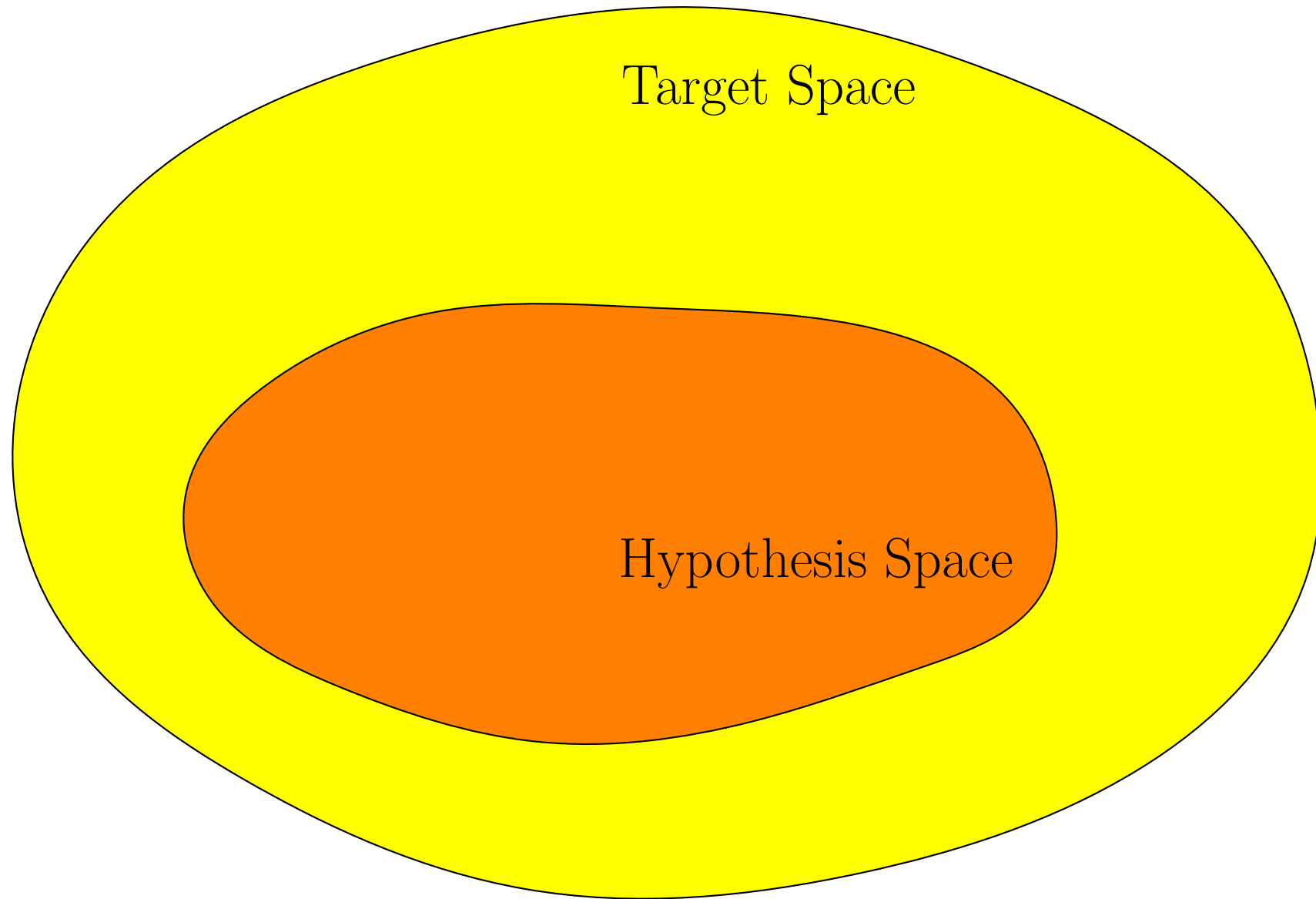
# Expected Generalisation Performance

- Our generalisation performance will depend on our training set,  $\mathcal{D}$
- To reason about generalisation we can ask what is the *expected generalisation*, that is, when we average over all different data sets of size  $m$  drawn independently from  $p(\mathbf{x})$
- For each data set,  $\mathcal{D}$ , we would learn a different approximator  $\hat{f}(\mathbf{x}|\mathcal{D})$  (usually through weights  $\mathbf{w}_{\mathcal{D}}$ )
- Note that in practice we only get one data set. We might be lucky and do better than the expected generalisation or we might be unlucky and do worse

# Expected Generalisation Performance

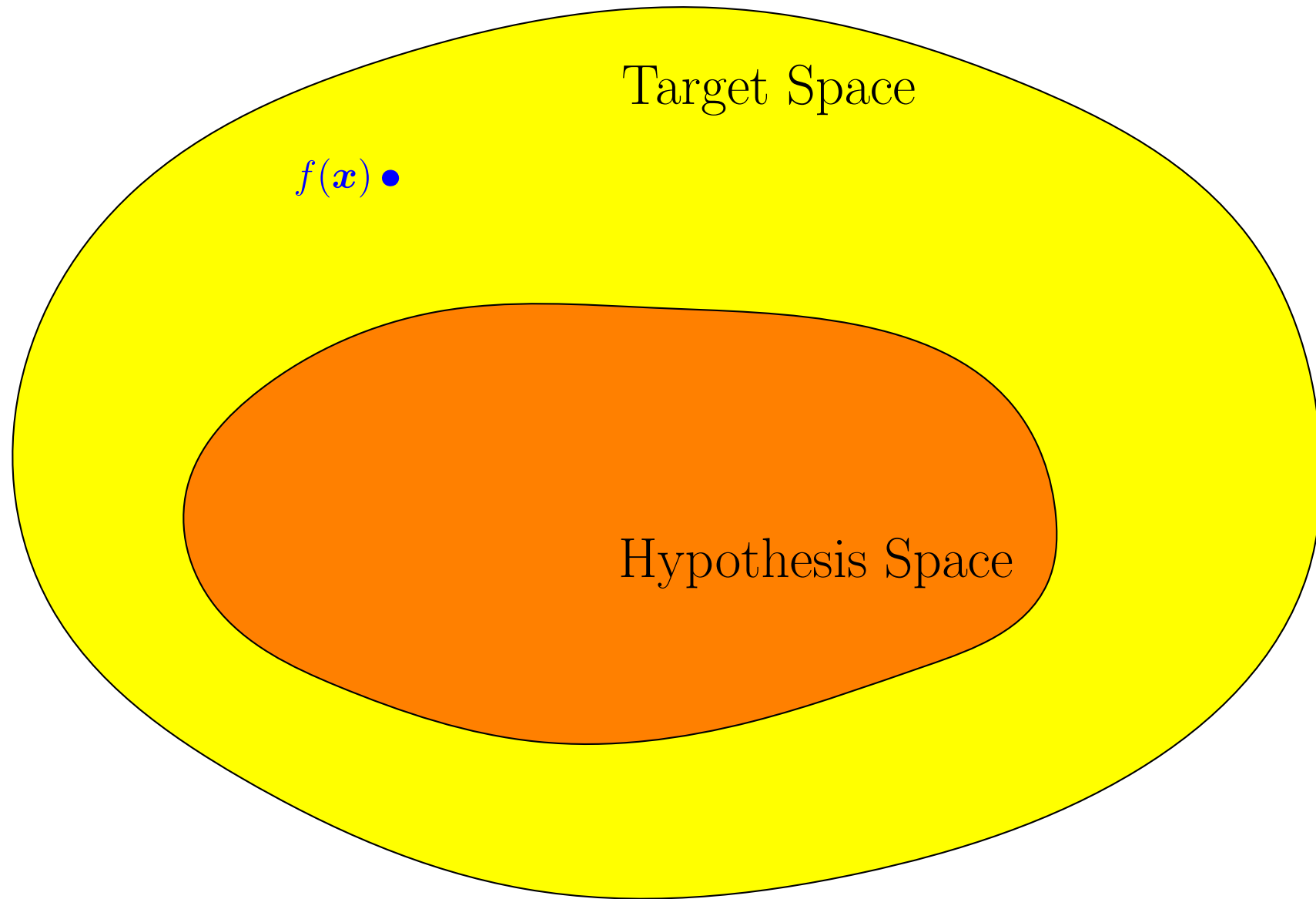
- Our generalisation performance will depend on our training set,  $\mathcal{D}$
- To reason about generalisation we can ask what is the *expected generalisation*, that is, when we average over all different data sets of size  $m$  drawn independently from  $p(\mathbf{x})$
- For each data set,  $\mathcal{D}$ , we would learn a different approximator  $\hat{f}(\mathbf{x}|\mathcal{D})$  (usually through weights  $\mathbf{w}_{\mathcal{D}}$ )
- Note that in practice we only get one data set. We might be lucky and do better than the expected generalisation or we might be unlucky and do worse

# Approximation and Estimation Errors

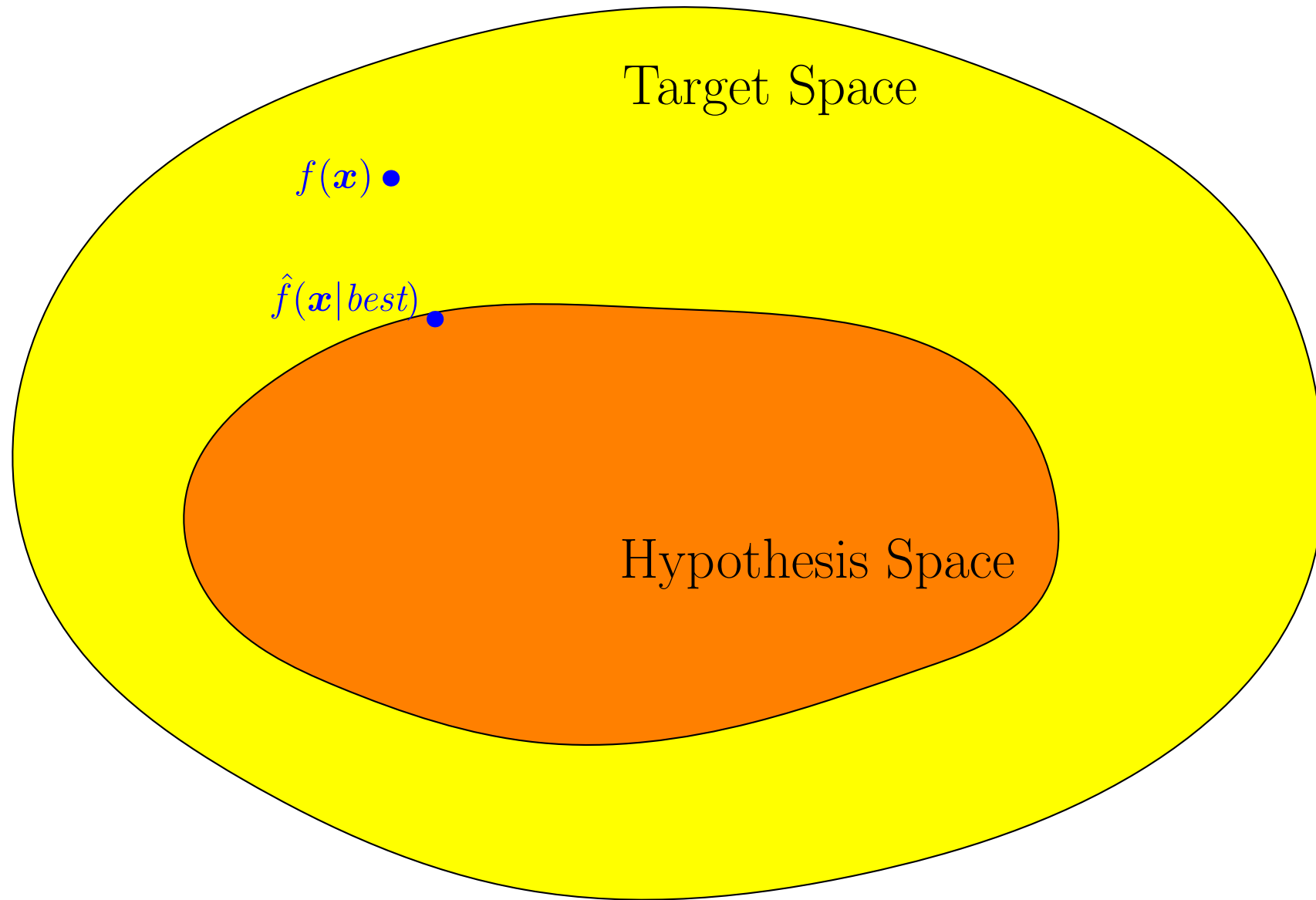




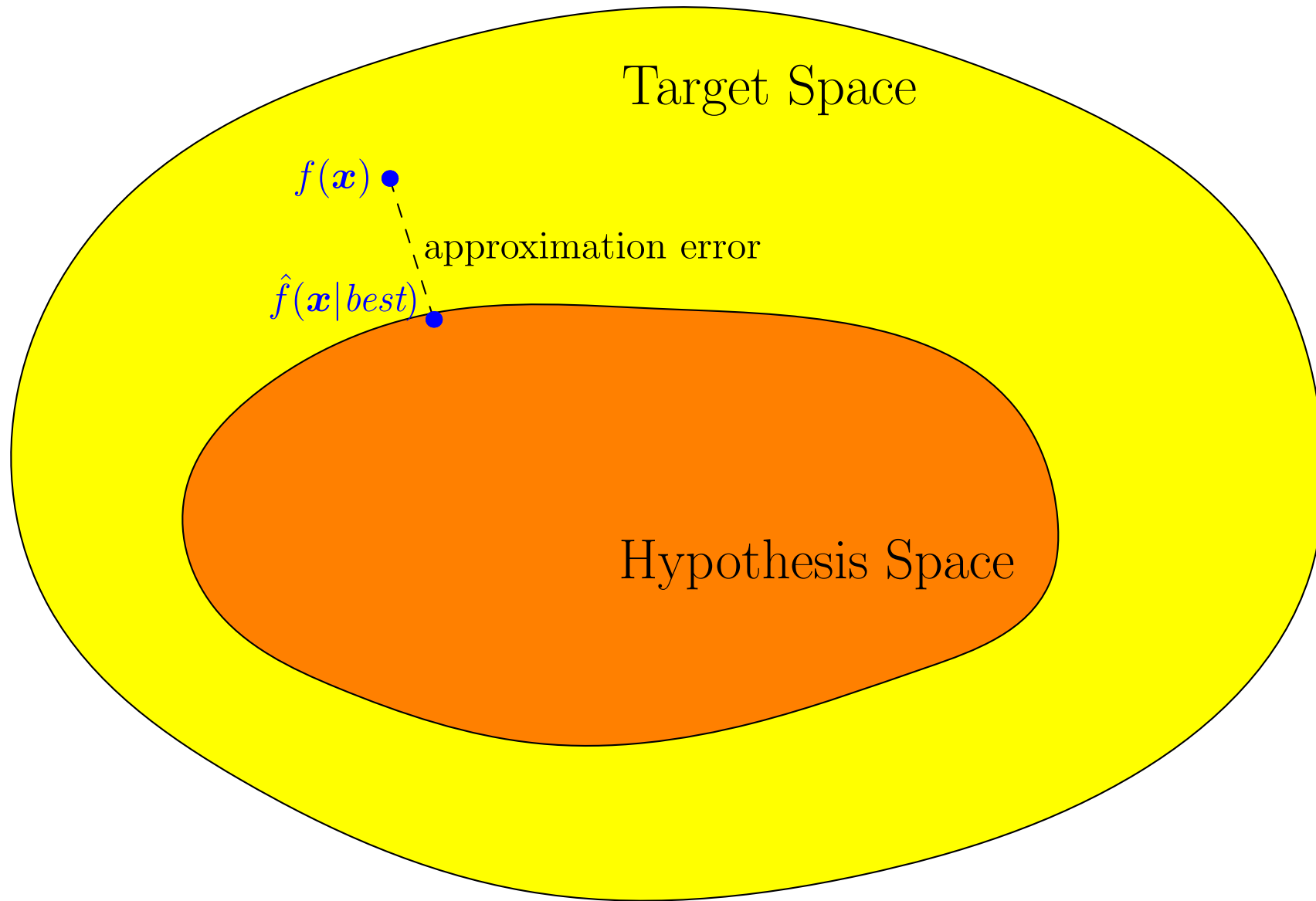
# Approximation and Estimation Errors



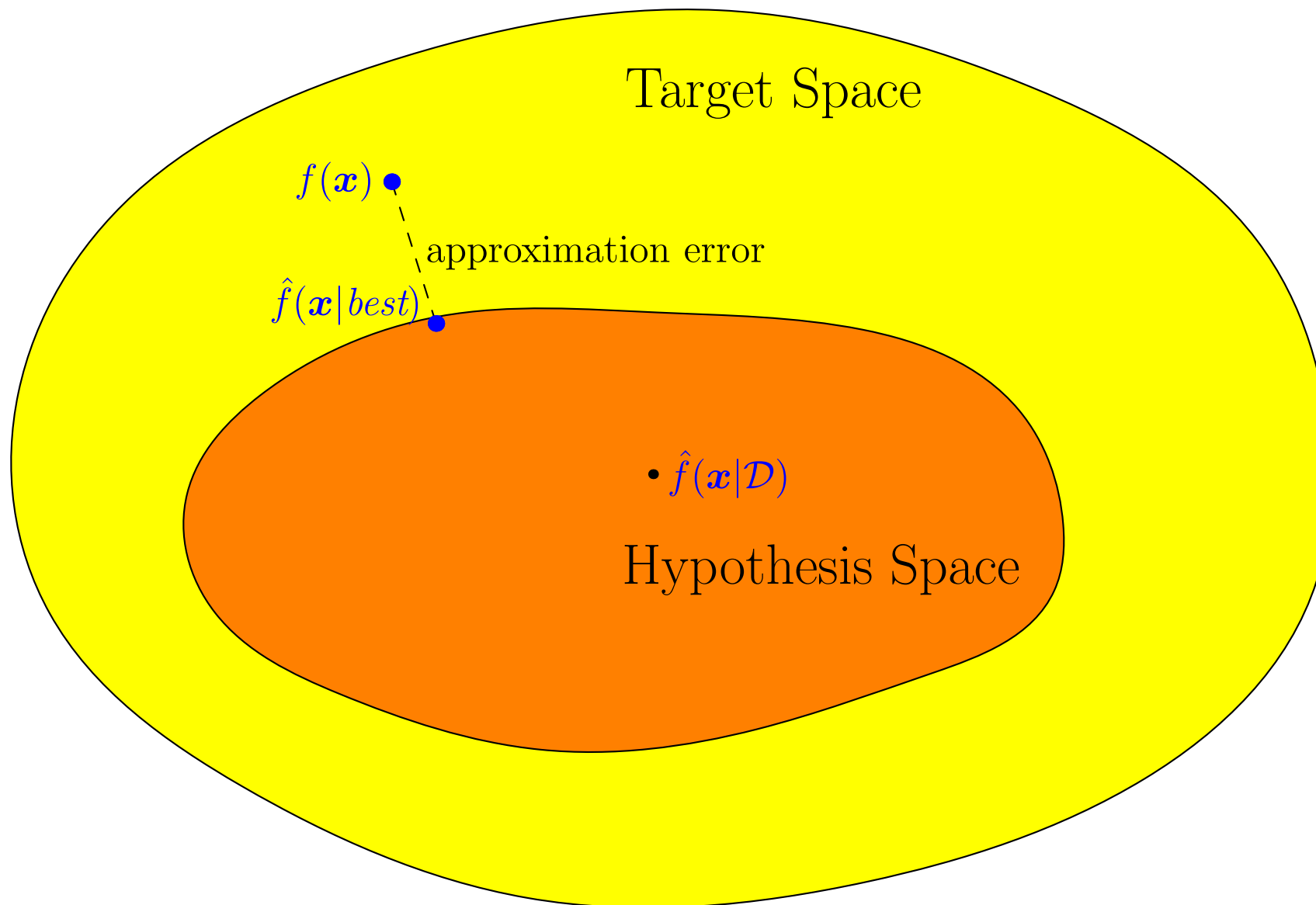
# Approximation and Estimation Errors



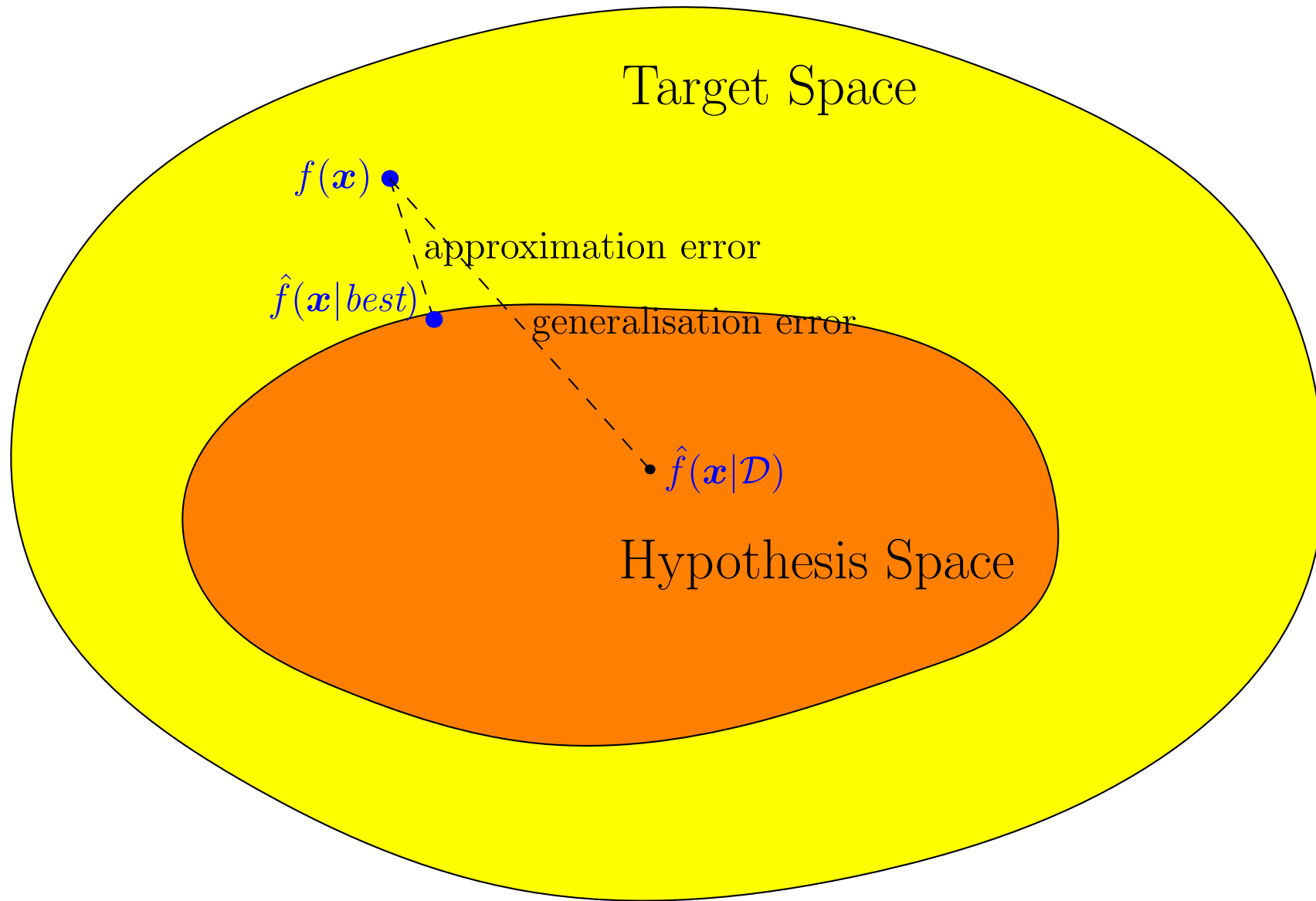
# Approximation and Estimation Errors



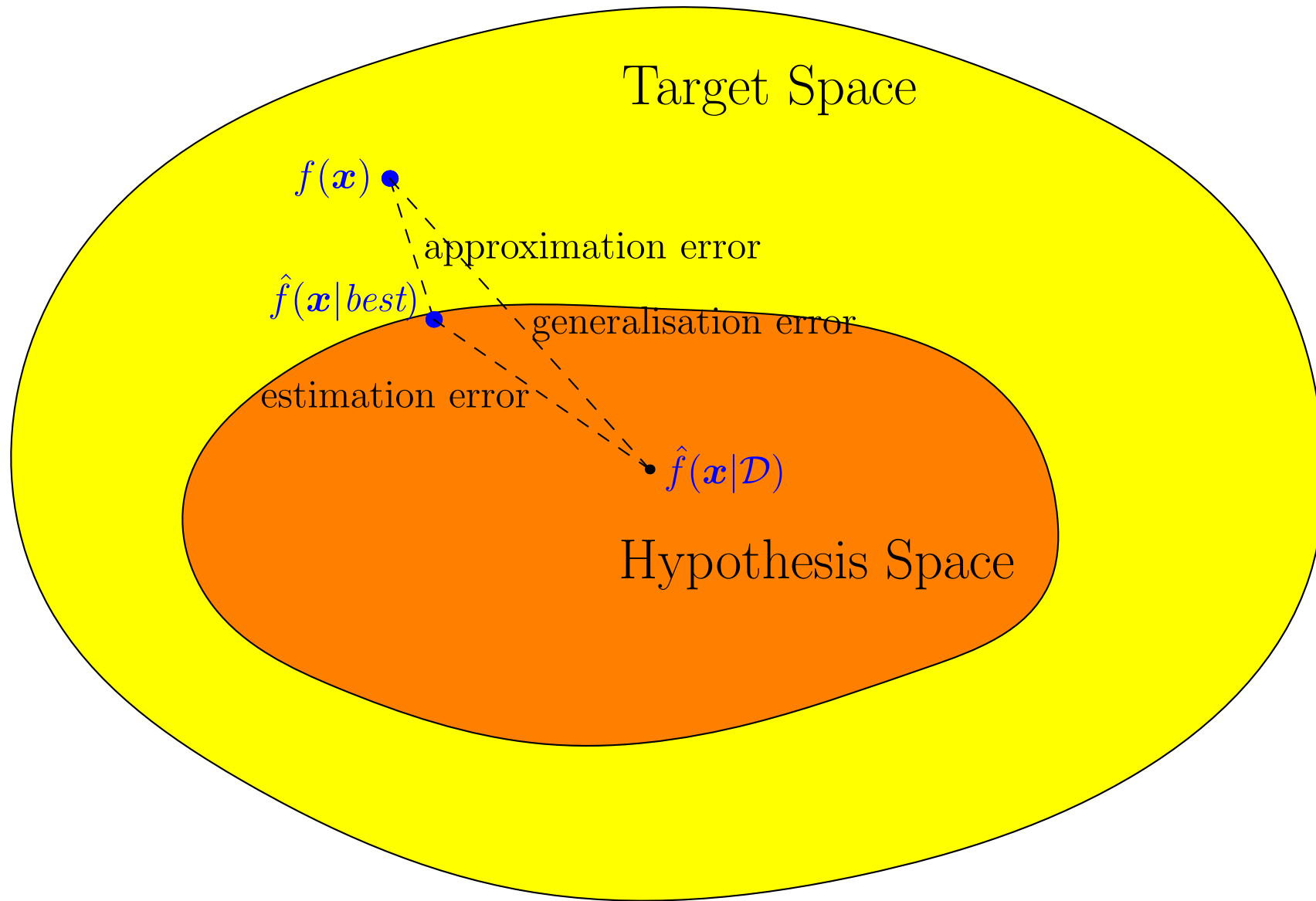
# Approximation and Estimation Errors



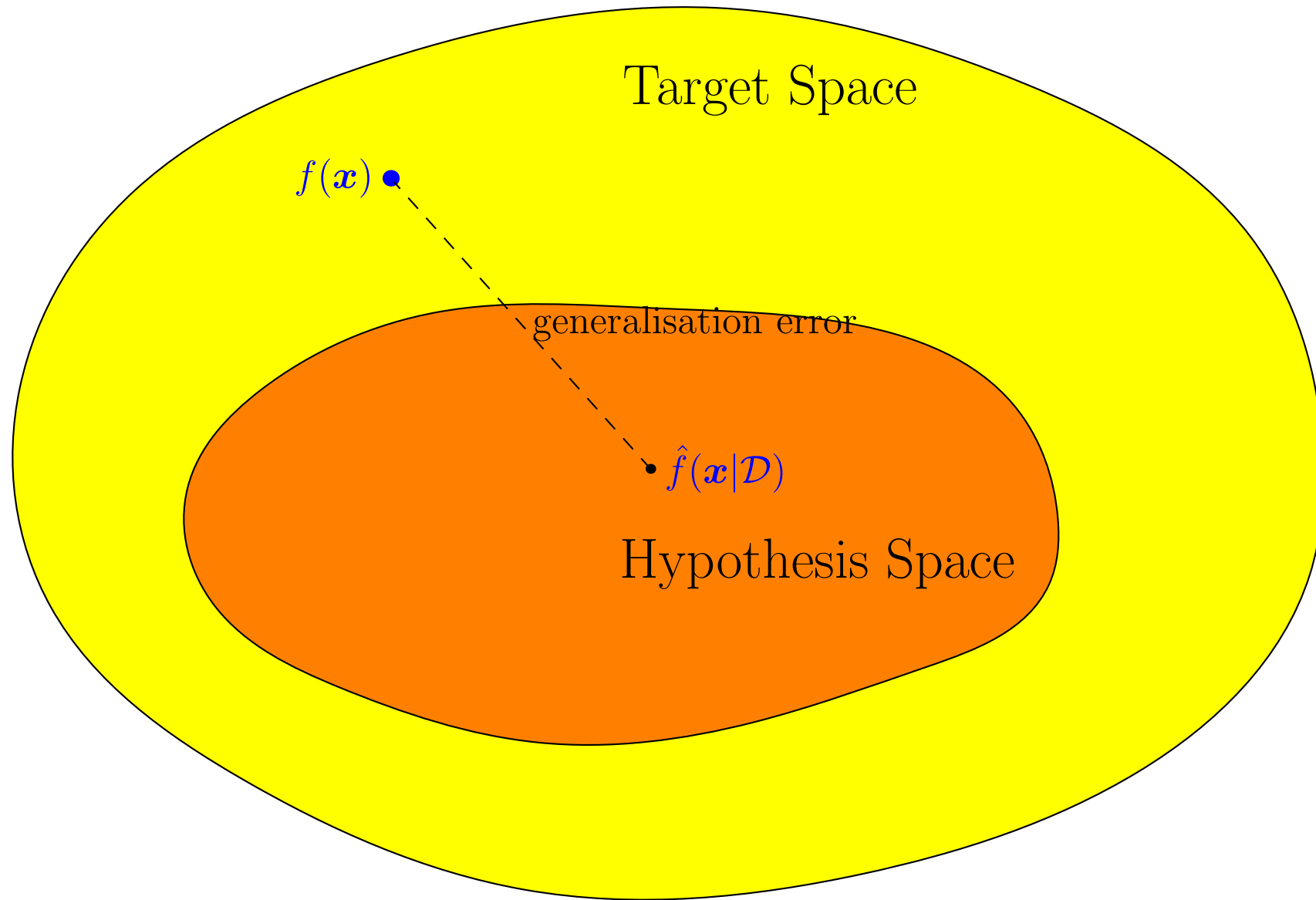
# Approximation and Estimation Errors



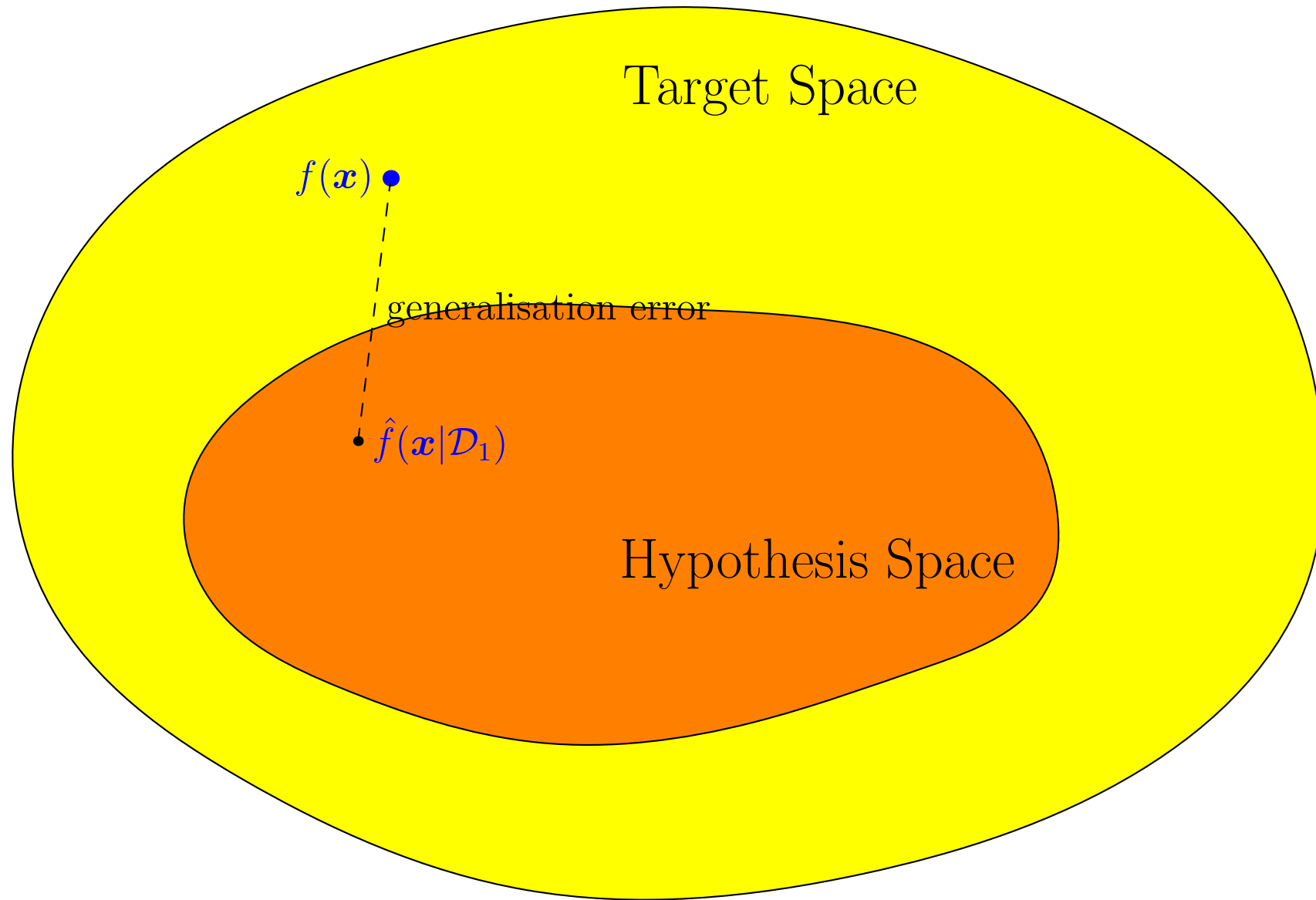
# Approximation and Estimation Errors



# Approximation and Estimation Errors

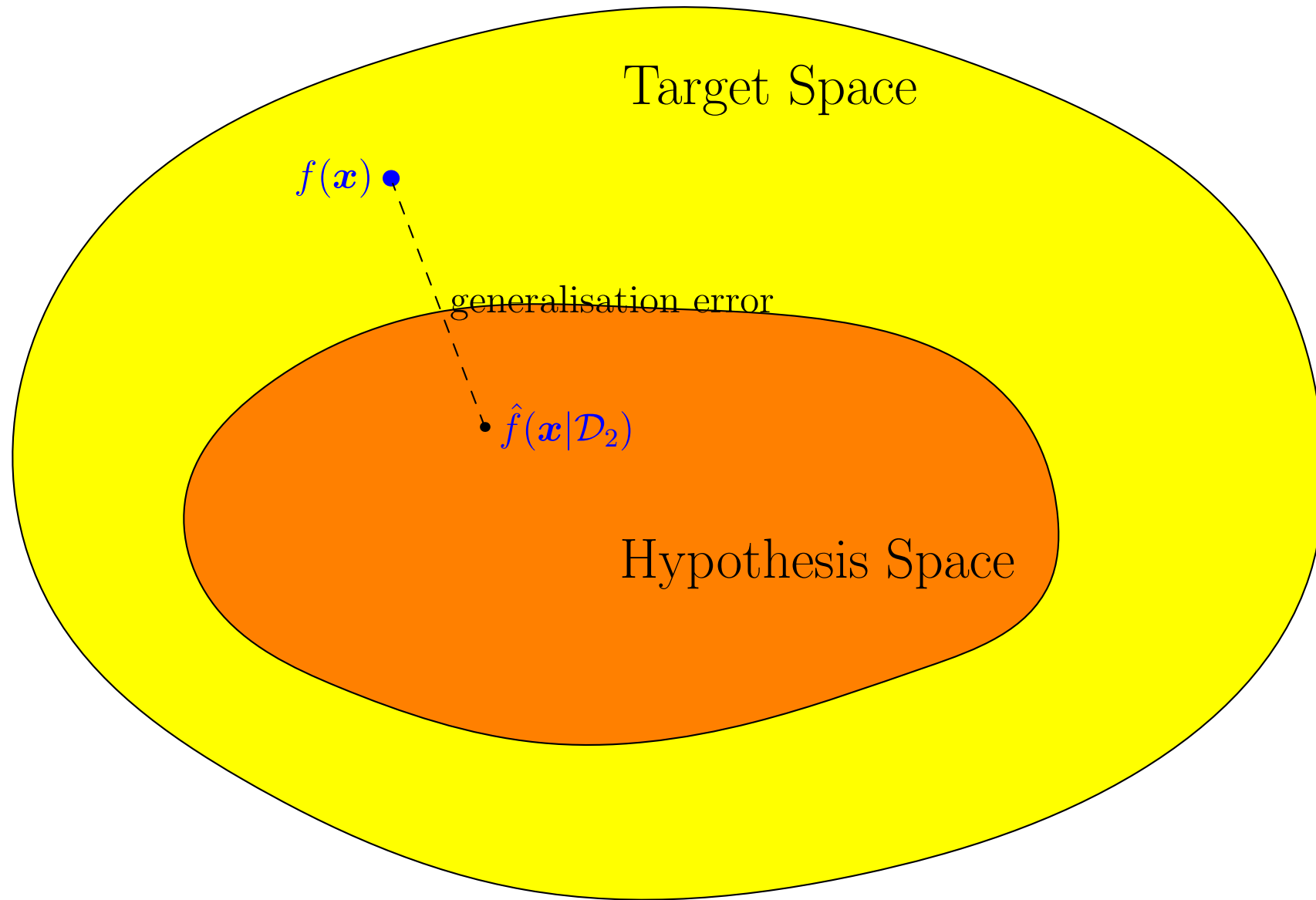


# Approximation and Estimation Errors

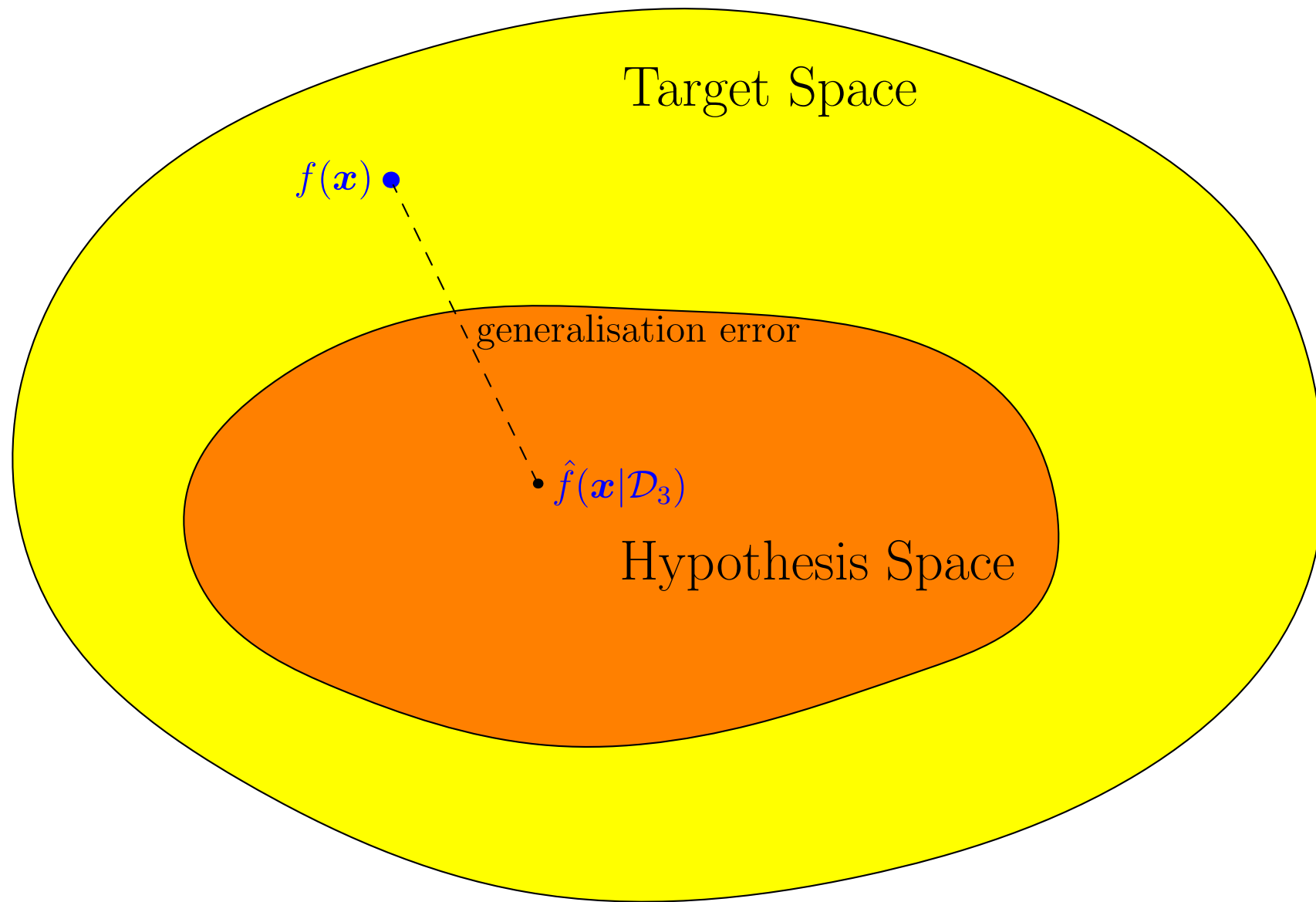




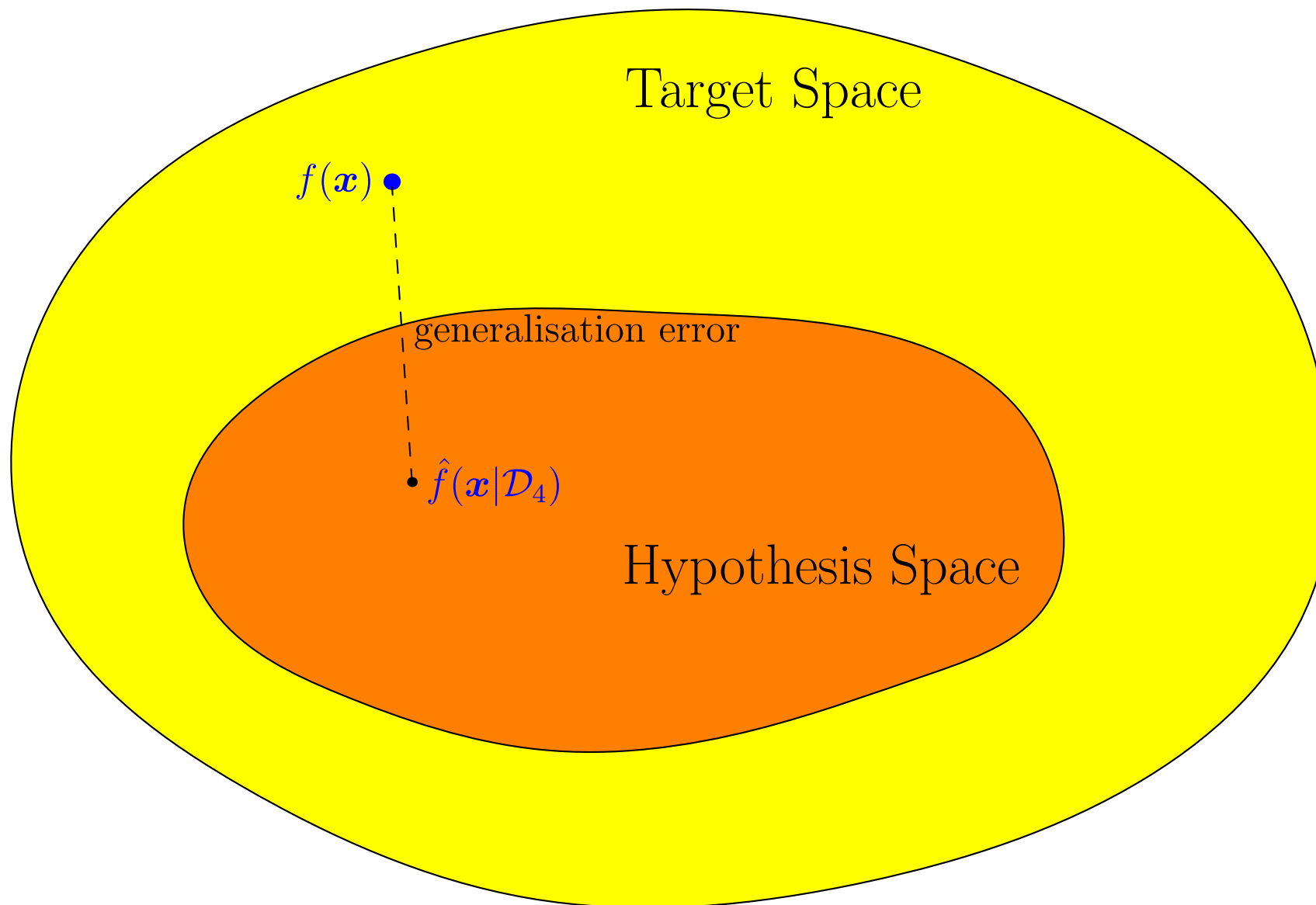
# Approximation and Estimation Errors



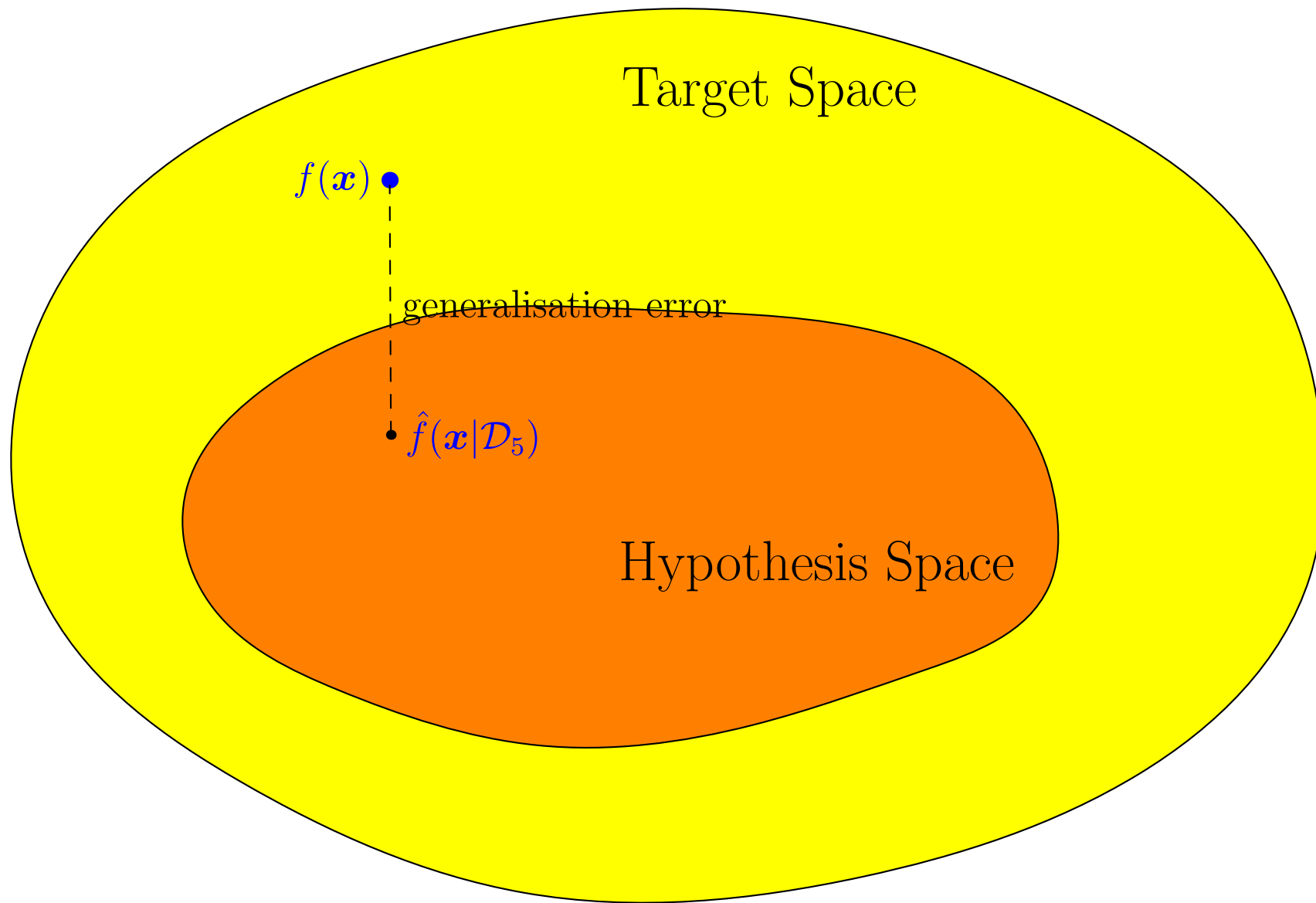
# Approximation and Estimation Errors



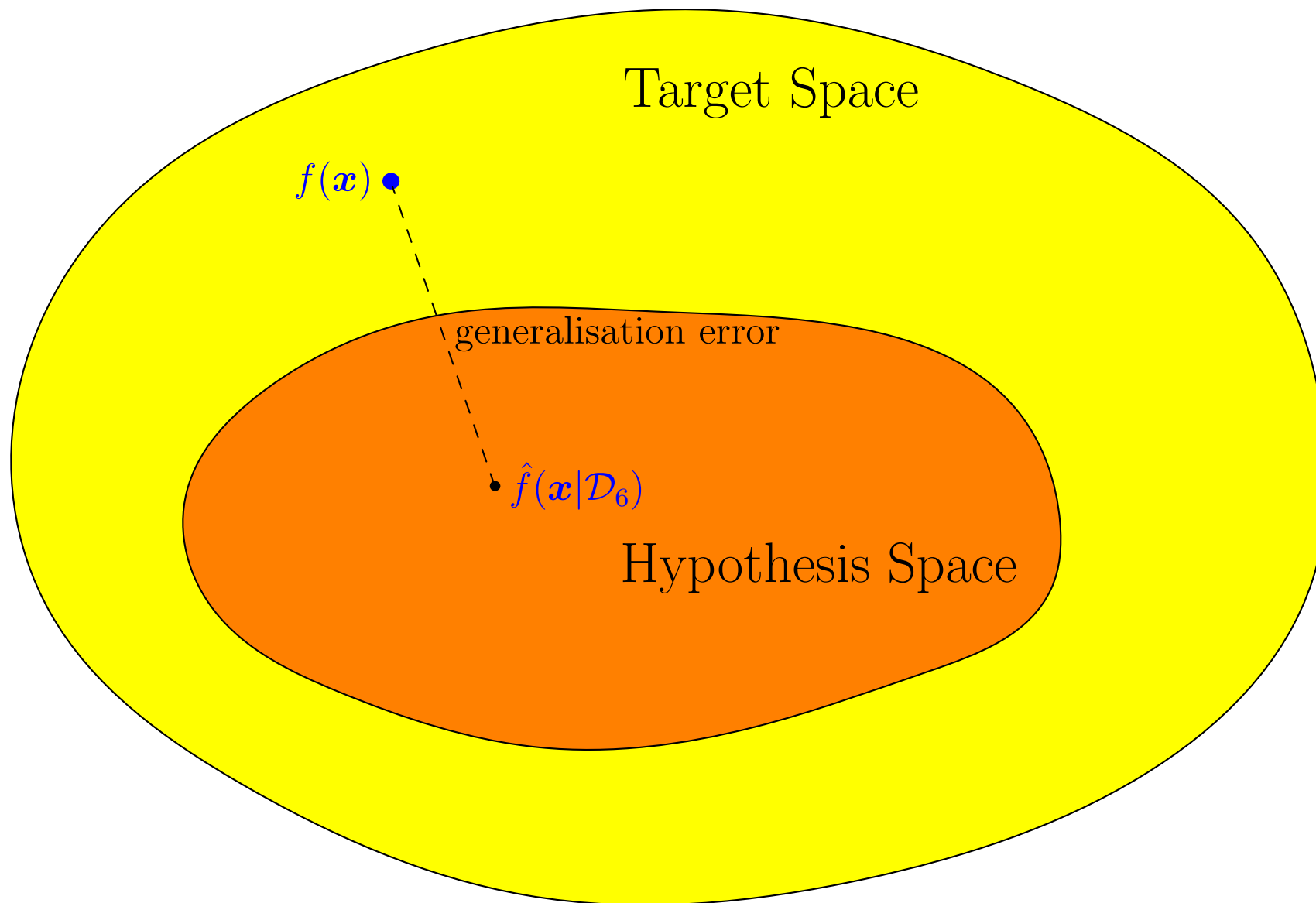
# Approximation and Estimation Errors



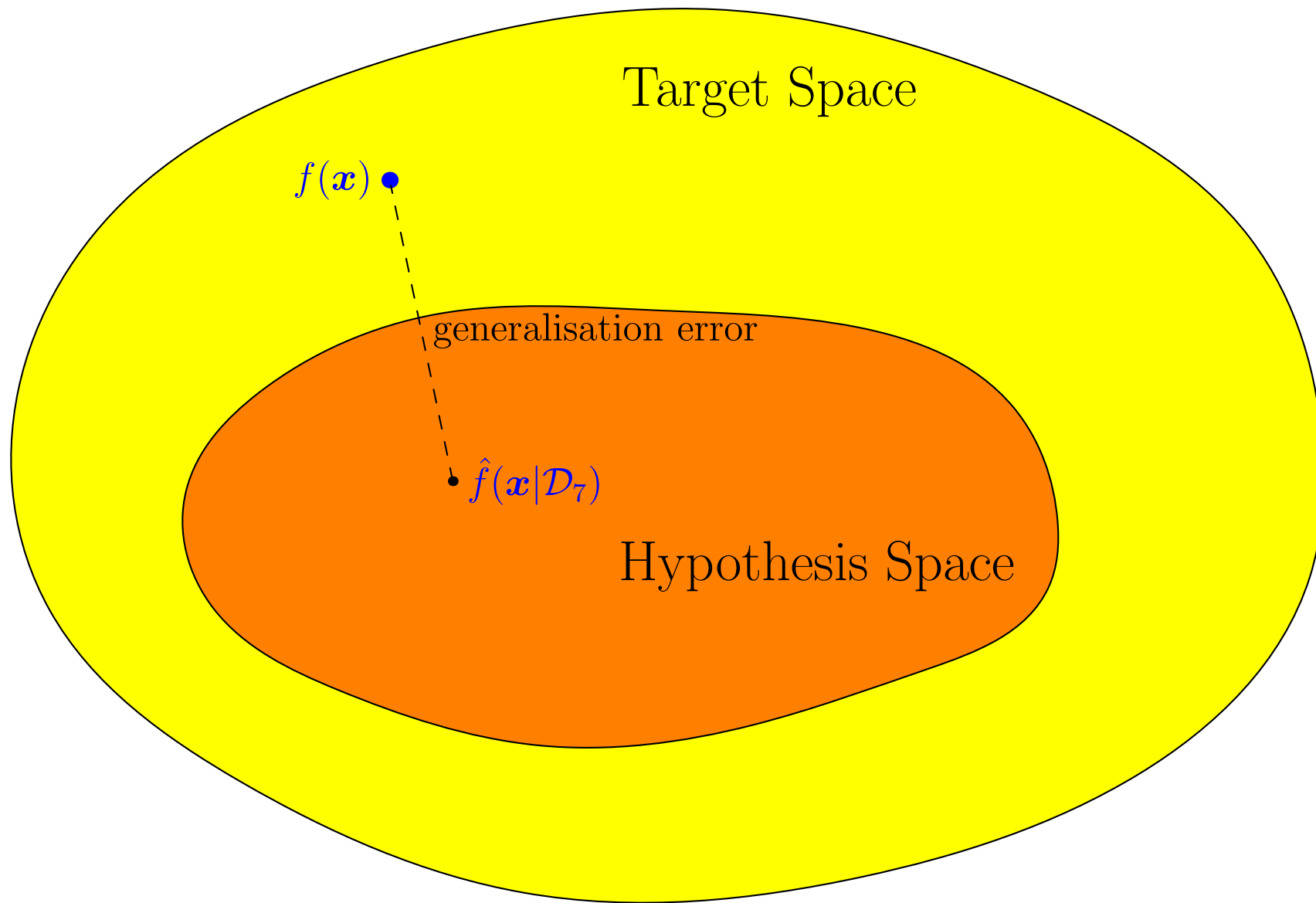
# Approximation and Estimation Errors



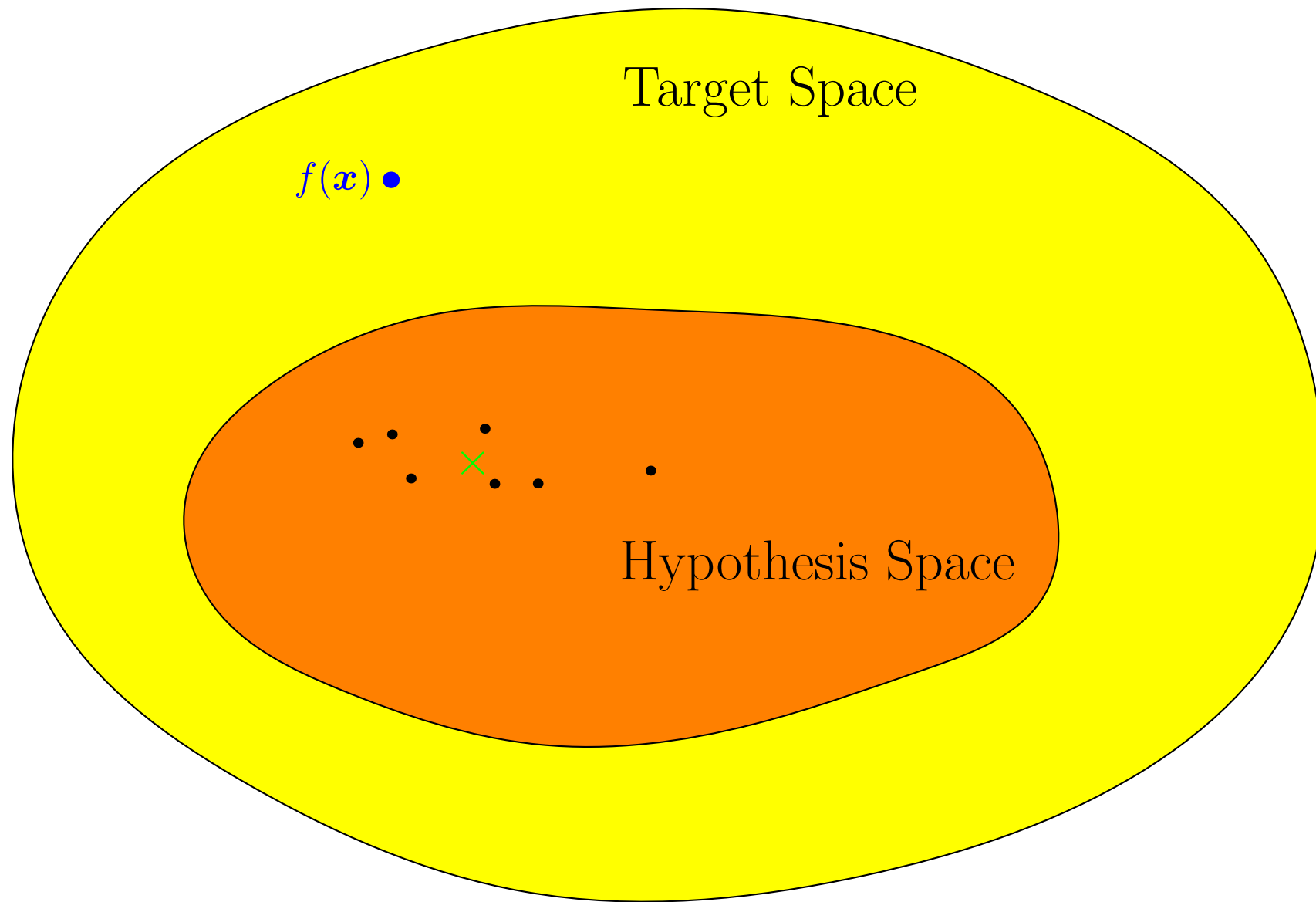
# Approximation and Estimation Errors



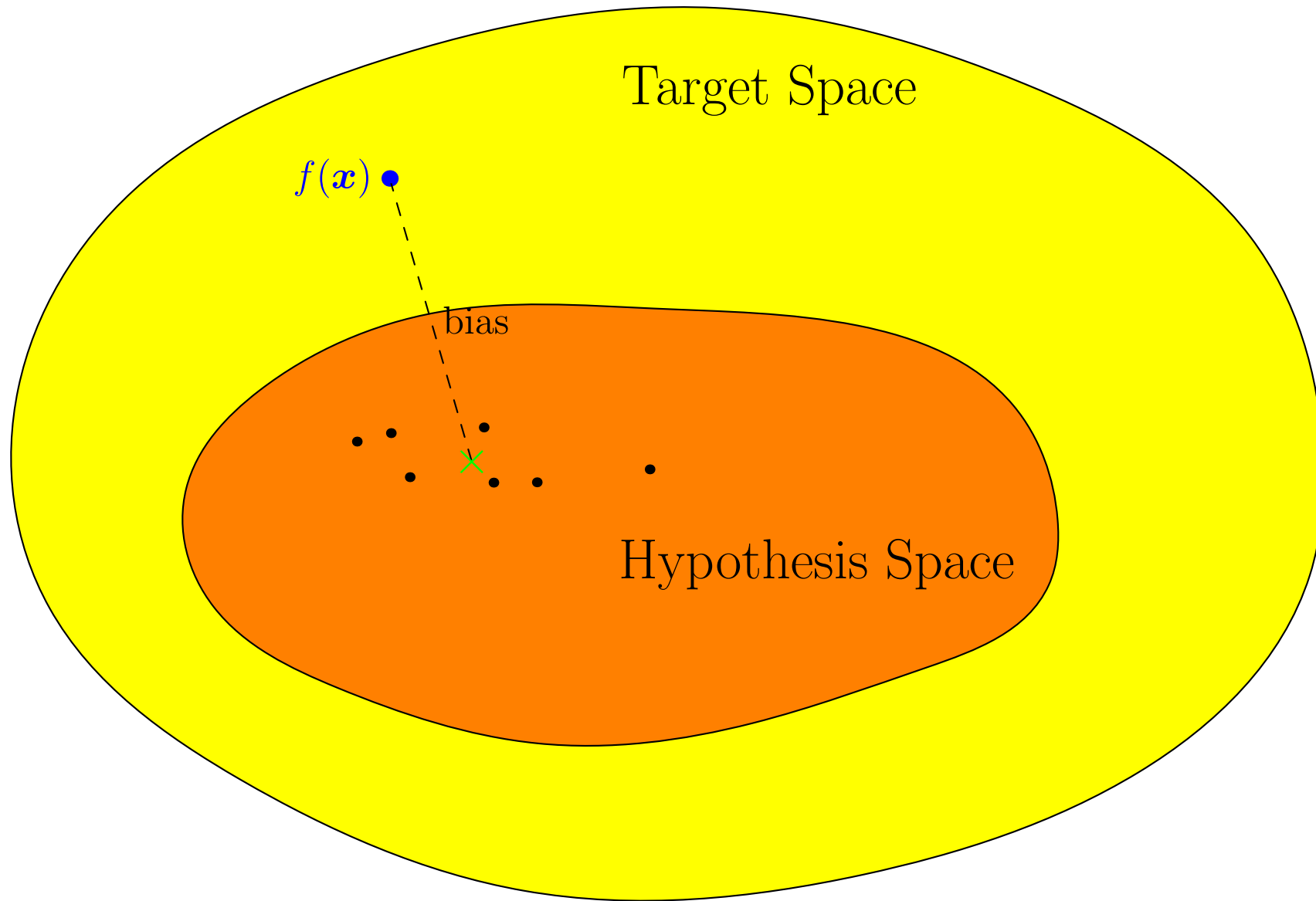
# Approximation and Estimation Errors



# Approximation and Estimation Errors

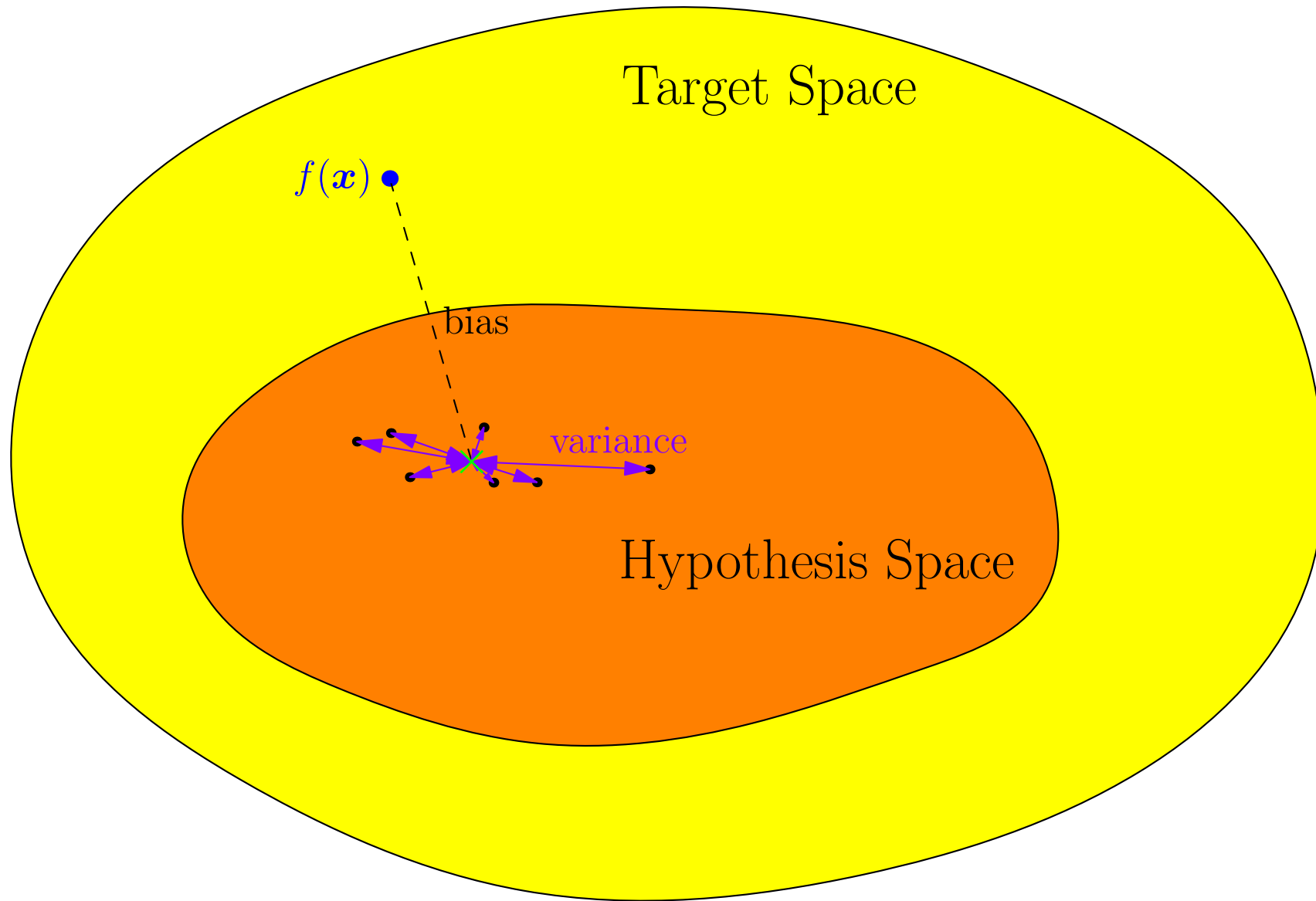


# Approximation and Estimation Errors





# Approximation and Estimation Errors



# Mean Machine

- To help understand generalisation we can consider the mean prediction with respect to machines trained with all data sets of size  $m$

$$\hat{f}_m(\mathbf{x}) = \mathbb{E}_{\mathcal{D}} \left[ \hat{f}(\mathbf{x}|\mathcal{D}) \right]$$

- We can define the **bias** to be generalisation performance of the mean machine

$$B = \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \left( \hat{f}_m(\mathbf{x}) - f(\mathbf{x}) \right)^2$$

# Mean Machine

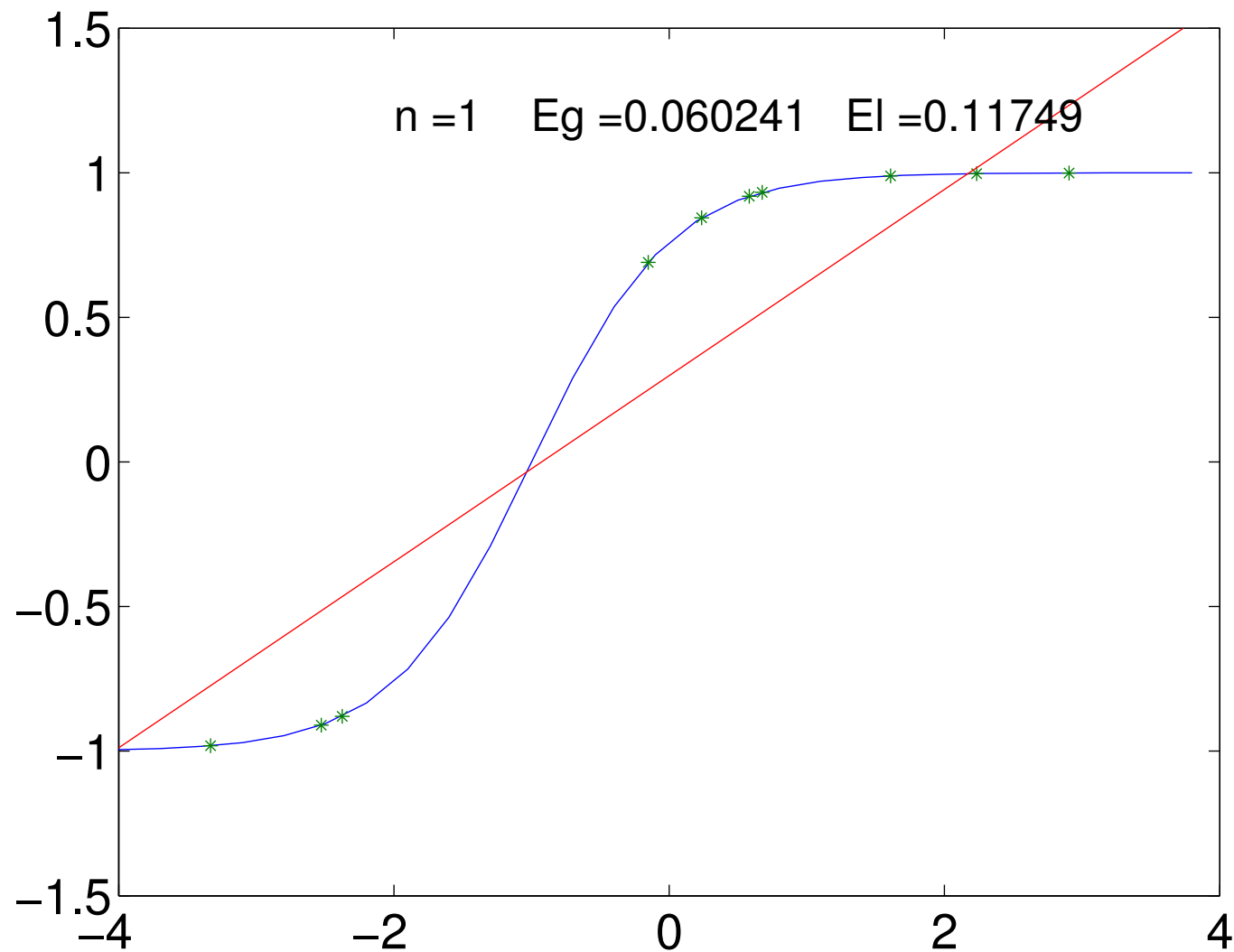
- To help understand generalisation we can consider the mean prediction with respect to machines trained with all data sets of size  $m$

$$\hat{f}_m(\mathbf{x}) = \mathbb{E}_{\mathcal{D}} \left[ \hat{f}(\mathbf{x}|\mathcal{D}) \right]$$

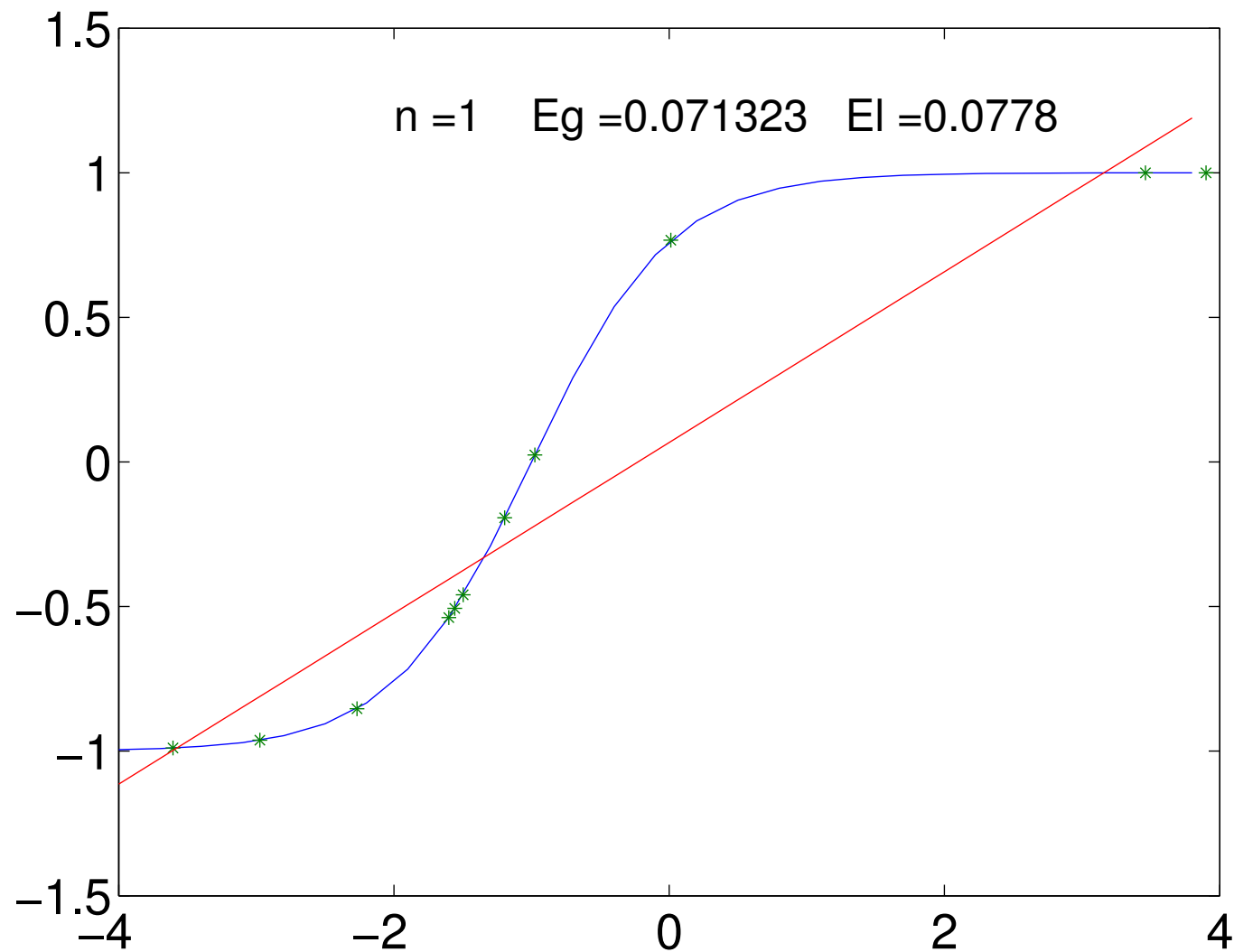
- We can define the **bias** to be generalisation performance of the mean machine

$$B = \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \left( \hat{f}_m(\mathbf{x}) - f(\mathbf{x}) \right)^2$$

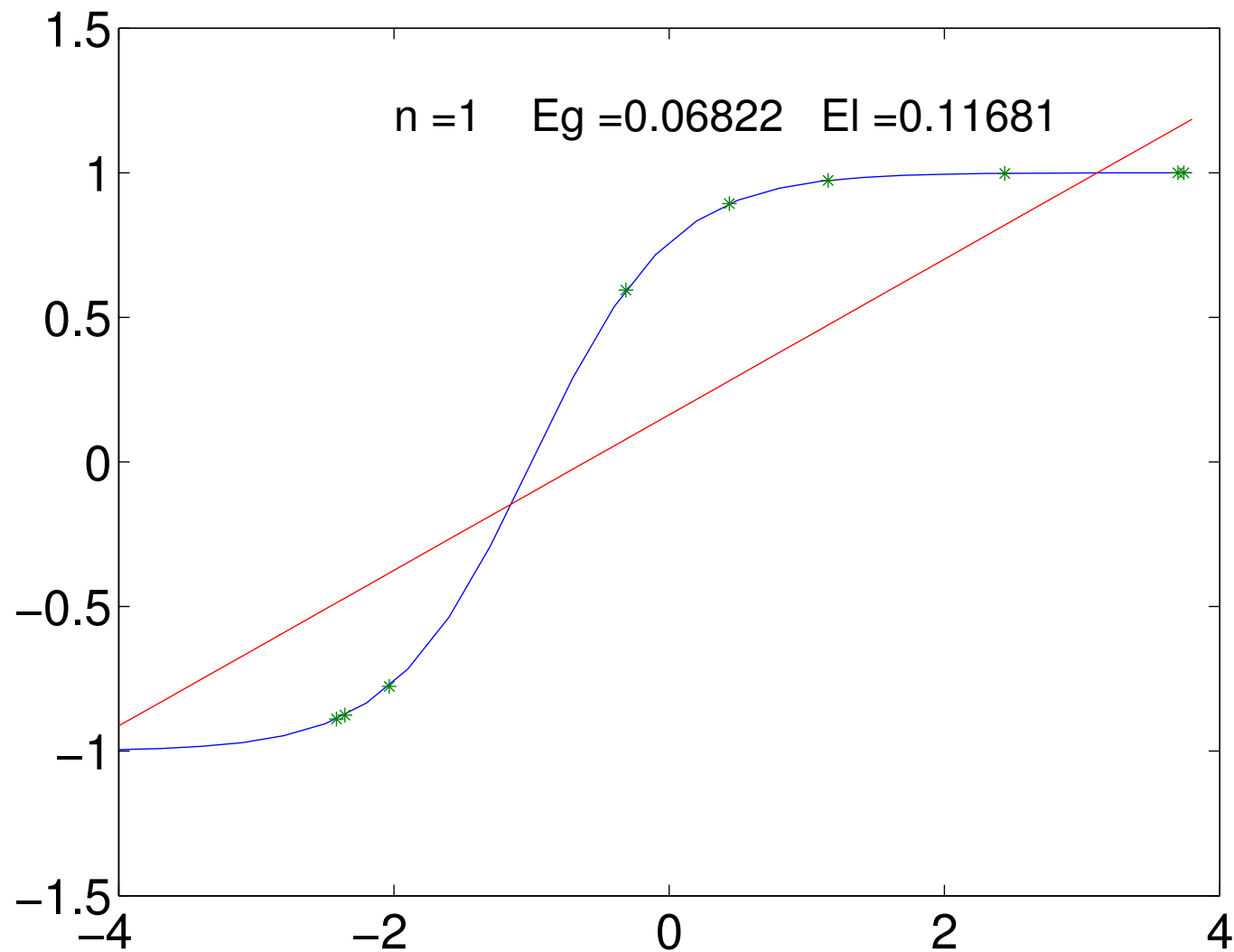
# Regression Example $n = 1$



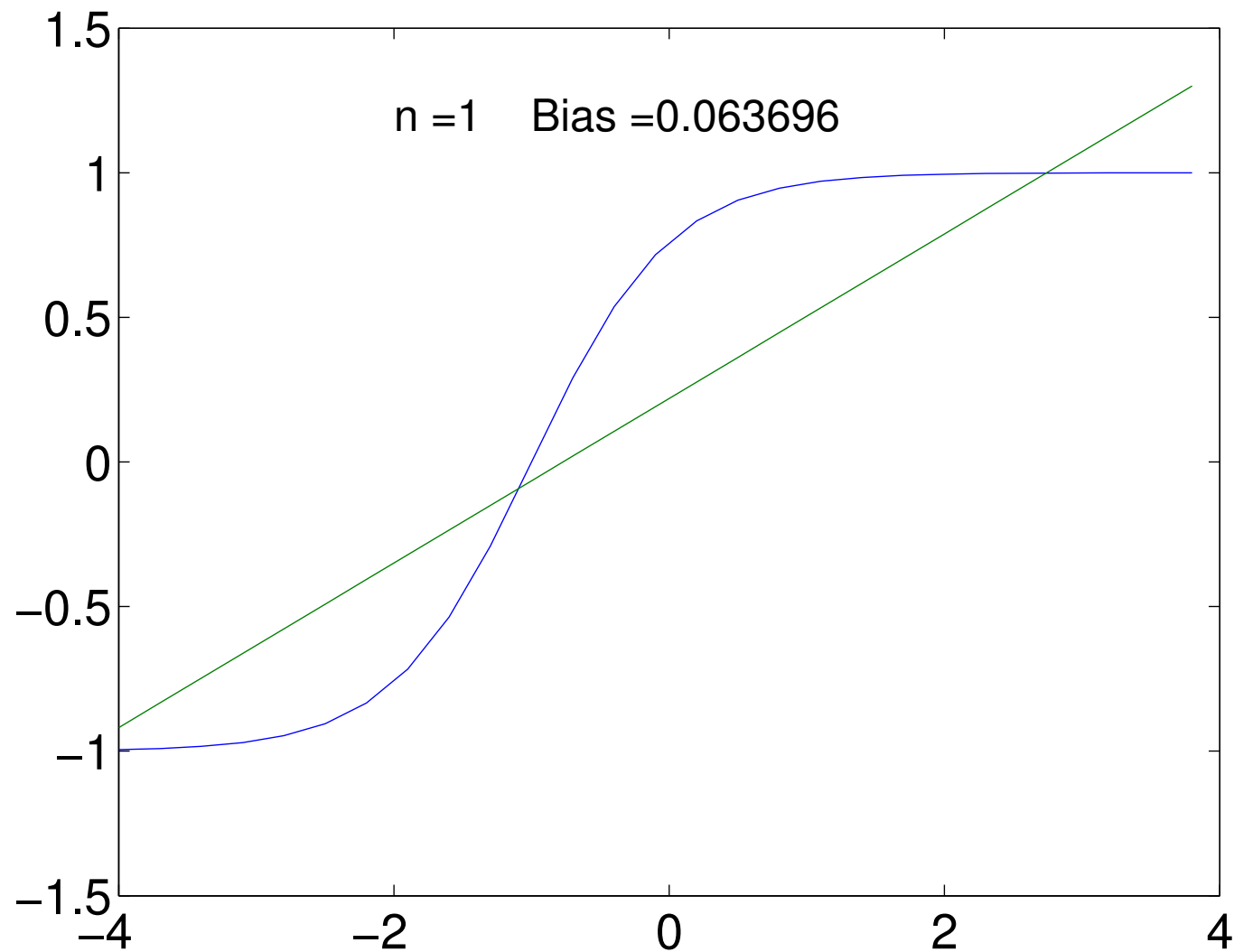
# Regression Example $n = 1$



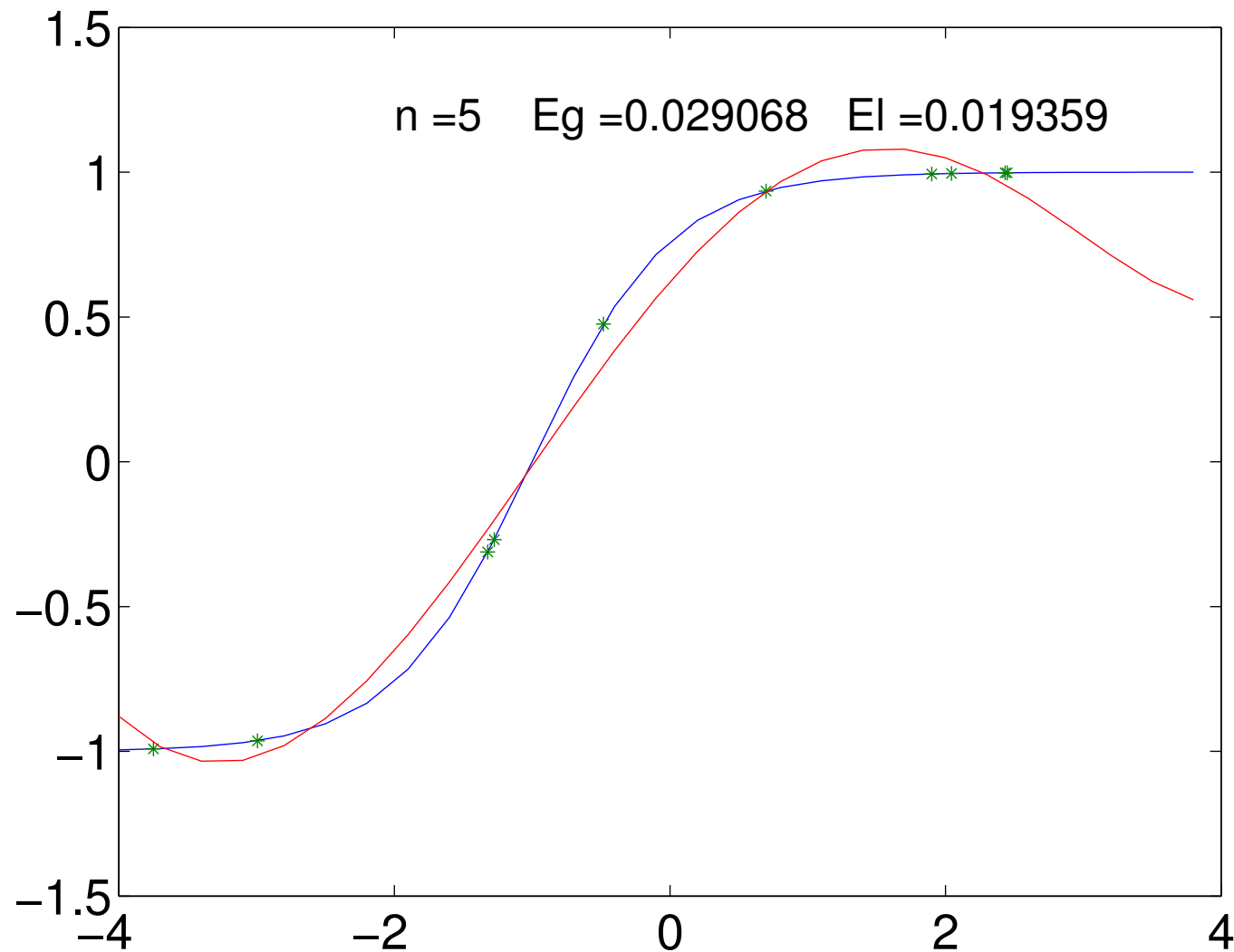
# Regression Example $n = 1$



# Regression Example $n = 1$

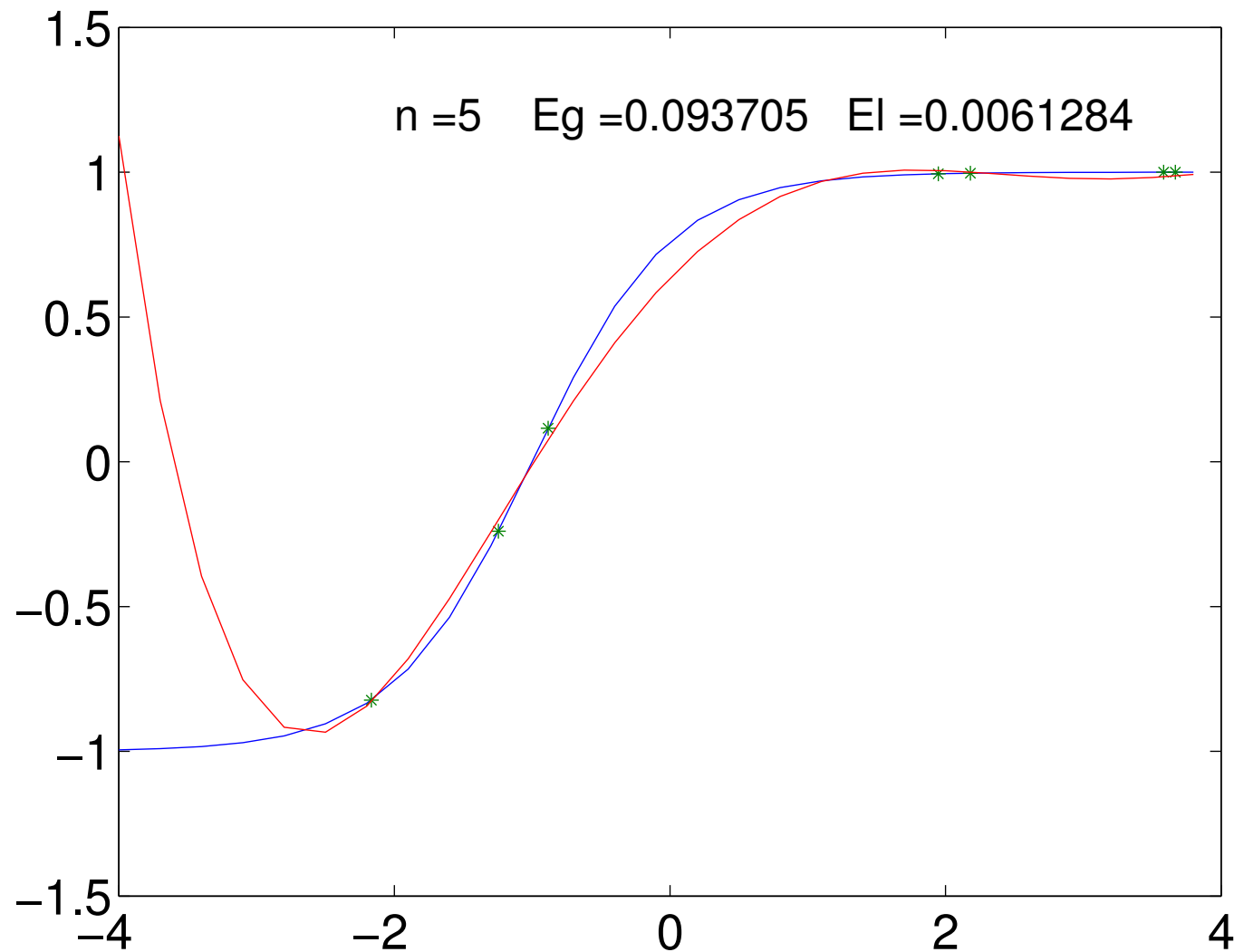


# Regression Example $n = 5$

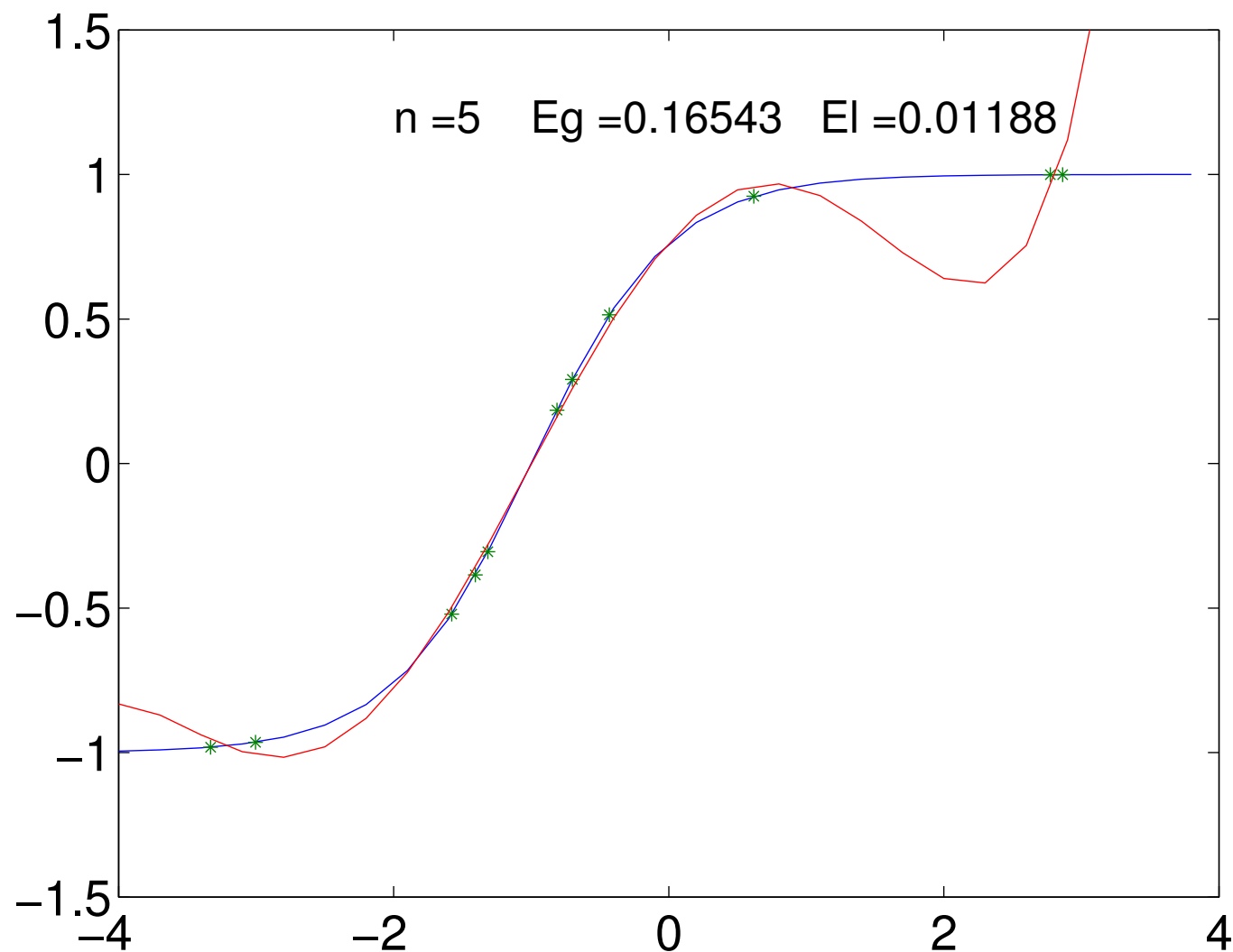




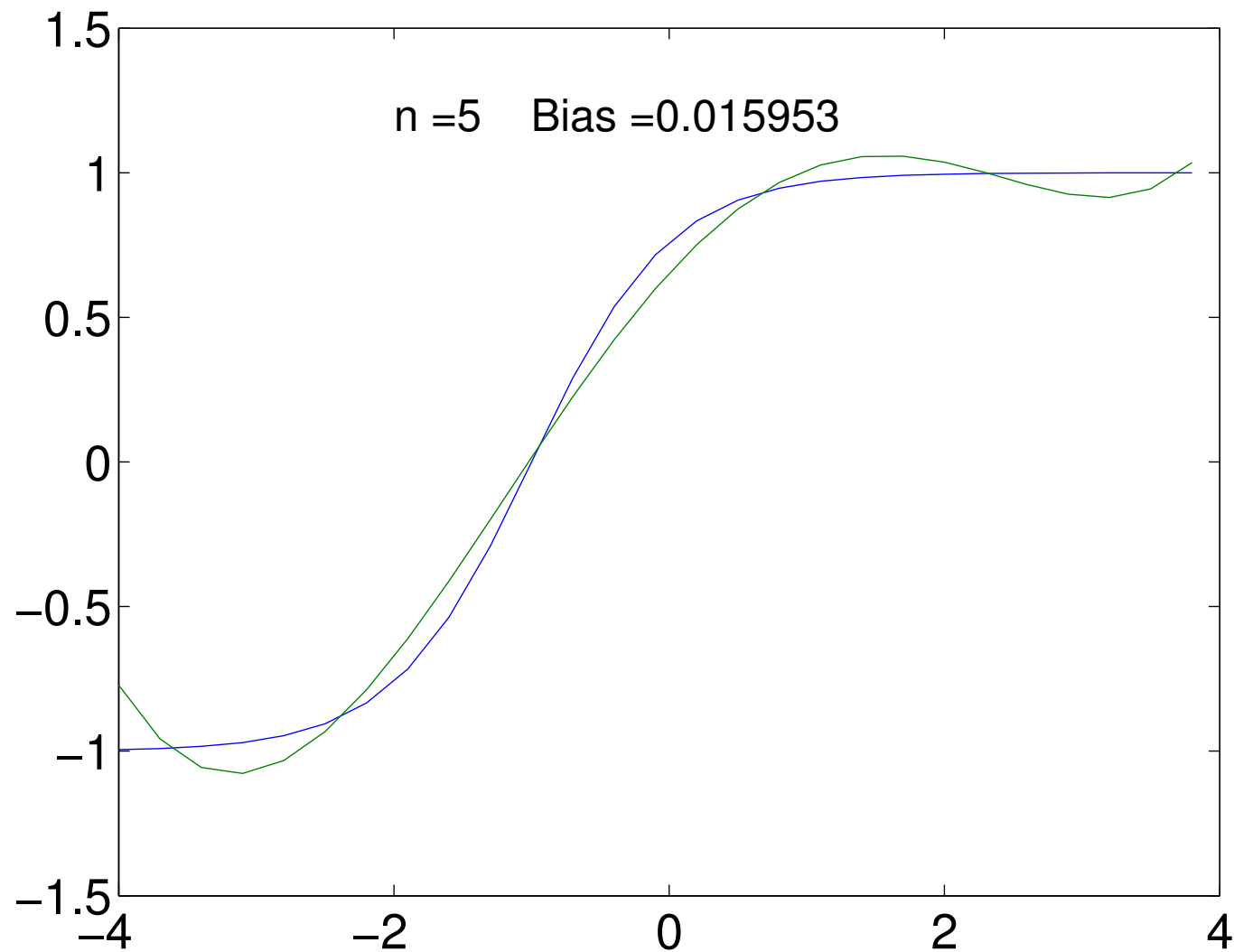
# Regression Example $n = 5$



# Regression Example $n = 5$



# Regression Example $n = 5$



# Bias and Variance

Consider the expected generalisation for data sets of size  $|\mathcal{D}| = m$

$$\bar{E}_G = \mathbb{E}_{\mathcal{D}}[E_G(\mathcal{D})]$$

# Bias and Variance

Consider the expected generalisation for data sets of size  $|\mathcal{D}| = m$

$$\bar{E}_G = \mathbb{E}_{\mathcal{D}}[E_G(\mathcal{D})] = \mathbb{E}_{\mathcal{D}} \left[ \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \left( \hat{f}(\mathbf{x}|\mathcal{D}) - f(\mathbf{x}) \right)^2 \right]$$

# Bias and Variance

Consider the expected generalisation for data sets of size  $|\mathcal{D}| = m$

$$\begin{aligned}\bar{E}_G &= \mathbb{E}_{\mathcal{D}}[E_G(\mathcal{D})] = \mathbb{E}_{\mathcal{D}} \left[ \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \left( \hat{f}(\mathbf{x}|\mathcal{D}) - f(\mathbf{x}) \right)^2 \right] \\ &= \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \mathbb{E}_{\mathcal{D}} \left[ \left( \hat{f}(\mathbf{x}|\mathcal{D}) - f(\mathbf{x}) \right)^2 \right]\end{aligned}$$

# Bias and Variance

Consider the expected generalisation for data sets of size  $|\mathcal{D}| = m$

$$\begin{aligned}\bar{E}_G &= \mathbb{E}_{\mathcal{D}}[E_G(\mathcal{D})] = \mathbb{E}_{\mathcal{D}} \left[ \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \left( \hat{f}(\mathbf{x}|\mathcal{D}) - f(\mathbf{x}) \right)^2 \right] \\ &= \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \mathbb{E}_{\mathcal{D}} \left[ \left( \hat{f}(\mathbf{x}|\mathcal{D}) - f(\mathbf{x}) \right)^2 \right] \\ &= \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \mathbb{E}_{\mathcal{D}} \left[ \left( \hat{f}(\mathbf{x}|\mathcal{D}) - \hat{f}_m(\mathbf{x}) + \hat{f}_m(\mathbf{x}) - f(\mathbf{x}) \right)^2 \right]\end{aligned}$$

# Bias and Variance

Consider the expected generalisation for data sets of size  $|\mathcal{D}| = m$

$$\begin{aligned}\bar{E}_G &= \mathbb{E}_{\mathcal{D}}[E_G(\mathcal{D})] = \mathbb{E}_{\mathcal{D}} \left[ \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \left( \hat{f}(\mathbf{x}|\mathcal{D}) - f(\mathbf{x}) \right)^2 \right] \\ &= \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \mathbb{E}_{\mathcal{D}} \left[ \left( \hat{f}(\mathbf{x}|\mathcal{D}) - f(\mathbf{x}) \right)^2 \right] \\ &= \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \mathbb{E}_{\mathcal{D}} \left[ \left( \left( \hat{f}(\mathbf{x}|\mathcal{D}) - \hat{f}_m(\mathbf{x}) \right) + \left( \hat{f}_m(\mathbf{x}) - f(\mathbf{x}) \right) \right)^2 \right]\end{aligned}$$



# Bias and Variance

Consider the expected generalisation for data sets of size  $|\mathcal{D}| = m$

$$\begin{aligned}\bar{E}_G &= \mathbb{E}_{\mathcal{D}}[E_G(\mathcal{D})] = \mathbb{E}_{\mathcal{D}} \left[ \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \left( \hat{f}(\mathbf{x}|\mathcal{D}) - f(\mathbf{x}) \right)^2 \right] \\&= \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \mathbb{E}_{\mathcal{D}} \left[ \left( \hat{f}(\mathbf{x}|\mathcal{D}) - f(\mathbf{x}) \right)^2 \right] \\&= \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \mathbb{E}_{\mathcal{D}} \left[ \left( \left( \hat{f}(\mathbf{x}|\mathcal{D}) - \hat{f}_m(\mathbf{x}) \right) + \left( \hat{f}_m(\mathbf{x}) - f(\mathbf{x}) \right) \right)^2 \right] \\&= \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \left( \mathbb{E}_{\mathcal{D}} \left[ \left( \hat{f}(\mathbf{x}|\mathcal{D}) - \hat{f}_m(\mathbf{x}) \right)^2 \right] + \left( \hat{f}_m(\mathbf{x}) - f(\mathbf{x}) \right)^2 \right. \\&\quad \left. + 2 \mathbb{E}_{\mathcal{D}} \left[ \left( \hat{f}(\mathbf{x}|\mathcal{D}) - \hat{f}_m(\mathbf{x}) \right) \left( \hat{f}_m(\mathbf{x}) - f(\mathbf{x}) \right) \right] \right)\end{aligned}$$

# Cross Term

- The cross term vanishes

$$C = \mathbb{E}_{\mathcal{D}} \left[ \left( \hat{f}(\mathbf{x}|\mathcal{D}) - \hat{f}_m(\mathbf{x}) \right) \left( \hat{f}_m(\mathbf{x}) - f(\mathbf{x}) \right) \right]$$

# Cross Term

- The cross term vanishes

$$\begin{aligned} C &= \mathbb{E}_{\mathcal{D}} \left[ \left( \hat{f}(\mathbf{x}|\mathcal{D}) - \hat{f}_m(\mathbf{x}) \right) \left( \hat{f}_m(\mathbf{x}) - f(\mathbf{x}) \right) \right] \\ &= \left( \mathbb{E}_{\mathcal{D}} \left[ \hat{f}(\mathbf{x}|\mathcal{D}) \right] - \hat{f}_m(\mathbf{x}) \right) \left( \hat{f}_m(\mathbf{x}) - f(\mathbf{x}) \right) \end{aligned}$$

# Cross Term

- The cross term vanishes

$$\begin{aligned} C &= \mathbb{E}_{\mathcal{D}} \left[ \left( \hat{f}(\mathbf{x}|\mathcal{D}) - \hat{f}_m(\mathbf{x}) \right) \left( \hat{f}_m(\mathbf{x}) - f(\mathbf{x}) \right) \right] \\ &= \left( \mathbb{E}_{\mathcal{D}} \left[ \hat{f}(\mathbf{x}|\mathcal{D}) \right] - \hat{f}_m(\mathbf{x}) \right) \left( \hat{f}_m(\mathbf{x}) - f(\mathbf{x}) \right) \\ &= \left( \hat{f}_m(\mathbf{x}) - \hat{f}_m(\mathbf{x}) \right) \left( \hat{f}_m(\mathbf{x}) - f(\mathbf{x}) \right) \end{aligned}$$

# Cross Term

- The cross term vanishes

$$\begin{aligned} C &= \mathbb{E}_{\mathcal{D}} \left[ \left( \hat{f}(\mathbf{x}|\mathcal{D}) - \hat{f}_m(\mathbf{x}) \right) \left( \hat{f}_m(\mathbf{x}) - f(\mathbf{x}) \right) \right] \\ &= \left( \mathbb{E}_{\mathcal{D}} \left[ \hat{f}(\mathbf{x}|\mathcal{D}) \right] - \hat{f}_m(\mathbf{x}) \right) \left( \hat{f}_m(\mathbf{x}) - f(\mathbf{x}) \right) \\ &= \left( \hat{f}_m(\mathbf{x}) - \hat{f}_m(\mathbf{x}) \right) \left( \hat{f}_m(\mathbf{x}) - f(\mathbf{x}) \right) \\ &= 0 \end{aligned}$$

# Cross Term

- The cross term vanishes

$$\begin{aligned} C &= \mathbb{E}_{\mathcal{D}} \left[ \left( \hat{f}(\mathbf{x}|\mathcal{D}) - \hat{f}_m(\mathbf{x}) \right) \left( \hat{f}_m(\mathbf{x}) - f(\mathbf{x}) \right) \right] \\ &= \left( \mathbb{E}_{\mathcal{D}} \left[ \hat{f}(\mathbf{x}|\mathcal{D}) \right] - \hat{f}_m(\mathbf{x}) \right) \left( \hat{f}_m(\mathbf{x}) - f(\mathbf{x}) \right) \\ &= \left( \hat{f}_m(\mathbf{x}) - \hat{f}_m(\mathbf{x}) \right) \left( \hat{f}_m(\mathbf{x}) - f(\mathbf{x}) \right) \\ &= 0 \end{aligned}$$

- Thus

$$\bar{E}_G = \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \mathbb{E}_{\mathcal{D}} \left[ \left( \hat{f}(\mathbf{x}|\mathcal{D}) - \hat{f}_m(\mathbf{x}) \right)^2 + \left( \hat{f}_m(\mathbf{x}) - f(\mathbf{x}) \right)^2 \right]$$

# Bias and Variance

- We can write the expected generalisation as

$$\begin{aligned}\mathbb{E}_{\mathcal{D}}[E_G(\mathcal{D})] &= \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \mathbb{E}_{\mathcal{D}} \left[ \left( \hat{f}(\mathbf{x}|\mathcal{D}) - \hat{f}_m(\mathbf{x}) \right)^2 \right] \\ &\quad + \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \left( \hat{f}_m(\mathbf{x}) - f(\mathbf{x}) \right)^2 = V + B\end{aligned}$$

- Where  $B$  is the bias and  $V$  is the variance defined by

$$V = \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \mathbb{E}_{\mathcal{D}} \left[ \left( \hat{f}(\mathbf{x}|\mathcal{D}) - \hat{f}_m(\mathbf{x}) \right)^2 \right]$$

# Bias and Variance

- We can write the expected generalisation as

$$\begin{aligned}\mathbb{E}_{\mathcal{D}}[E_G(\mathcal{D})] &= \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \mathbb{E}_{\mathcal{D}} \left[ \left( \hat{f}(\mathbf{x}|\mathcal{D}) - \hat{f}_m(\mathbf{x}) \right)^2 \right] \\ &\quad + \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \left( \hat{f}_m(\mathbf{x}) - f(\mathbf{x}) \right)^2 = V + B\end{aligned}$$

- Where  $B$  is the bias and  $V$  is the variance defined by

$$V = \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \mathbb{E}_{\mathcal{D}} \left[ \left( \hat{f}(\mathbf{x}|\mathcal{D}) - \hat{f}_m(\mathbf{x}) \right)^2 \right]$$



# Bias and Variance

- We can write the expected generalisation as

$$\begin{aligned}\mathbb{E}_{\mathcal{D}}[E_G(\mathcal{D})] &= \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \mathbb{E}_{\mathcal{D}} \left[ \left( \hat{f}(\mathbf{x}|\mathcal{D}) - \hat{f}_m(\mathbf{x}) \right)^2 \right] \\ &\quad + \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \left( \hat{f}_m(\mathbf{x}) - f(\mathbf{x}) \right)^2 = V + B\end{aligned}$$

- Where  $B$  is the bias and  $V$  is the variance defined by

$$V = \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \mathbb{E}_{\mathcal{D}} \left[ \left( \hat{f}(\mathbf{x}|\mathcal{D}) - \hat{f}_m(\mathbf{x}) \right)^2 \right]$$

# Bias-Variance Dilemma

- The bias measure the generalisation performance of the *mean machine* and is large if the machine is too simple to capture the changes in the function we want to learn
- The variance measures the variation in the prediction of the machine as we change the data set we train on

$$V = \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \mathbb{E}_{\mathcal{D}} \left[ \left( \hat{f}(\mathbf{x}|\mathcal{D}) - \hat{f}_m(\mathbf{x}) \right)^2 \right]$$

- The variance is usually large if we have a complex machine
- Striking the right balance is often the key to getting good results

# Bias-Variance Dilemma

- The bias measure the generalisation performance of the *mean machine* and is large if the machine is too simple to capture the changes in the function we want to learn
- The variance measures the variation in the prediction of the machine as we change the data set we train on

$$V = \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \mathbb{E}_{\mathcal{D}} \left[ \left( \hat{f}(\mathbf{x}|\mathcal{D}) - \hat{f}_m(\mathbf{x}) \right)^2 \right]$$

- The variance is usually large if we have a complex machine
- Striking the right balance is often the key to getting good results

# Bias-Variance Dilemma

- The bias measure the generalisation performance of the *mean machine* and is large if the machine is too simple to capture the changes in the function we want to learn
- The variance measures the variation in the prediction of the machine as we change the data set we train on

$$V = \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \mathbb{E}_{\mathcal{D}} \left[ \left( \hat{f}(\mathbf{x}|\mathcal{D}) - \hat{f}_m(\mathbf{x}) \right)^2 \right]$$

- The variance is usually large if we have a complex machine
- Striking the right balance is often the key to getting good results

# Bias-Variance Dilemma

- The bias measure the generalisation performance of the *mean machine* and is large if the machine is too simple to capture the changes in the function we want to learn
- The variance measures the variation in the prediction of the machine as we change the data set we train on

$$V = \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \mathbb{E}_{\mathcal{D}} \left[ \left( \hat{f}(\mathbf{x}|\mathcal{D}) - \hat{f}_m(\mathbf{x}) \right)^2 \right]$$

- The variance is usually large if we have a complex machine
- Striking the right balance is often the key to getting good results

# Balancing Bias and Variance

- We want to choose a learning machine that is complex enough to capture the underlying function we are trying to learn, but otherwise as simple as possible
- There are a number of tricks to achieve this balance
- Some require us to preprocess the data to reduce the number of inputs
- Some machines cleverly adjust their own complexity
- Today we will look at machines that achieve this balance

# Balancing Bias and Variance

- We want to choose a learning machine that is complex enough to capture the underlying function we are trying to learn, but otherwise as simple as possible
- There are a number of tricks to achieve this balance
- Some require us to preprocess the data to reduce the number of inputs
- Some machines cleverly adjust their own complexity
- Today we will look at machines that achieve this balance

# Balancing Bias and Variance

- We want to choose a learning machine that is complex enough to capture the underlying function we are trying to learn, but otherwise as simple as possible
- There are a number of tricks to achieve this balance
- Some require us to preprocess the data to reduce the number of inputs
- Some machines cleverly adjust their own complexity
- Today we will look at machines that achieve this balance



# Balancing Bias and Variance

- We want to choose a learning machine that is complex enough to capture the underlying function we are trying to learn, but otherwise as simple as possible
- There are a number of tricks to achieve this balance
- Some require us to preprocess the data to reduce the number of inputs
- Some machines cleverly adjust their own complexity
- Today we will look at machines that achieve this balance

# Balancing Bias and Variance

- We want to choose a learning machine that is complex enough to capture the underlying function we are trying to learn, but otherwise as simple as possible
- There are a number of tricks to achieve this balance
- Some require us to preprocess the data to reduce the number of inputs
- Some machines cleverly adjust their own complexity
- Today we will look at machines that achieve this balance

# Over-fitting

- Complex machine can **over-fitting**

***over-fitting:** fitting the training data well at the cost of getting poorer generalisation performance*

- Three red cars. . .
- Models that fit the training data are necessarily predictive

# Over-fitting

- Complex machine can **over-fitting**

***over-fitting:** fitting the training data well at the cost of getting poorer generalisation performance*

- Three red cars. . .
- Models that fit the training data are necessarily predictive

# Over-fitting

- Complex machine can **over-fitting**

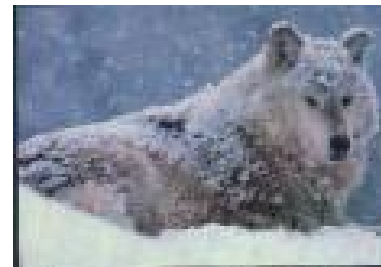
***over-fitting:** fitting the training data well at the cost of getting poorer generalisation performance*

- Three red cars. . .
- Models that fit the training data are necessarily predictive

# Binary Classification Task for You



Class 1



Class 2

# Which Category?

- Which category does the following image belong to?



# Spurious Rules

- You ask a learning machine to solve a task based on data
- It will find a rule that does this, but not necessary the rule you had in mind
- Infinitely flexible machines have an infinity of spurious rules they can learn



# Spurious Rules

- You ask a learning machine to solve a task based on data
- It will find a rule that does this, but not necessary the rule you had in mind
- Infinitely flexible machines have an infinity of spurious rules they can learn

# Spurious Rules

- You ask a learning machine to solve a task based on data
- It will find a rule that does this, but not necessary the rule you had in mind—machine learning isn't magic, it can't read your mind
- Infinitely flexible machines have an infinity of spurious rules they can learn

# Spurious Rules

- You ask a learning machine to solve a task based on data
- It will find a rule that does this, but not necessary the rule you had in mind—machine learning isn't magic, it can't read your mind
- Infinitely flexible machines have an infinity of spurious rules they can learn

# Spurious Rules

- You ask a learning machine to solve a task based on data
- It will find a rule that does this, but not necessary the rule you had in mind—machine learning isn't magic, it can't read your mind
- Infinitely flexible machines have an infinity of spurious rules they can learn—they are useless

# All Binary Functions

$$\mathbf{x}_0 = 000 \quad y_0 = \begin{cases} 0 \\ 1 \end{cases}$$

$$\mathbf{x}_1 = 100 \quad y_1 = \begin{cases} 0 \\ 1 \end{cases}$$

$$\mathbf{x}_2 = 010 \quad y_2 = \begin{cases} 0 \\ 1 \end{cases}$$

$$\mathbf{x}_3 = 110 \quad y_3 = \begin{cases} 0 \\ 1 \end{cases}$$

$$\mathbf{x}_4 = 001 \quad y_4 = \begin{cases} 0 \\ 1 \end{cases}$$

$$\mathbf{x}_5 = 101 \quad y_5 = \begin{cases} 0 \\ 1 \end{cases}$$

$$\mathbf{x}_6 = 011 \quad y_6 = \begin{cases} 0 \\ 1 \end{cases}$$

$$\mathbf{x}_7 = 111 \quad y_7 = \begin{cases} 0 \\ 1 \end{cases}$$

# All Binary Functions

$$\mathbf{x}_0 = 000 \quad y_0 = \begin{cases} 0 \\ \textcolor{blue}{\times} \end{cases}$$

$$\mathbf{x}_1 = 100 \quad y_1 = \begin{cases} 0 \\ 1 \end{cases}$$

$$\mathbf{x}_2 = 010 \quad y_2 = \begin{cases} \textcolor{blue}{\times} \\ 1 \end{cases}$$

$$\mathbf{x}_3 = 110 \quad y_3 = \begin{cases} \textcolor{blue}{\times} \\ 1 \end{cases}$$

$$\mathbf{x}_4 = 001 \quad y_4 = \begin{cases} 0 \\ \textcolor{blue}{\times} \end{cases}$$

$$\mathbf{x}_5 = 101 \quad y_5 = \begin{cases} 0 \\ \textcolor{blue}{\times} \end{cases}$$

$$\mathbf{x}_6 = 011 \quad y_6 = \begin{cases} 0 \\ 1 \end{cases}$$

$$\mathbf{x}_7 = 111 \quad y_7 = \begin{cases} 0 \\ 1 \end{cases}$$

$$\mathcal{D} = \{(000, 0), (010, 1), (110, 1), (001, 0), (101, 0)\}$$

# All Binary Functions

$$\mathbf{x}_0 = 000 \quad y_0 = \begin{cases} 0 \\ \textcolor{blue}{\times} \end{cases}$$

$$\mathbf{x}_1 = 100 \quad y_1 = \begin{cases} 0 \\ 1 \end{cases} \quad \textcolor{red}{\text{unseen}}$$

$$\mathbf{x}_2 = 010 \quad y_2 = \begin{cases} \textcolor{blue}{\times} \\ 1 \end{cases}$$

$$\mathbf{x}_3 = 110 \quad y_3 = \begin{cases} \textcolor{blue}{\times} \\ 1 \end{cases}$$

$$\mathbf{x}_4 = 001 \quad y_4 = \begin{cases} 0 \\ \textcolor{blue}{\times} \end{cases}$$

$$\mathbf{x}_5 = 101 \quad y_5 = \begin{cases} 0 \\ \textcolor{blue}{\times} \end{cases}$$

$$\mathbf{x}_6 = 011 \quad y_6 = \begin{cases} 0 \\ 1 \end{cases}$$

$$\mathbf{x}_7 = 111 \quad y_7 = \begin{cases} 0 \\ 1 \end{cases}$$

$$\mathcal{D} = \{(000, 0), (010, 1), (110, 1), (001, 0), (101, 0)\}$$

# All Binary Functions

$$\mathbf{x}_0 = 000 \quad y_0 = \begin{cases} 0 \\ \text{✖} \end{cases}$$

$$\mathbf{x}_1 = 100 \quad y_1 = \begin{cases} 0 \\ 1 \end{cases}$$

$$\mathbf{x}_2 = 010 \quad y_2 = \begin{cases} \text{✖} \\ 1 \end{cases}$$

$$\mathbf{x}_3 = 110 \quad y_3 = \begin{cases} \text{✖} \\ 1 \end{cases}$$

$$\mathbf{x}_4 = 001 \quad y_4 = \begin{cases} 0 \\ \text{✖} \end{cases}$$

$$\mathbf{x}_5 = 101 \quad y_5 = \begin{cases} 0 \\ \text{✖} \end{cases}$$

$$\mathbf{x}_6 = 011 \quad y_6 = \begin{cases} 0 \\ 1 \end{cases}$$

$$\mathbf{x}_7 = 111 \quad y_7 = \begin{cases} 0 \\ 1 \end{cases} \quad \text{unseen}$$

$$\mathcal{D} = \{(000, 0), (010, 1), (110, 1), (001, 0), (101, 0)\}$$



# All Binary Functions

$$\mathbf{x}_0 = 000 \quad y_0 = \begin{cases} 0 \\ \textcolor{blue}{\times} \end{cases}$$

$$\mathbf{x}_1 = 100 \quad y_1 = \begin{cases} 0 \\ 1 \end{cases}$$

$$\mathbf{x}_2 = 010 \quad y_2 = \begin{cases} \textcolor{blue}{\times} \\ 1 \end{cases}$$

$$\mathbf{x}_3 = 110 \quad y_3 = \begin{cases} \textcolor{blue}{\times} \\ 1 \end{cases}$$

$$\mathbf{x}_4 = 001 \quad y_4 = \begin{cases} 0 \\ \textcolor{blue}{\times} \end{cases} \text{ seen}$$

$$\mathbf{x}_5 = 101 \quad y_5 = \begin{cases} 0 \\ \textcolor{blue}{\times} \end{cases}$$

$$\mathbf{x}_6 = 011 \quad y_6 = \begin{cases} 0 \\ 1 \end{cases}$$

$$\mathbf{x}_7 = 111 \quad y_7 = \begin{cases} 0 \\ 1 \end{cases}$$

$$\mathcal{D} = \{(000, 0), (010, 1), (110, 1), (001, 0), (101, 0)\}$$

# Controlling Complexity

- Infinitely flexible machine don't generalise
- Machine learning only works because there is some structure in the data
- A successful machine should capture this structure
- Even deep learning machines with millions of parameters only work because they successfully capture the structure of images or text
- Different learning machines have different performance on different problems because the data has different structure

# Controlling Complexity

- Infinitely flexible machine don't generalise (because any unseen data could have any value)
- **Machine learning only works because there is some structure in the data**
- A successful machine should capture this structure
- Even deep learning machines with millions of parameters only work because they successfully capture the structure of images or text
- Different learning machines have different performance on different problems because the data has different structure

# Controlling Complexity

- Infinitely flexible machine don't generalise (because any unseen data could have any value)
- **Machine learning only works because there is some structure in the data**
- A successful machine should capture this structure
- Even deep learning machines with millions of parameters only work because they successfully capture the structure of images or text
- Different learning machines have different performance on different problems because the data has different structure

# Controlling Complexity

- Infinitely flexible machine don't generalise (because any unseen data could have any value)
- **Machine learning only works because there is some structure in the data**
- A successful machine should capture this structure
- Even deep learning machines with millions of parameters only work because they successfully capture the structure of images or text
- Different learning machines have different performance on different problems because the data has different structure

# Controlling Complexity

- Infinitely flexible machine don't generalise (because any unseen data could have any value)
- **Machine learning only works because there is some structure in the data**
- A successful machine should capture this structure
- Even deep learning machines with millions of parameters only work because they successfully capture the structure of images or text
- Different learning machines have different performance on different problems because the data has different structure

# Controlling Complexity

- Infinitely flexible machine don't generalise (because any unseen data could have any value)
- **Machine learning only works because there is some structure in the data**
- A successful machine should capture this structure
- Even deep learning machines with millions of parameters only work because they successfully capture the structure of images or text
- Different learning machines have different performance on different problems because the data has different structure

# Training Examples

- As we increase the number of training examples, we make it hard to find a spurious rule
- Bigger data sets allow us to use more complicated machines
- Part of the success of deep learning is because they use huge training sets
- (Labelled) data is often expensive to collect so we sometimes have no choice but to use a small training set
- One of the limitations of using deep learning comes because we often have limited data



# Training Examples

- As we increase the number of training examples, we make it hard to find a spurious rule
- Bigger data sets allow us to use more complicated machines
- Part of the success of deep learning is because they use huge training sets
- (Labelled) data is often expensive to collect so we sometimes have no choice but to use a small training set
- One of the limitations of using deep learning comes because we often have limited data

# Training Examples

- As we increase the number of training examples, we make it hard to find a spurious rule
- Bigger data sets allow us to use more complicated machines
- Part of the success of deep learning is because they use huge training sets
- (Labelled) data is often expensive to collect so we sometimes have no choice but to use a small training set
- One of the limitations of using deep learning comes because we often have limited data

# Training Examples

- As we increase the number of training examples, we make it hard to find a spurious rule
- Bigger data sets allow us to use more complicated machines
- Part of the success of deep learning is because they use huge training sets—but this is only a part of their success
- (Labelled) data is often expensive to collect so we sometimes have no choice but to use a small training set
- One of the limitations of using deep learning comes because we often have limited data

# Training Examples

- As we increase the number of training examples, we make it hard to find a spurious rule
- Bigger data sets allow us to use more complicated machines
- Part of the success of deep learning is because they use huge training sets—but this is only a part of their success
- (Labelled) data is often expensive to collect so we sometimes have no choice but to use a small training set
- One of the limitations of using deep learning comes because we often have limited data

# Training Examples

- As we increase the number of training examples, we make it hard to find a spurious rule
- Bigger data sets allow us to use more complicated machines
- Part of the success of deep learning is because they use huge training sets—but this is only a part of their success
- (Labelled) data is often expensive to collect so we sometimes have no choice but to use a small training set
- One of the limitations of using deep learning comes because we often have limited data

# Dimensionality Reduction

- We can simplify our machines by using less features
- We can project our data onto a lower dimensional sub-space (e.g. one with the maximum variation in the data PCA)
- We can use clustering to find exemplars and recode our data in terms of differences from the exemplars (radial basis functions)
- Whether this helps depends on whether the information we discard is pertinent to the task we are trying to perform

# Dimensionality Reduction

- We can simplify our machines by using less features
- We can project our data onto a lower dimensional sub-space (e.g. one with the maximum variation in the data PCA)
- We can use clustering to find exemplars and recode our data in terms of differences from the exemplars (radial basis functions)
- Whether this helps depends on whether the information we discard is pertinent to the task we are trying to perform

# Dimensionality Reduction

- We can simplify our machines by using less features
- We can project our data onto a lower dimensional sub-space (e.g. one with the maximum variation in the data PCA)
- We can use clustering to find exemplars and recode our data in terms of differences from the exemplars (radial basis functions)
- Whether this helps depends on whether the information we discard is pertinent to the task we are trying to perform



# Dimensionality Reduction

- We can simplify our machines by using less features
- We can project our data onto a lower dimensional sub-space (e.g. one with the maximum variation in the data PCA)
- We can use clustering to find exemplars and recode our data in terms of differences from the exemplars (radial basis functions)
- Whether this helps depends on whether the information we discard is pertinent to the task we are trying to perform

# Feature Selection

- Spurious features will allow us to find spurious rules (**over-fitting**)
- We can try different combinations of features to find the best set, although it rapidly becomes intractable to do this in all ways
- We can use various heuristics to decide which features to keep
- Feature selection however can be powerful, often we can get very good results by keeping only a few variables
- As well as possibly improving generalisation we also get a more **interpretable** rule

# Feature Selection

- Spurious features will allow us to find spurious rules (**over-fitting**)
- We can try different combinations of features to find the best set, although it rapidly becomes intractable to do this in all ways
- We can use various heuristics to decide which features to keep
- Feature selection however can be powerful, often we can get very good results by keeping only a few variables
- As well as possibly improving generalisation we also get a more **interpretable** rule

# Feature Selection

- Spurious features will allow us to find spurious rules (**over-fitting**)
- We can try different combinations of features to find the best set, although it rapidly becomes intractable to do this in all ways
- We can use various heuristics to decide which features to keep
- Feature selection however can be powerful, often we can get very good results by keeping only a few variables
- As well as possibly improving generalisation we also get a more **interpretable** rule

# Feature Selection

- Spurious features will allow us to find spurious rules (**over-fitting**)
- We can try different combinations of features to find the best set, although it rapidly becomes intractable to do this in all ways
- We can use various heuristics to decide which features to keep, **but no heuristic is fail-safe**
- Feature selection however can be powerful, often we can get very good results by keeping only a few variables
- As well as possibly improving generalisation we also get a more **interpretable** rule

# Feature Selection

- Spurious features will allow us to find spurious rules (**over-fitting**)
- We can try different combinations of features to find the best set, although it rapidly becomes intractable to do this in all ways
- We can use various heuristics to decide which features to keep, but no heuristic is fail-safe
- Feature selection however can be powerful, often we can get very good results by keeping only a few variables
- As well as possibly improving generalisation we also get a more **interpretable** rule

# Feature Selection

- Spurious features will allow us to find spurious rules (**over-fitting**)
- We can try different combinations of features to find the best set, although it rapidly becomes intractable to do this in all ways
- We can use various heuristics to decide which features to keep, but no heuristic is fail-safe
- Feature selection however can be powerful, often we can get very good results by keeping only a few variables
- As well as possibly improving generalisation we also get a more **interpretable** rule

# Explicit Regularisation

- As Niranjan showed us we can modify our error function to choose smoother functions

$$E = \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - y_n)^2$$

- Second term is minimised when  $w_i = 0$
- If  $w_i$  is large then

$$f(\mathbf{x}|\mathbf{w}) = \mathbf{w}^\top \mathbf{x}_n$$

varies rapidly as we change  $x_i$



# Explicit Regularisation

- As Niranjan showed us we can modify our error function to choose smoother functions

$$E = \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - y_n)^2 + \nu \|\mathbf{w}\|^2$$

- Second term is minimised when  $w_i = 0$
- If  $w_i$  is large then

$$f(\mathbf{x}|\mathbf{w}) = \mathbf{w}^\top \mathbf{x}_n$$

varies rapidly as we change  $x_i$

# Explicit Regularisation

- As Niranjan showed us we can modify our error function to choose smoother functions

$$E = \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - y_n)^2 + \nu \|\mathbf{w}\|^2$$

(Good to normalise data)

- Second term is minimised when  $w_i = 0$
- If  $w_i$  is large then

$$f(\mathbf{x}|\mathbf{w}) = \mathbf{w}^\top \mathbf{x}_n$$

varies rapidly as we change  $x_i$

# Explicit Regularisation

- As Niranjan showed us we can modify our error function to choose smoother functions

$$E = \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - y_n)^2 + \nu \|\mathbf{w}\|^2$$

(Good to normalise data)

- Second term is minimised when  $w_i = 0$
- If  $w_i$  is large then

$$f(\mathbf{x}|\mathbf{w}) = \mathbf{w}^\top \mathbf{x}_n$$

varies rapidly as we change  $x_i$

# Explicit Regularisation

- As Niranjan showed us we can modify our error function to choose smoother functions

$$E = \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - y_n)^2 + \nu \|\mathbf{w}\|^2$$

(Good to normalise data)

- Second term is minimised when  $w_i = 0$
- If  $w_i$  is large then

$$f(\mathbf{x}|\mathbf{w}) = \mathbf{w}^\top \mathbf{x}_n$$

varies rapidly as we change  $x_i$

# Explicit Regularisation

- As Niranjan showed us we can modify our error function to choose smoother functions

$$E = \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - y_n)^2 + \nu \|\mathbf{w}\|^2$$

(Good to normalise data)

- Second term is minimised when  $w_i = 0$
- If  $w_i$  is large then

$$f(\mathbf{x}|\mathbf{w}) = \mathbf{w}^\top \mathbf{x}_n = \sum_{i=1}^p w_i x_i$$

varies rapidly as we change  $x_i$

# Explicit Regularisation

- As Niranjan showed us we can modify our error function to choose smoother functions

$$E = \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - y_n)^2 + \nu \|\mathbf{w}\|^2$$

(Good to normalise data)

- Second term is minimised when  $w_i = 0$
- If  $w_i$  is large then

$$f(\mathbf{x}|\mathbf{w}) = \mathbf{w}^\top \mathbf{x}_n = \sum_{i=1}^p w_i x_i$$

varies rapidly as we change  $x_i$

# Lasso

- We can use other regularisers

$$E = \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - y_n)^2$$

- Spurious features (e.g. shoe size) will give us a small improvement in training error

# Lasso

- We can use other regularisers

$$E = \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - y_n)^2 + \nu \sum_{i=1}^p |w_i|$$

- Spurious features (e.g. shoe size) will give us a small improvement in training error



# Lasso

- We can use other regularisers

$$E = \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - y_n)^2 + \nu \sum_{i=1}^p |w_i|$$

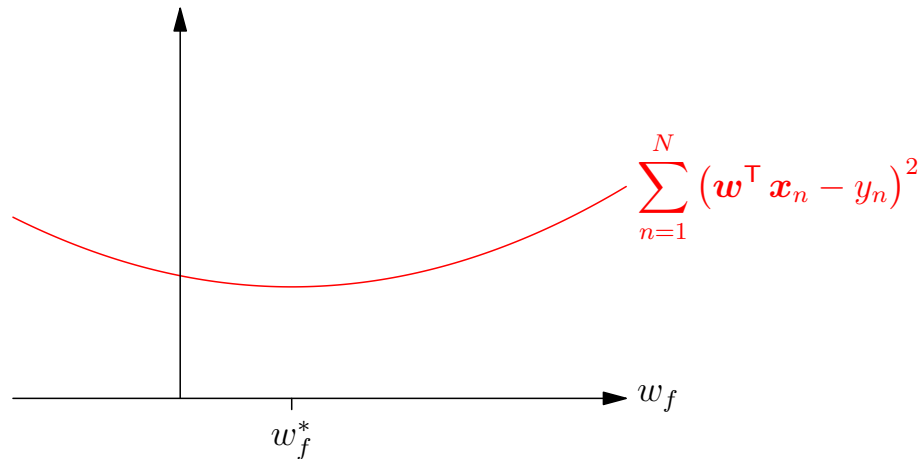
- Spurious features (e.g. shoe size) will give us a small improvement in training error

# Lasso

- We can use other regularisers

$$E = \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - y_n)^2 + \nu \sum_{i=1}^p |w_i|$$

- Spurious features (e.g. shoe size) will give us a small improvement in training error

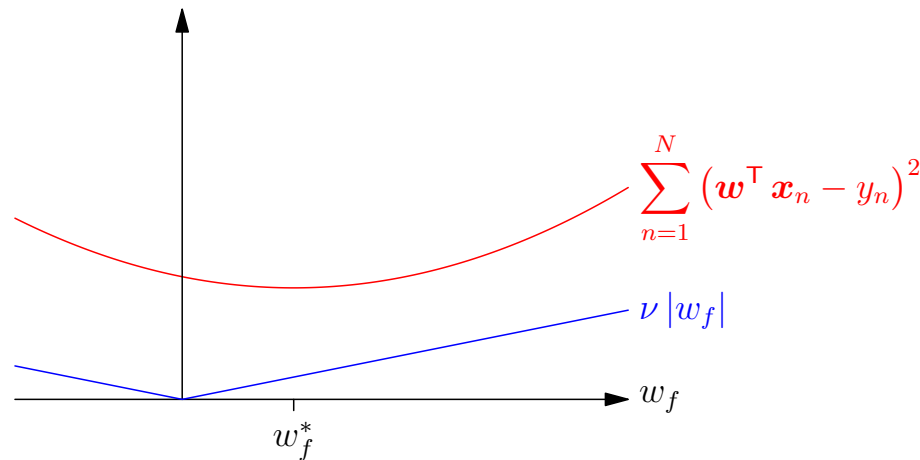


# Lasso

- We can use other regularisers

$$E = \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - y_n)^2 + \nu \sum_{i=1}^p |w_i|$$

- Spurious features (e.g. shoe size) will give us a small improvement in training error

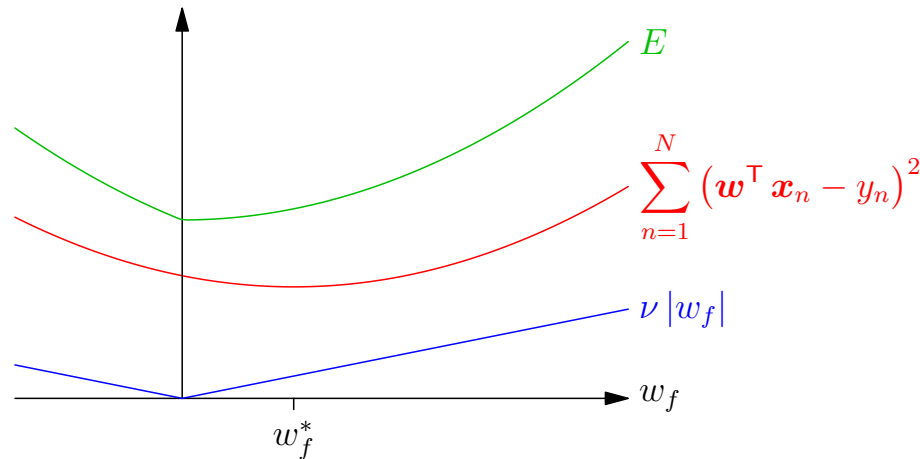


# Lasso

- We can use other regularisers

$$E = \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - y_n)^2 + \nu \sum_{i=1}^p |w_i|$$

- Spurious features (e.g. shoe size) will give us a small improvement in training error

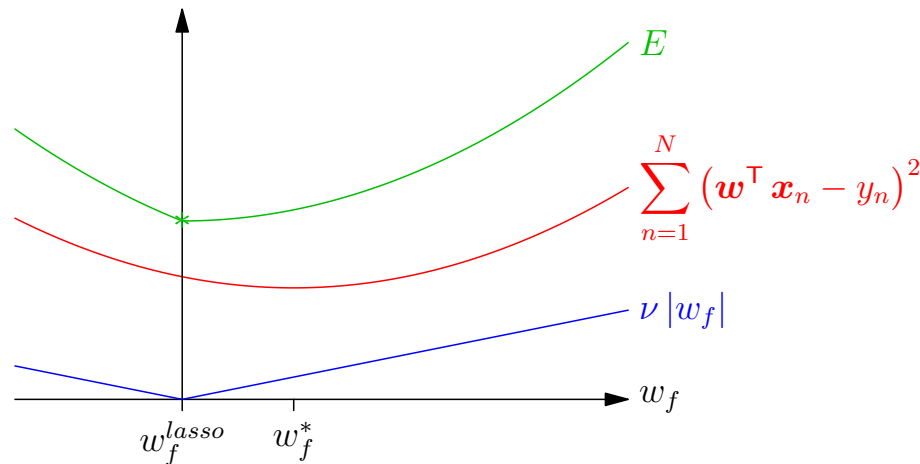


# Lasso

- We can use other regularisers

$$E = \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - y_n)^2 + \nu \sum_{i=1}^p |w_i|$$

- Spurious features (e.g. shoe size) will give us a small improvement in training error

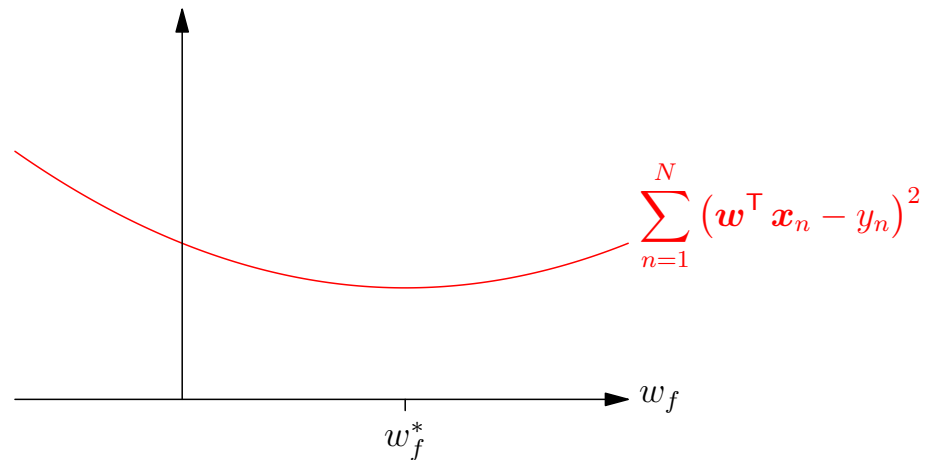


# Lasso

- We can use other regularisers

$$E = \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - y_n)^2 + \nu \sum_{i=1}^p |w_i|$$

- Spurious features (e.g. shoe size) will give us a small improvement in training error

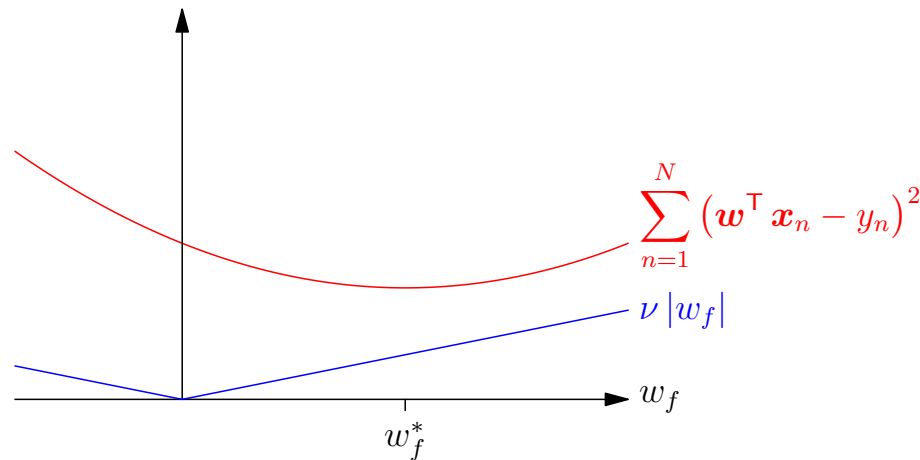


# Lasso

- We can use other regularisers

$$E = \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - y_n)^2 + \nu \sum_{i=1}^p |w_i|$$

- Spurious features (e.g. shoe size) will give us a small improvement in training error

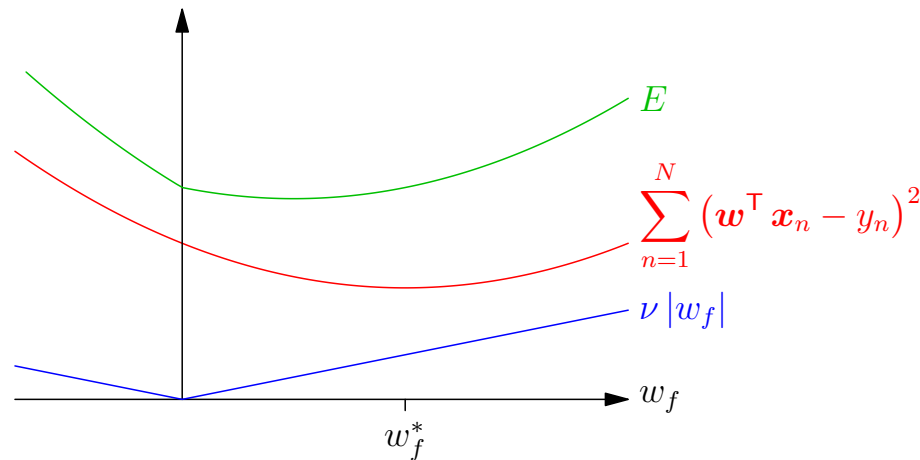


# Lasso

- We can use other regularisers

$$E = \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - y_n)^2 + \nu \sum_{i=1}^p |w_i|$$

- Spurious features (e.g. shoe size) will give us a small improvement in training error



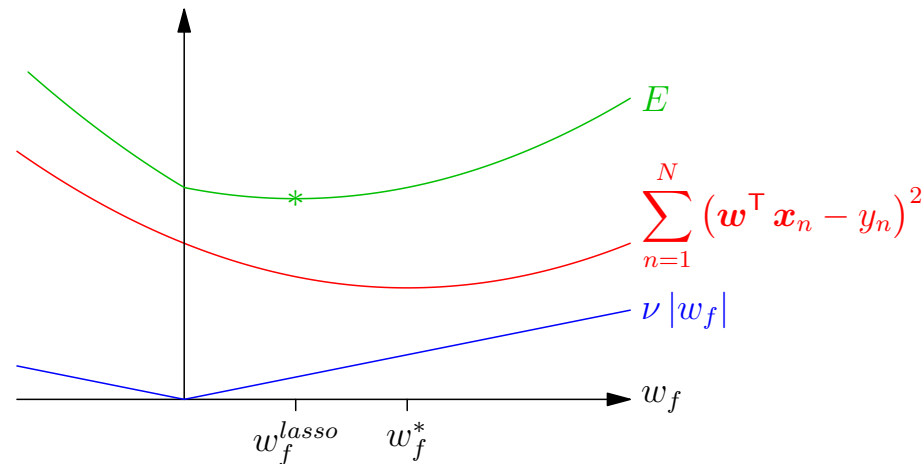


# Lasso

- We can use other regularisers

$$E = \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - y_n)^2 + \nu \sum_{i=1}^p |w_i|$$

- Spurious features (e.g. shoe size) will give us a small improvement in training error

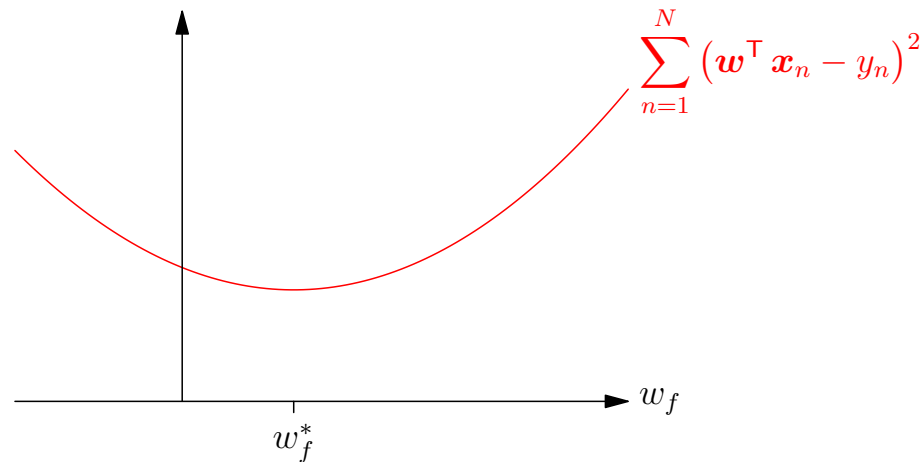


# Lasso

- We can use other regularisers

$$E = \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - y_n)^2 + \nu \sum_{i=1}^p |w_i|$$

- Spurious features (e.g. shoe size) will give us a small improvement in training error

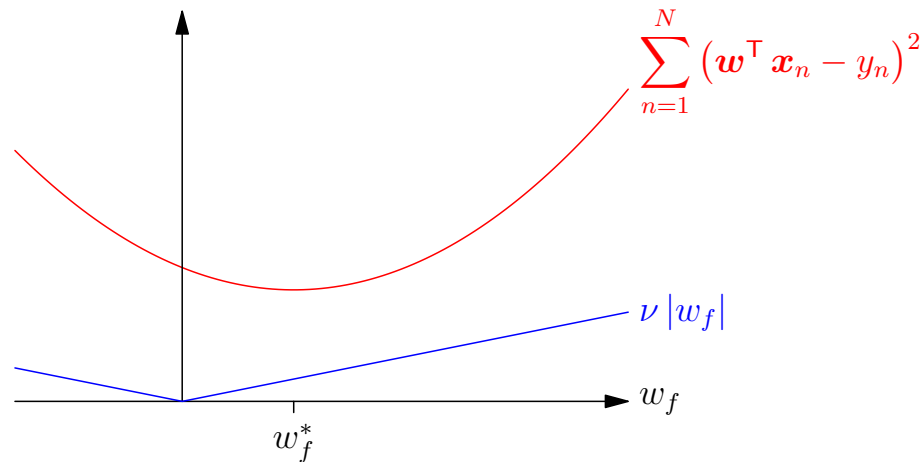


# Lasso

- We can use other regularisers

$$E = \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - y_n)^2 + \nu \sum_{i=1}^p |w_i|$$

- Spurious features (e.g. shoe size) will give us a small improvement in training error

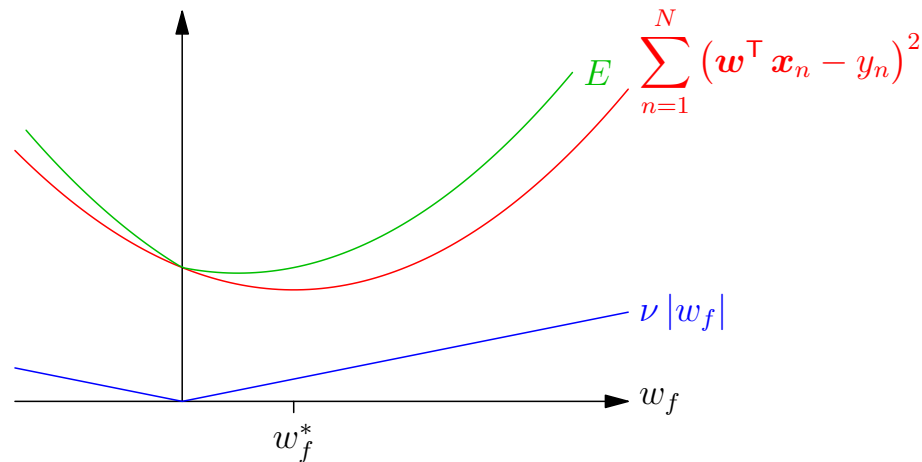


# Lasso

- We can use other regularisers

$$E = \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - y_n)^2 + \nu \sum_{i=1}^p |w_i|$$

- Spurious features (e.g. shoe size) will give us a small improvement in training error

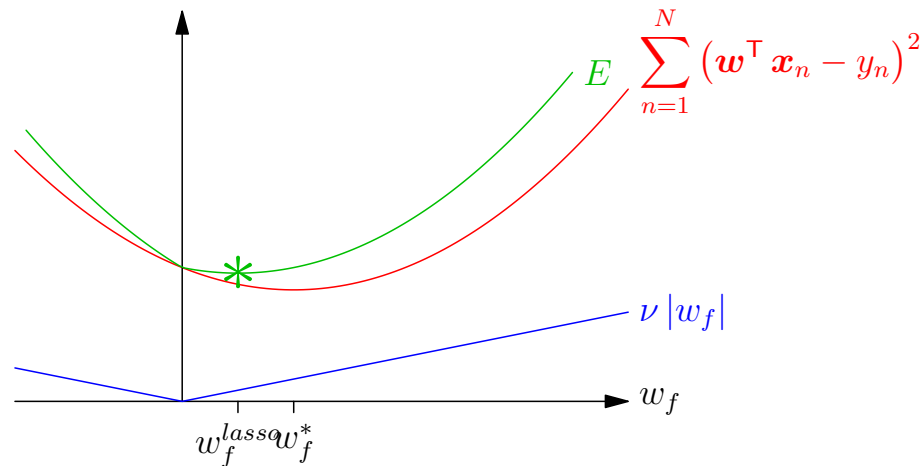


# Lasso

- We can use other regularisers

$$E = \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - y_n)^2 + \nu \sum_{i=1}^p |w_i|$$

- Spurious features (e.g. shoe size) will give us a small improvement in training error

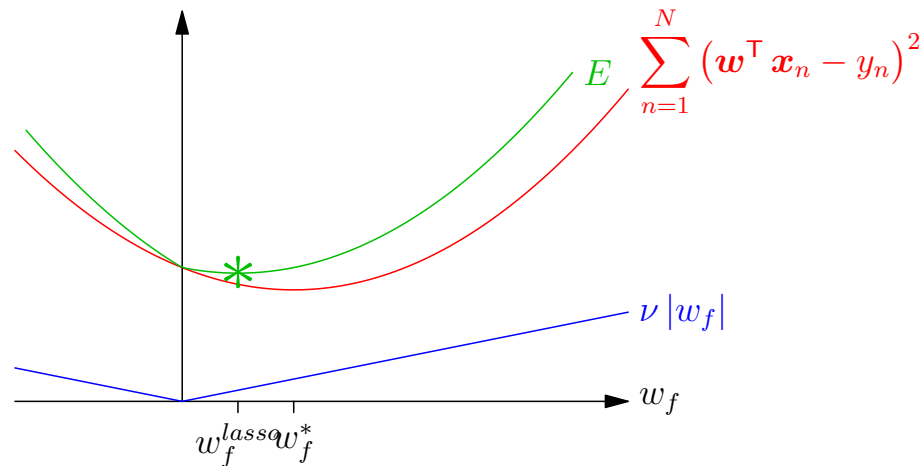


# Lasso

- We can use other regularisers

$$E = \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - y_n)^2 + \nu \sum_{i=1}^p |w_i|$$

- Spurious features (e.g. shoe size) will give us a small improvement in training error



- Does automatic feature selection

# Implicit Regularisation

- In the last two examples we added an explicit regularisation term that made the function we learnt simpler
- Some learning machines do this less explicitly
- Some deep learning architectures do subtle averaging
- Sometimes the architecture biases the machine to find a simple solution
- We will see this in support vector machines shortly

# Implicit Regularisation

- In the last two examples we added an explicit regularisation term that made the function we learnt simpler
- Some learning machines do this less explicitly
- Some deep learning architectures do subtle averaging
- Sometimes the architecture biases the machine to find a simple solution
- We will see this in support vector machines shortly



# Implicit Regularisation

- In the last two examples we added an explicit regularisation term that made the function we learnt simpler
- Some learning machines do this less explicitly
- Some deep learning architectures do subtle averaging
- Sometimes the architecture biases the machine to find a simple solution
- We will see this in support vector machines shortly

# Implicit Regularisation

- In the last two examples we added an explicit regularisation term that made the function we learnt simpler
- Some learning machines do this less explicitly
- Some deep learning architectures do subtle averaging
- Sometimes the architecture biases the machine to find a simple solution
- We will see this in support vector machines shortly

# Implicit Regularisation

- In the last two examples we added an explicit regularisation term that made the function we learnt simpler
- Some learning machines do this less explicitly
- Some deep learning architectures do subtle averaging
- Sometimes the architecture biases the machine to find a simple solution
- We will see this in support vector machines shortly

# Choosing Machines

- We often have to choose between learning machines
- Or select “*hyper-parameters*”
- Decide what features to use
- We don't want to choose a model that over-fits the training data
- Need an unbiased estimate of the generalisation performance

# Choosing Machines

- We often have to choose between learning machines
- Or select “*hyper-parameters*”
- Decide what features to use
- We don't want to choose a model that over-fits the training data
- Need an unbiased estimate of the generalisation performance

# Choosing Machines

- We often have to choose between learning machines
- Or select “*hyper-parameters*”
- Decide what features to use
- We don't want to choose a model that over-fits the training data
- Need an unbiased estimate of the generalisation performance

# Choosing Machines

- We often have to choose between learning machines
- Or select “*hyper-parameters*”
- Decide what features to use
- We don't want to choose a model that over-fits the training data
- Need an unbiased estimate of the generalisation performance

# Choosing Machines

- We often have to choose between learning machines
- Or select “*hyper-parameters*”
- Decide what features to use
- We don't want to choose a model that over-fits the training data
- Need an unbiased estimate of the generalisation performance



# Measuring Generalisation Performance

- Recall, we want to predict **unseen** data
- **You cannot use data that you have trained on!**
- Need to split your data up into training and validation set
- Use the validation set to choose the hyper-parameters
- Because your validation set is small, you only get an estimate of your generalisation error

# Measuring Generalisation Performance

- Recall, we want to predict **unseen** data
- **You cannot use data that you have trained on!**
- Need to split your data up into training and validation set
- Use the validation set to choose the hyper-parameters
- Because your validation set is small, you only get an estimate of your generalisation error

# Measuring Generalisation Performance

- Recall, we want to predict **unseen** data
- **You cannot use data that you have trained on!—you will overfit**
- Need to split your data up into training and validation set
- Use the validation set to choose the hyper-parameters
- Because your validation set is small, you only get an estimate of your generalisation error

# Measuring Generalisation Performance

- Recall, we want to predict **unseen** data
- **You cannot use data that you have trained on!**—you will overfit
- Need to split your data up into training and validation set
- Use the validation set to choose the hyper-parameters
- Because your validation set is small, you only get an estimate of your generalisation error

# Measuring Generalisation Performance

- Recall, we want to predict **unseen** data
- **You cannot use data that you have trained on!**—you will overfit
- Need to split your data up into training and validation set
- Use the validation set to choose the hyper-parameters
- Because your validation set is small, you only get an estimate of your generalisation error

# Measuring Generalisation Performance

- Recall, we want to predict **unseen** data
- **You cannot use data that you have trained on!**—you will overfit
- Need to split your data up into training and validation set
- Use the validation set to choose the hyper-parameters
- Because your validation set is small, you only get an estimate of your generalisation error

# Cross Validation

- If you want to use more data for training then you can use cross validation
- $K$ -fold cross validation splits the data into  $K$  groups

# Cross Validation

- If you want to use more data for training then you can use cross validation
- $K$ -fold cross validation splits the data into  $K$  groups



# Cross Validation

- If you want to use more data for training then you can use cross validation
- $K$ -fold cross validation splits the data into  $K$  groups

$$\mathcal{D} = \{D_i\}_{i=1}^P \quad D_i = (\mathbf{x}_i, y_i)$$

$D_1$	$D_2$	$D_3$	$D_4$	$D_5$	$D_6$	$D_7$	$D_8$	$D_9$	$D_{10}$	$D_{11}$	$D_{12}$	$D_{13}$	$D_{14}$	$D_{15}$	$D_{16}$	$D_{17}$	$D_{18}$	$D_{19}$	$D_{20}$
-------	-------	-------	-------	-------	-------	-------	-------	-------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------

# Cross Validation

- If you want to use more data for training then you can use cross validation
- $K$ -fold cross validation splits the data into  $K$  groups

$$\mathcal{D} = \{D_i\}_{i=1}^P \quad D_i = (\mathbf{x}_i, y_i)$$

$D_1$	$D_2$	$D_3$	$D_4$	$D_5$	$D_6$	$D_7$	$D_8$	$D_9$	$D_{10}$	$D_{11}$	$D_{12}$	$D_{13}$	$D_{14}$	$D_{15}$	$D_{16}$	$D_{17}$	$D_{18}$	$D_{19}$	$D_{20}$
-------	-------	-------	-------	-------	-------	-------	-------	-------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------

Test Set

Training Set

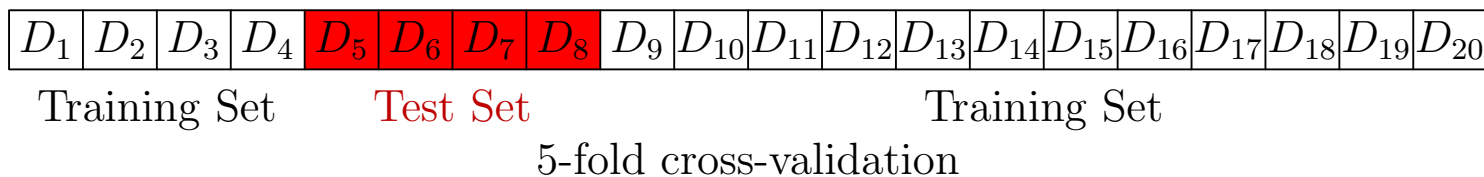
5-fold cross-validation

$$E_g = 5.1$$

# Cross Validation

- If you want to use more data for training then you can use cross validation
- $K$ -fold cross validation splits the data into  $K$  groups

$$\mathcal{D} = \{D_i\}_{i=1}^P \quad D_i = (\mathbf{x}_i, y_i)$$

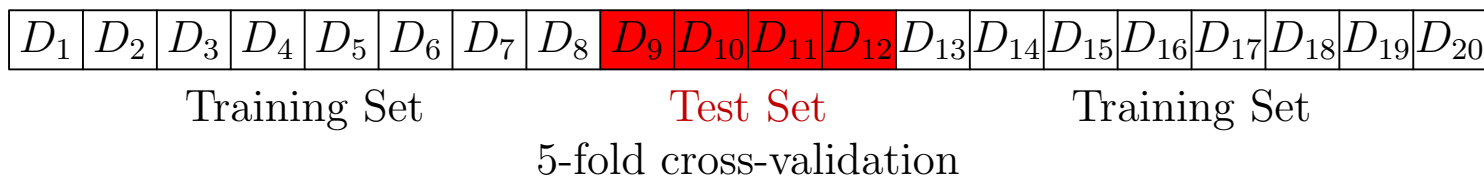


$$E_g = 3.7$$

# Cross Validation

- If you want to use more data for training then you can use cross validation
- $K$ -fold cross validation splits the data into  $K$  groups

$$\mathcal{D} = \{D_i\}_{i=1}^P \quad D_i = (\mathbf{x}_i, y_i)$$

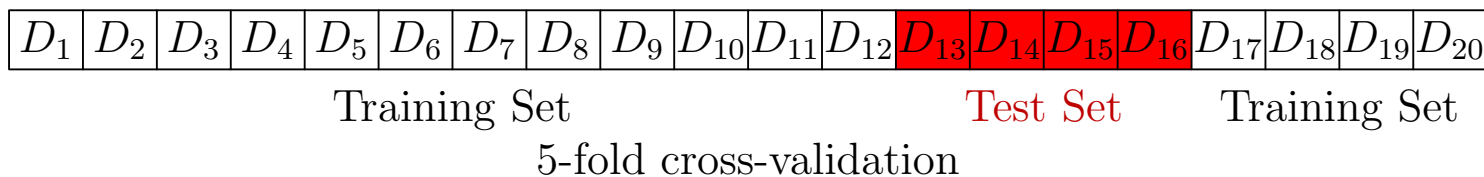


$$E_g = 4.6$$

# Cross Validation

- If you want to use more data for training then you can use cross validation
- $K$ -fold cross validation splits the data into  $K$  groups

$$\mathcal{D} = \{D_i\}_{i=1}^P \quad D_i = (\mathbf{x}_i, y_i)$$



$E_g =$

4.6

# Cross Validation

- If you want to use more data for training then you can use cross validation
- $K$ -fold cross validation splits the data into  $K$  groups

$$\mathcal{D} = \{D_i\}_{i=1}^P \quad D_i = (\mathbf{x}_i, y_i)$$

$D_1$	$D_2$	$D_3$	$D_4$	$D_5$	$D_6$	$D_7$	$D_8$	$D_9$	$D_{10}$	$D_{11}$	$D_{12}$	$D_{13}$	$D_{14}$	$D_{15}$	$D_{16}$	$D_{17}$	$D_{18}$	$D_{19}$	$D_{20}$
-------	-------	-------	-------	-------	-------	-------	-------	-------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------

Training Set

Test Set

5-fold cross-validation

$E_g =$

3.3

# Cross Validation

- If you want to use more data for training then you can use cross validation
- $K$ -fold cross validation splits the data into  $K$  groups

$$\mathcal{D} = \{D_i\}_{i=1}^P \quad D_i = (\mathbf{x}_i, y_i)$$

$D_1$	$D_2$	$D_3$	$D_4$	$D_5$	$D_6$	$D_7$	$D_8$	$D_9$	$D_{10}$	$D_{11}$	$D_{12}$	$D_{13}$	$D_{14}$	$D_{15}$	$D_{16}$	$D_{17}$	$D_{18}$	$D_{19}$	$D_{20}$
-------	-------	-------	-------	-------	-------	-------	-------	-------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------

$$\langle E_g \rangle = \frac{5.1 + 3.7 + 4.6 + 4.6 + 3.3}{5} = 4.3$$

# Cross Validation

- If you want to use more data for training then you can use cross validation
- $K$ -fold cross validation splits the data into  $K$  groups

$$\mathcal{D} = \{D_i\}_{i=1}^P \quad D_i = (\mathbf{x}_i, y_i)$$

$D_1$	$D_2$	$D_3$	$D_4$	$D_5$	$D_6$	$D_7$	$D_8$	$D_9$	$D_{10}$	$D_{11}$	$D_{12}$	$D_{13}$	$D_{14}$	$D_{15}$	$D_{16}$	$D_{17}$	$D_{18}$	$D_{19}$	$D_{20}$
-------	-------	-------	-------	-------	-------	-------	-------	-------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------

Test Set

Training Set

10-fold cross-validation

$$E_g = 5.4$$



# Cross Validation

- If you want to use more data for training then you can use cross validation
- $K$ -fold cross validation splits the data into  $K$  groups

$$\mathcal{D} = \{D_i\}_{i=1}^P \quad D_i = (\mathbf{x}_i, y_i)$$

$D_1$	$D_2$	$D_3$	$D_4$	$D_5$	$D_6$	$D_7$	$D_8$	$D_9$	$D_{10}$	$D_{11}$	$D_{12}$	$D_{13}$	$D_{14}$	$D_{15}$	$D_{16}$	$D_{17}$	$D_{18}$	$D_{19}$	$D_{20}$
-------	-------	-------	-------	-------	-------	-------	-------	-------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------

Test Set

Training Set

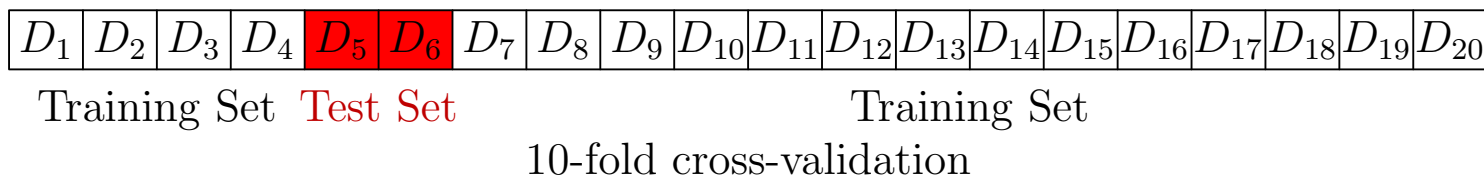
10-fold cross-validation

$$E_g = 1.4$$

# Cross Validation

- If you want to use more data for training then you can use cross validation
- $K$ -fold cross validation splits the data into  $K$  groups

$$\mathcal{D} = \{D_i\}_{i=1}^P \quad D_i = (\mathbf{x}_i, y_i)$$

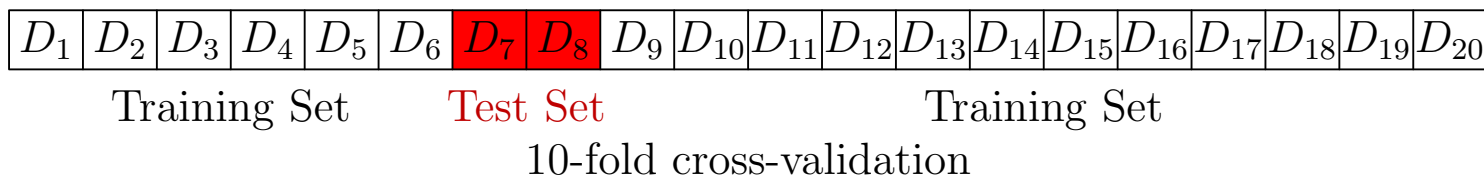


$$E_g = 4.4$$

# Cross Validation

- If you want to use more data for training then you can use cross validation
- $K$ -fold cross validation splits the data into  $K$  groups

$$\mathcal{D} = \{D_i\}_{i=1}^P \quad D_i = (\mathbf{x}_i, y_i)$$



$$E_g = 3.2$$

# Cross Validation

- If you want to use more data for training then you can use cross validation
- $K$ -fold cross validation splits the data into  $K$  groups

$$\mathcal{D} = \{D_i\}_{i=1}^P \quad D_i = (\mathbf{x}_i, y_i)$$

$D_1$	$D_2$	$D_3$	$D_4$	$D_5$	$D_6$	$D_7$	$D_8$	$D_9$	$D_{10}$	$D_{11}$	$D_{12}$	$D_{13}$	$D_{14}$	$D_{15}$	$D_{16}$	$D_{17}$	$D_{18}$	$D_{19}$	$D_{20}$
-------	-------	-------	-------	-------	-------	-------	-------	-------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------

Training Set

Test Set

Training Set

10-fold cross-validation

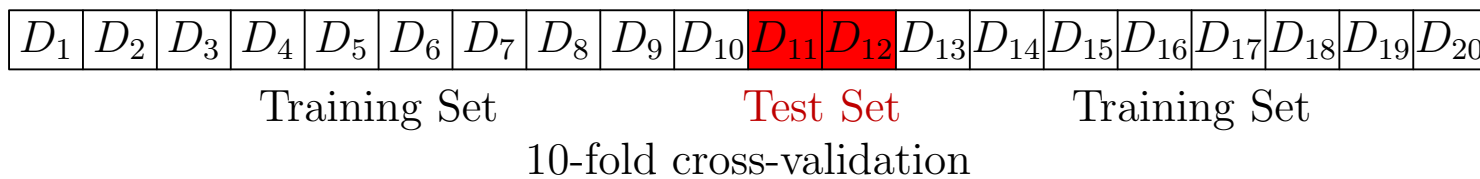
$E_g =$

7

# Cross Validation

- If you want to use more data for training then you can use cross validation
- $K$ -fold cross validation splits the data into  $K$  groups

$$\mathcal{D} = \{D_i\}_{i=1}^P \quad D_i = (\mathbf{x}_i, y_i)$$

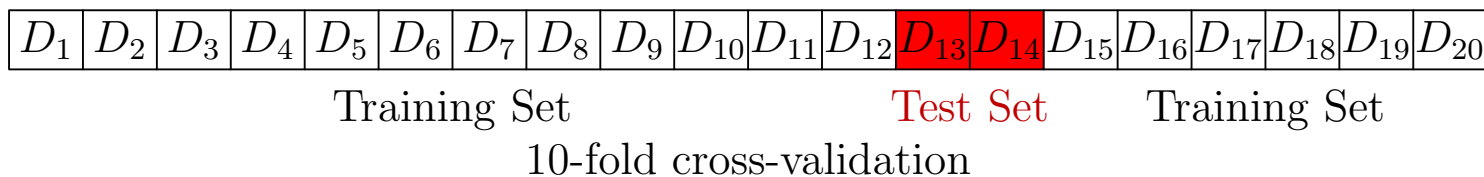


$$E_g = 0.59$$

# Cross Validation

- If you want to use more data for training then you can use cross validation
- $K$ -fold cross validation splits the data into  $K$  groups

$$\mathcal{D} = \{D_i\}_{i=1}^P \quad D_i = (\mathbf{x}_i, y_i)$$



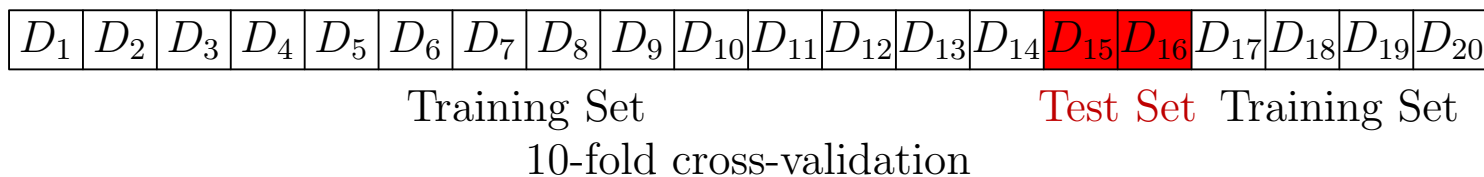
$E_g =$

4.1

# Cross Validation

- If you want to use more data for training then you can use cross validation
- $K$ -fold cross validation splits the data into  $K$  groups

$$\mathcal{D} = \{D_i\}_{i=1}^P \quad D_i = (\mathbf{x}_i, y_i)$$



$E_g =$

5

# Cross Validation

- If you want to use more data for training then you can use cross validation
- $K$ -fold cross validation splits the data into  $K$  groups

$$\mathcal{D} = \{D_i\}_{i=1}^P \quad D_i = (\mathbf{x}_i, y_i)$$

$D_1$	$D_2$	$D_3$	$D_4$	$D_5$	$D_6$	$D_7$	$D_8$	$D_9$	$D_{10}$	$D_{11}$	$D_{12}$	$D_{13}$	$D_{14}$	$D_{15}$	$D_{16}$	$D_{17}$	$D_{18}$	$D_{19}$	$D_{20}$
-------	-------	-------	-------	-------	-------	-------	-------	-------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------

Training Set

Test Set

10-fold cross-validation

$E_g =$

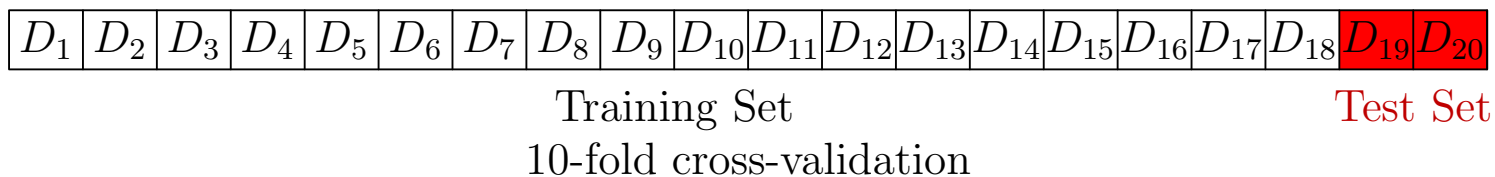
5.8



# Cross Validation

- If you want to use more data for training then you can use cross validation
- $K$ -fold cross validation splits the data into  $K$  groups

$$\mathcal{D} = \{D_i\}_{i=1}^P \quad D_i = (\mathbf{x}_i, y_i)$$



$E_g =$

2.3

# Cross Validation

- If you want to use more data for training then you can use cross validation
- $K$ -fold cross validation splits the data into  $K$  groups

$$\mathcal{D} = \{D_i\}_{i=1}^P \quad D_i = (\mathbf{x}_i, y_i)$$

$D_1$	$D_2$	$D_3$	$D_4$	$D_5$	$D_6$	$D_7$	$D_8$	$D_9$	$D_{10}$	$D_{11}$	$D_{12}$	$D_{13}$	$D_{14}$	$D_{15}$	$D_{16}$	$D_{17}$	$D_{18}$	$D_{19}$	$D_{20}$
-------	-------	-------	-------	-------	-------	-------	-------	-------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------

$$\langle E_g \rangle = \frac{5.4 + 1.4 + 4.4 + 3.2 + 7 + 0.59 + 4.1 + 5 + 5.8 + 2.3}{10} = 3.9$$

# Cross Validation

- If you want to use more data for training then you can use cross validation
- $K$ -fold cross validation splits the data into  $K$  groups

$$\mathcal{D} = \{D_i\}_{i=1}^P \quad D_i = (\mathbf{x}_i, y_i)$$

$D_1$	$D_2$	$D_3$	$D_4$	$D_5$	$D_6$	$D_7$	$D_8$	$D_9$	$D_{10}$	$D_{11}$	$D_{12}$	$D_{13}$	$D_{14}$	$D_{15}$	$D_{16}$	$D_{17}$	$D_{18}$	$D_{19}$	$D_{20}$
-------	-------	-------	-------	-------	-------	-------	-------	-------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------

Test

Leave-one-out cross-validation

$$E_g = 4.2$$

# Cross Validation

- If you want to use more data for training then you can use cross validation
- $K$ -fold cross validation splits the data into  $K$  groups

$$\mathcal{D} = \{D_i\}_{i=1}^P \quad D_i = (\mathbf{x}_i, y_i)$$

$D_1$	$D_2$	$D_3$	$D_4$	$D_5$	$D_6$	$D_7$	$D_8$	$D_9$	$D_{10}$	$D_{11}$	$D_{12}$	$D_{13}$	$D_{14}$	$D_{15}$	$D_{16}$	$D_{17}$	$D_{18}$	$D_{19}$	$D_{20}$
-------	-------	-------	-------	-------	-------	-------	-------	-------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------

Test

Leave-one-out cross-validation

$$E_g = 2.9$$

# Cross Validation

- If you want to use more data for training then you can use cross validation
- $K$ -fold cross validation splits the data into  $K$  groups

$$\mathcal{D} = \{D_i\}_{i=1}^P \quad D_i = (\mathbf{x}_i, y_i)$$

$D_1$	$D_2$	$D_3$	$D_4$	$D_5$	$D_6$	$D_7$	$D_8$	$D_9$	$D_{10}$	$D_{11}$	$D_{12}$	$D_{13}$	$D_{14}$	$D_{15}$	$D_{16}$	$D_{17}$	$D_{18}$	$D_{19}$	$D_{20}$
-------	-------	-------	-------	-------	-------	-------	-------	-------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------

Test

Leave-one-out cross-validation

$$E_g = 4.2$$

# Cross Validation

- If you want to use more data for training then you can use cross validation
- $K$ -fold cross validation splits the data into  $K$  groups

$$\mathcal{D} = \{D_i\}_{i=1}^P \quad D_i = (\mathbf{x}_i, y_i)$$

$D_1$	$D_2$	$D_3$	$D_4$	$D_5$	$D_6$	$D_7$	$D_8$	$D_9$	$D_{10}$	$D_{11}$	$D_{12}$	$D_{13}$	$D_{14}$	$D_{15}$	$D_{16}$	$D_{17}$	$D_{18}$	$D_{19}$	$D_{20}$
-------	-------	-------	-------	-------	-------	-------	-------	-------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------

Test

Leave-one-out cross-validation

$$E_g = 1.4$$

# Cross Validation

- If you want to use more data for training then you can use cross validation
- $K$ -fold cross validation splits the data into  $K$  groups

$$\mathcal{D} = \{D_i\}_{i=1}^P \quad D_i = (\mathbf{x}_i, y_i)$$

$D_1$	$D_2$	$D_3$	$D_4$	$D_5$	$D_6$	$D_7$	$D_8$	$D_9$	$D_{10}$	$D_{11}$	$D_{12}$	$D_{13}$	$D_{14}$	$D_{15}$	$D_{16}$	$D_{17}$	$D_{18}$	$D_{19}$	$D_{20}$
-------	-------	-------	-------	-------	-------	-------	-------	-------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------

Test

Leave-one-out cross-validation

$$E_g = 3$$

# Cross Validation

- If you want to use more data for training then you can use cross validation
- $K$ -fold cross validation splits the data into  $K$  groups

$$\mathcal{D} = \{D_i\}_{i=1}^P \quad D_i = (\mathbf{x}_i, y_i)$$

$D_1$	$D_2$	$D_3$	$D_4$	$D_5$	$D_6$	$D_7$	$D_8$	$D_9$	$D_{10}$	$D_{11}$	$D_{12}$	$D_{13}$	$D_{14}$	$D_{15}$	$D_{16}$	$D_{17}$	$D_{18}$	$D_{19}$	$D_{20}$
-------	-------	-------	-------	-------	-------	-------	-------	-------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------

Test

Leave-one-out cross-validation

$$E_g = 3.7$$



# Cross Validation

- If you want to use more data for training then you can use cross validation
- $K$ -fold cross validation splits the data into  $K$  groups

$$\mathcal{D} = \{D_i\}_{i=1}^P \quad D_i = (\mathbf{x}_i, y_i)$$

$D_1$	$D_2$	$D_3$	$D_4$	$D_5$	$D_6$	$D_7$	$D_8$	$D_9$	$D_{10}$	$D_{11}$	$D_{12}$	$D_{13}$	$D_{14}$	$D_{15}$	$D_{16}$	$D_{17}$	$D_{18}$	$D_{19}$	$D_{20}$
-------	-------	-------	-------	-------	-------	-------	-------	-------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------

$$\langle E_g \rangle = 3.98$$

- Leave-one-out cross-validation is extreme case

# Testing and Validation

- The *validation set* is usually used for choosing hyper-parameters
- This also can include ways of preprocessing, how long you train for, different ways of handling missing data, etc.
- Having chosen your best machine you now have a **biased** estimate of your generalisation performance
- You need a third independent set, *the test set*, to obtain an unbiased estimate of the generalisation error

# Testing and Validation

- The *validation set* is usually used for choosing hyper-parameters
- This also can include ways of preprocessing, how long you train for, different ways of handling missing data, etc.
- Having chosen your best machine you now have a **biased** estimate of your generalisation performance
- You need a third independent set, *the test set*, to obtain an unbiased estimate of the generalisation error

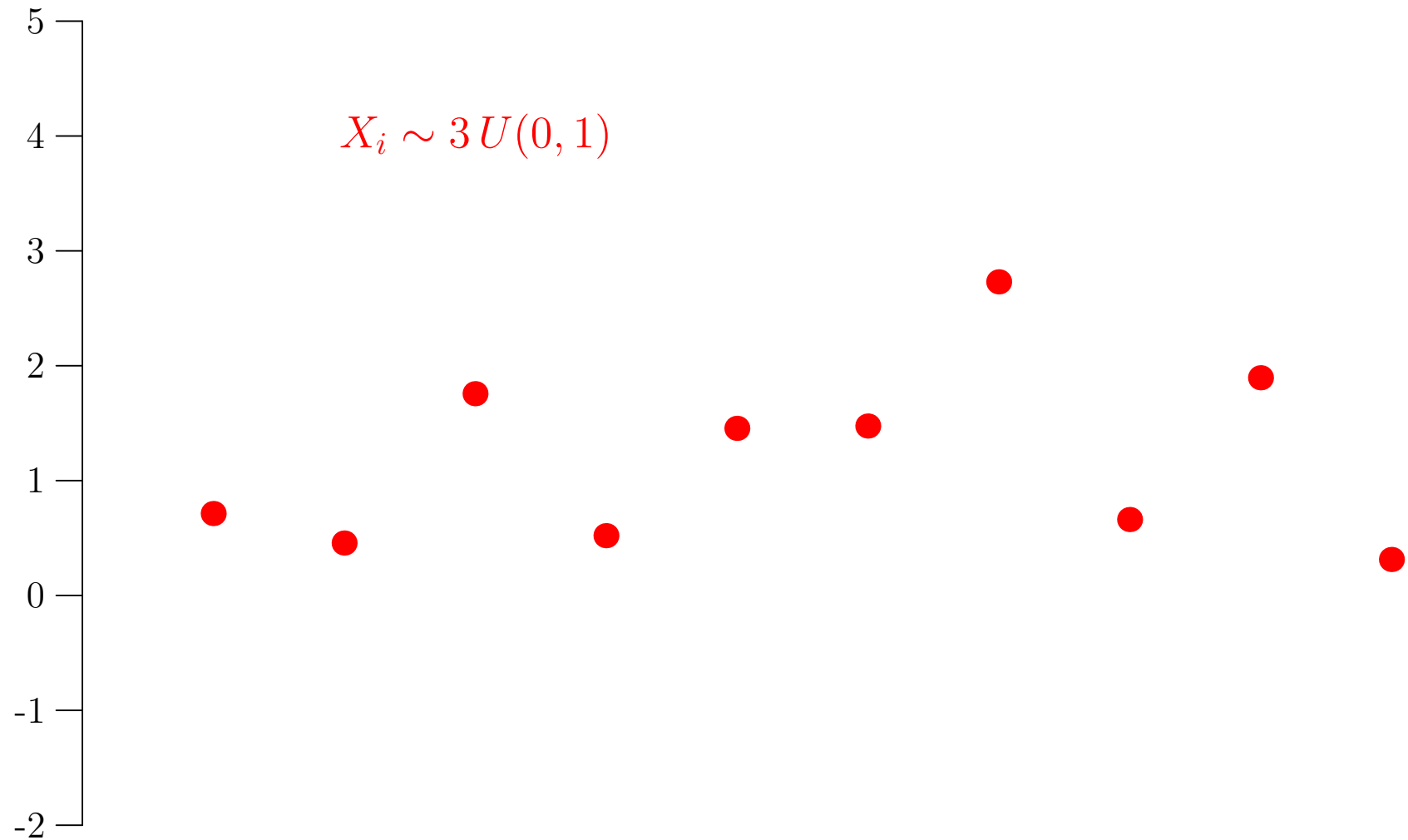
# Testing and Validation

- The *validation set* is usually used for choosing hyper-parameters
- This also can include ways of preprocessing, how long you train for, different ways of handling missing data, etc.
- Having chosen your best machine you now have a **biased** estimate of your generalisation performance
- You need a third independent set, *the test set*, to obtain an unbiased estimate of the generalisation error

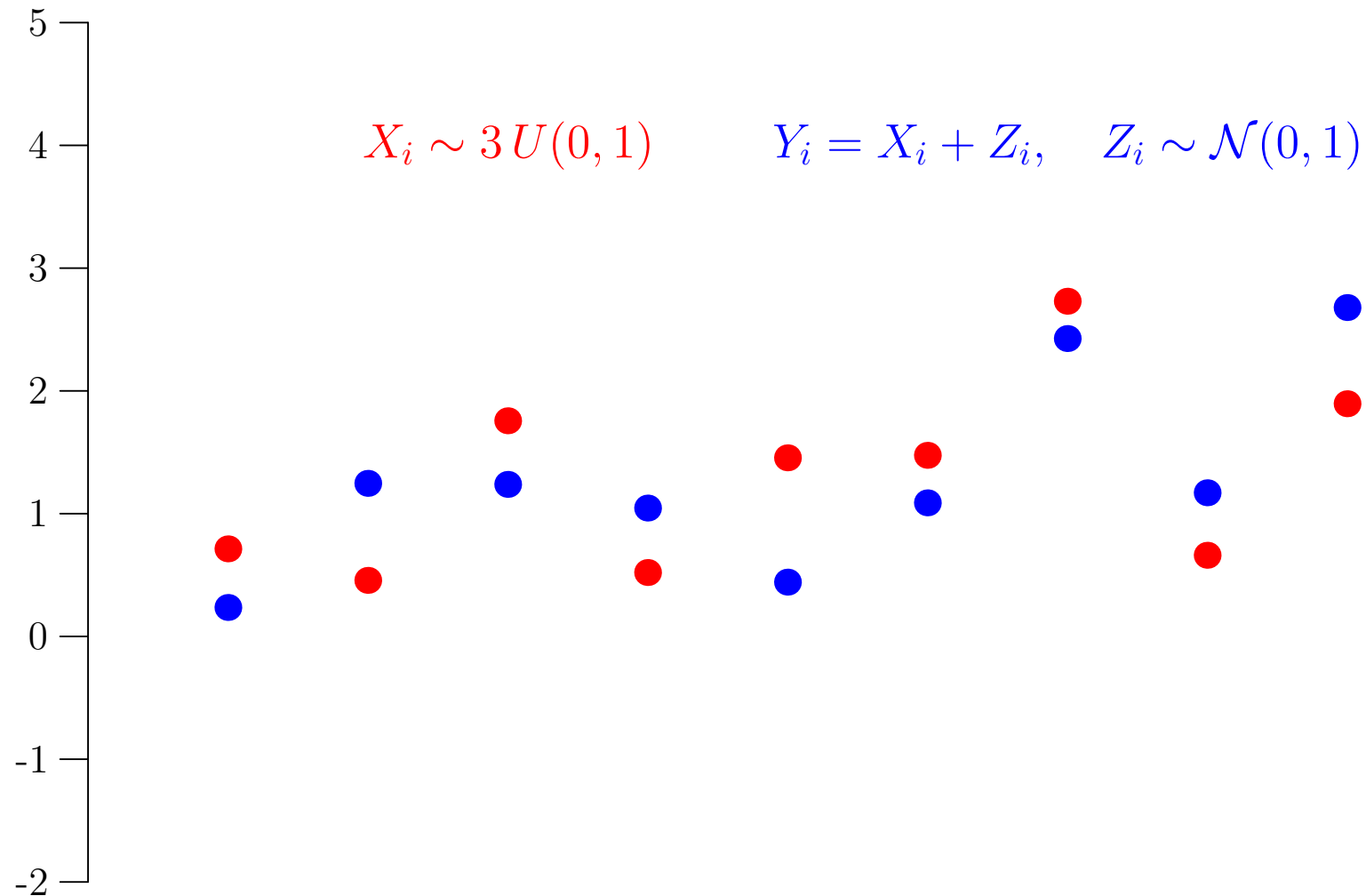
# Testing and Validation

- The *validation set* is usually used for choosing hyper-parameters
- This also can include ways of preprocessing, how long you train for, different ways of handling missing data, etc.
- Having chosen your best machine you now have a **biased** estimate of your generalisation performance
- You need a third independent set, *the test set*, to obtain an unbiased estimate of the generalisation error

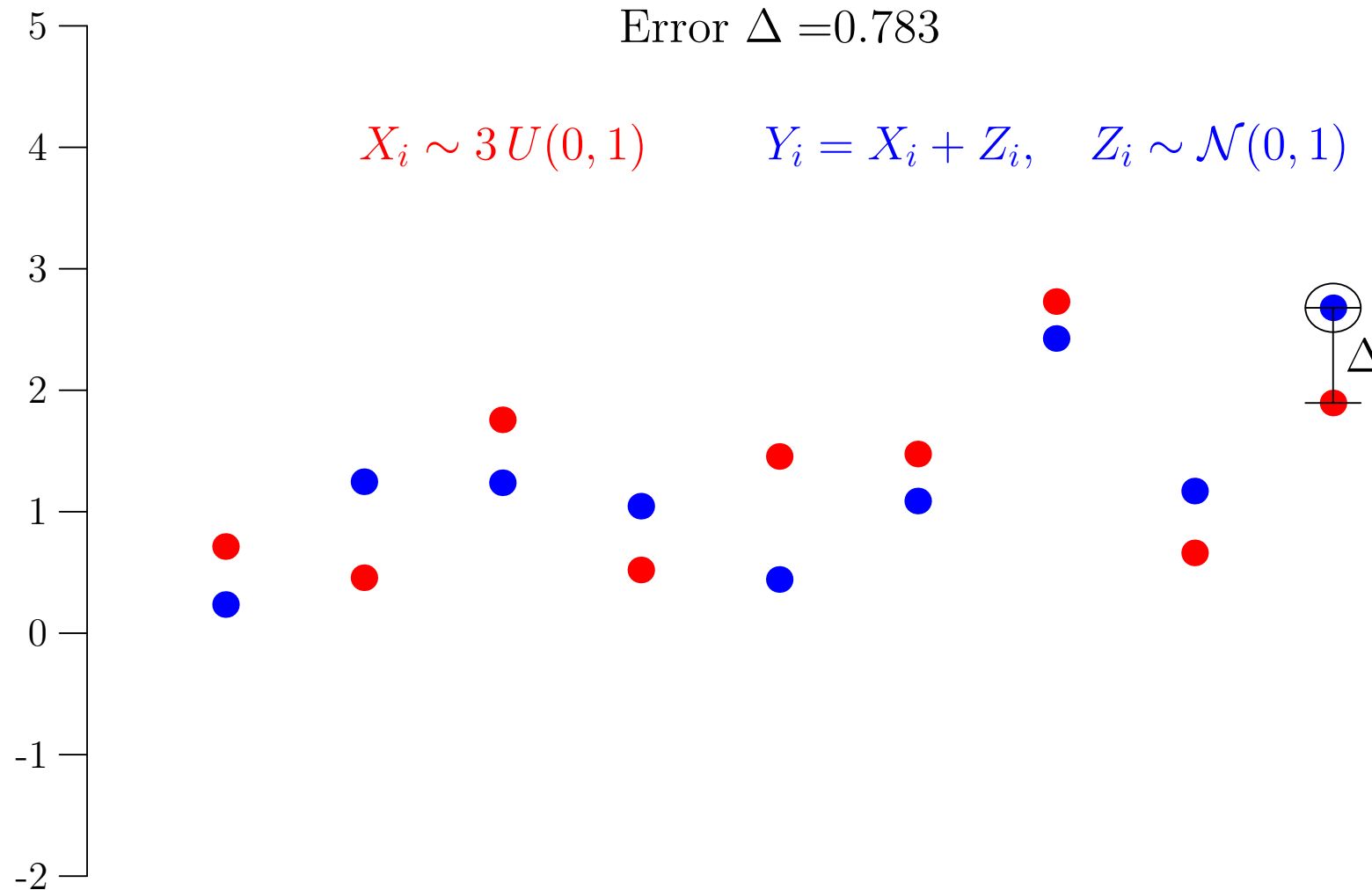
# The Over-fitting Game



# The Over-fitting Game

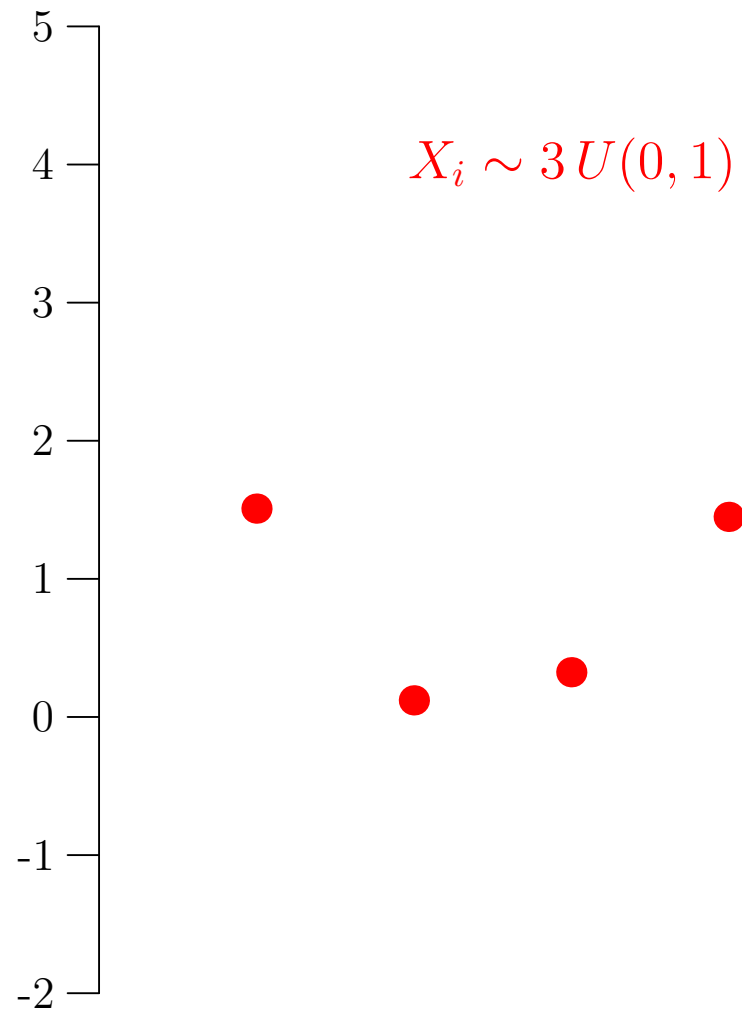


# The Over-fitting Game

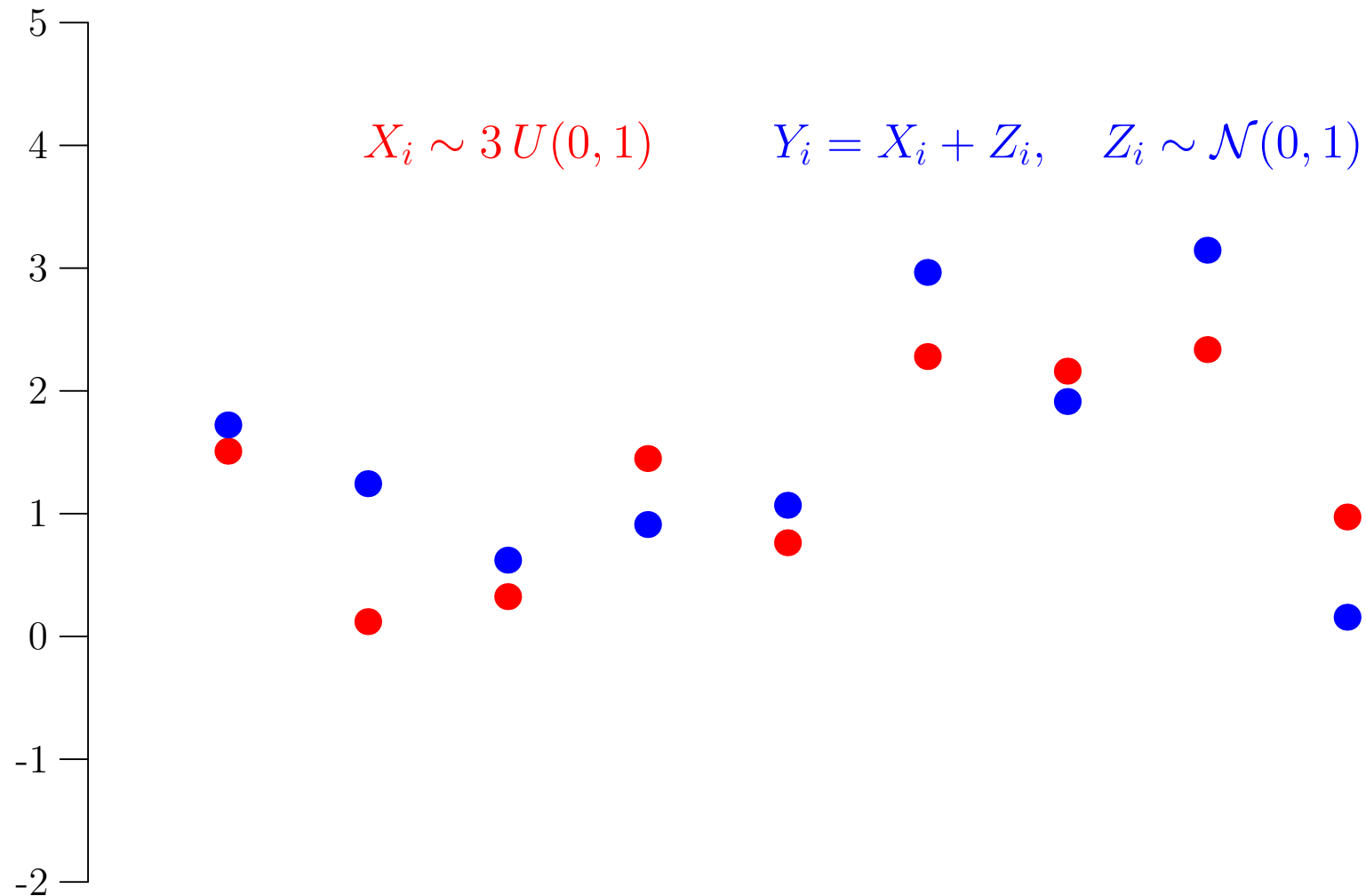




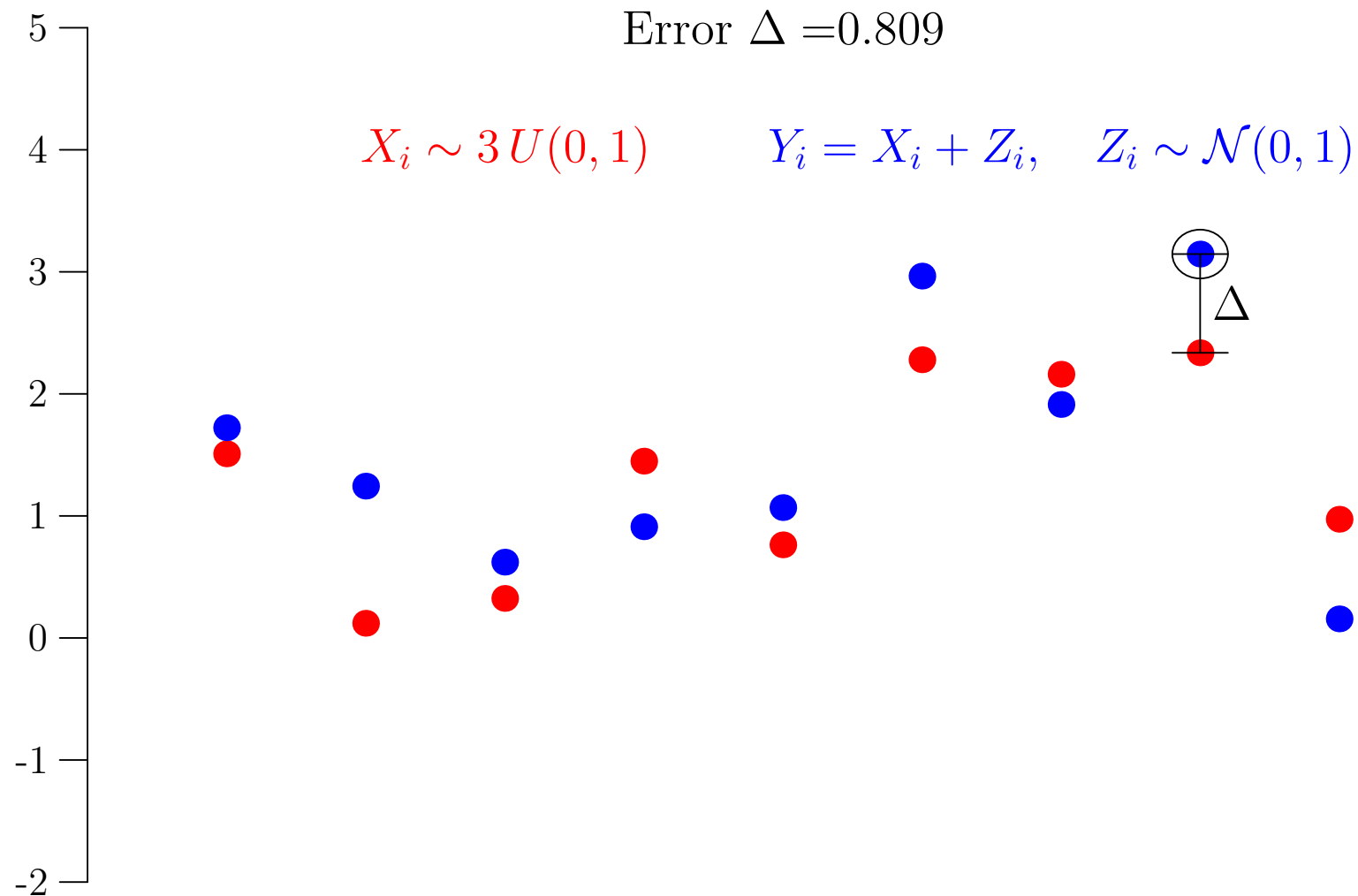
# The Over-fitting Game



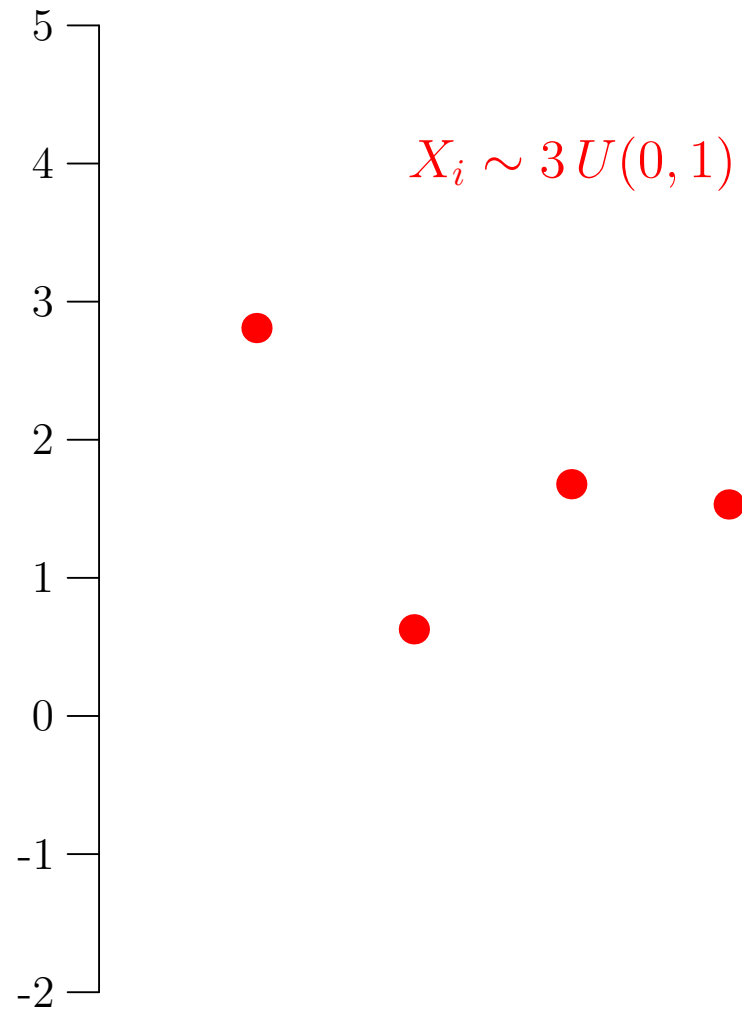
# The Over-fitting Game



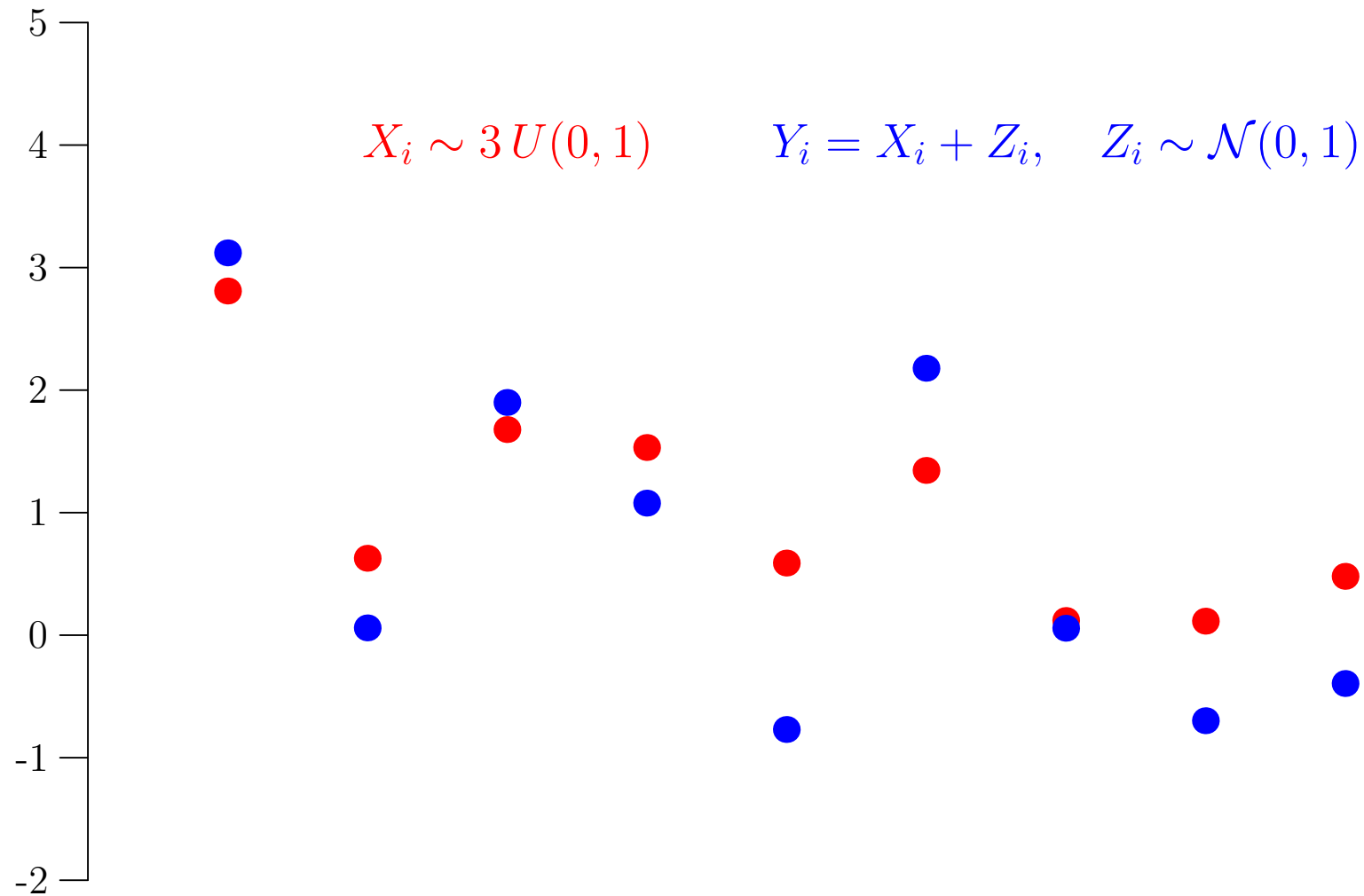
# The Over-fitting Game



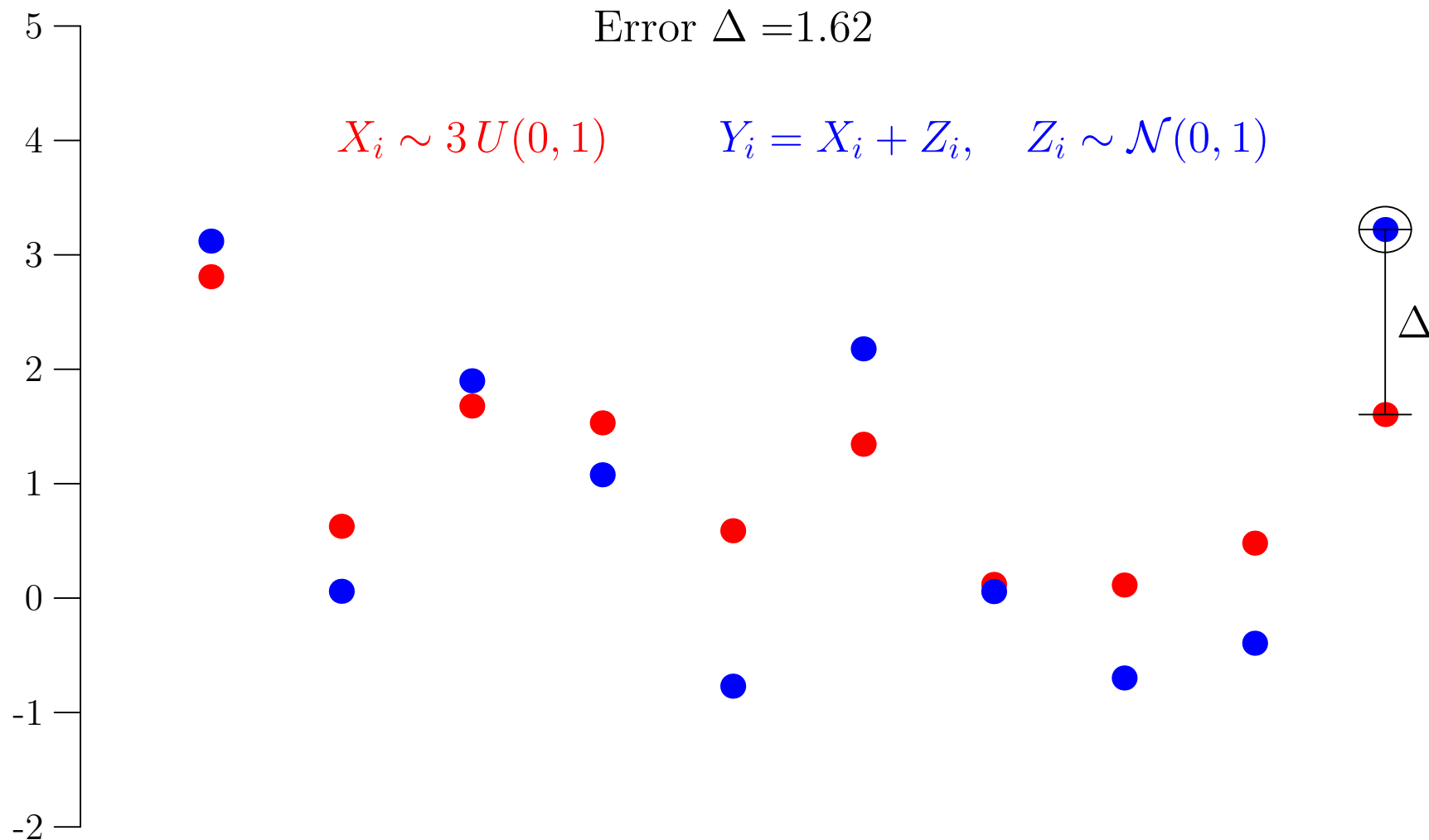
# The Over-fitting Game



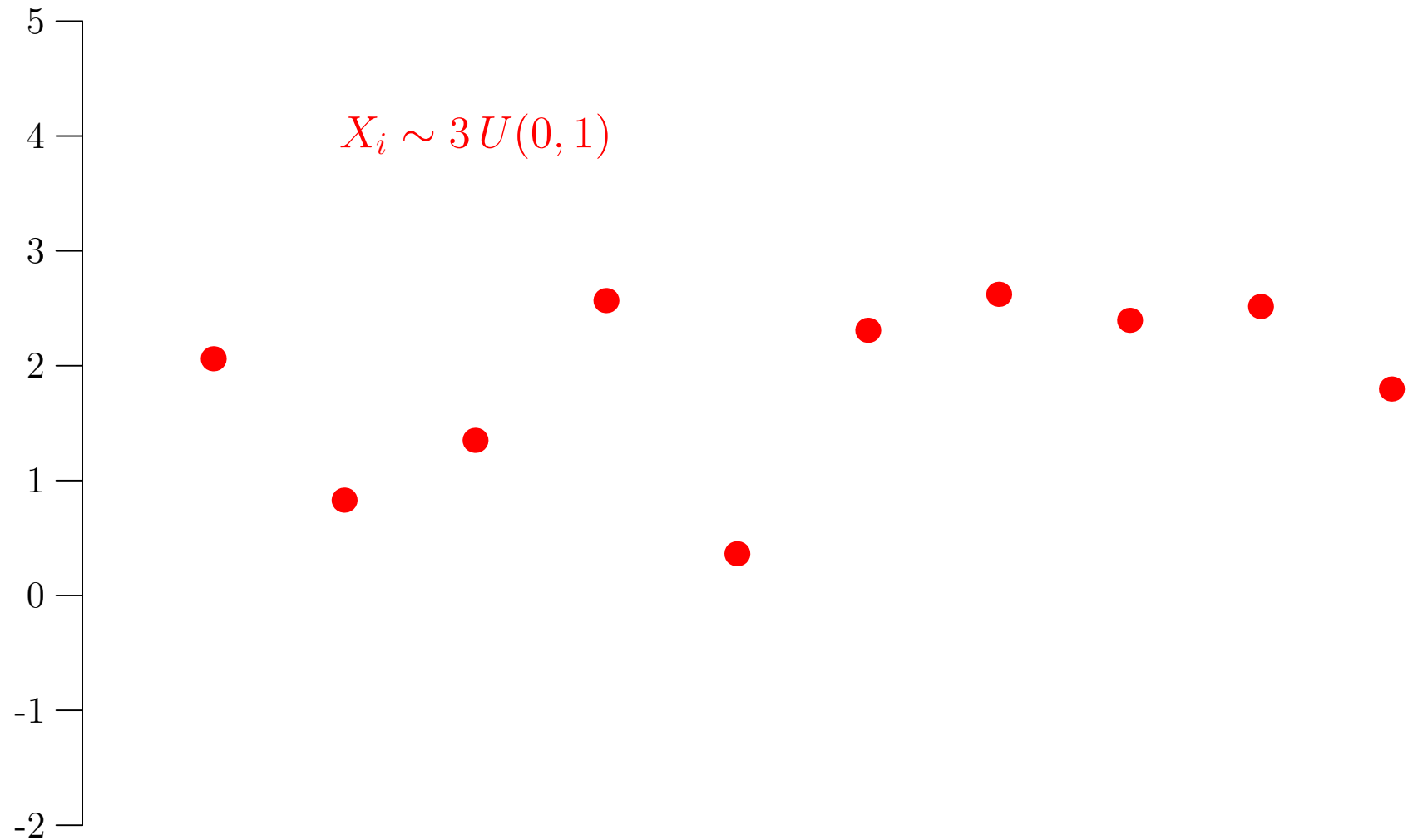
# The Over-fitting Game



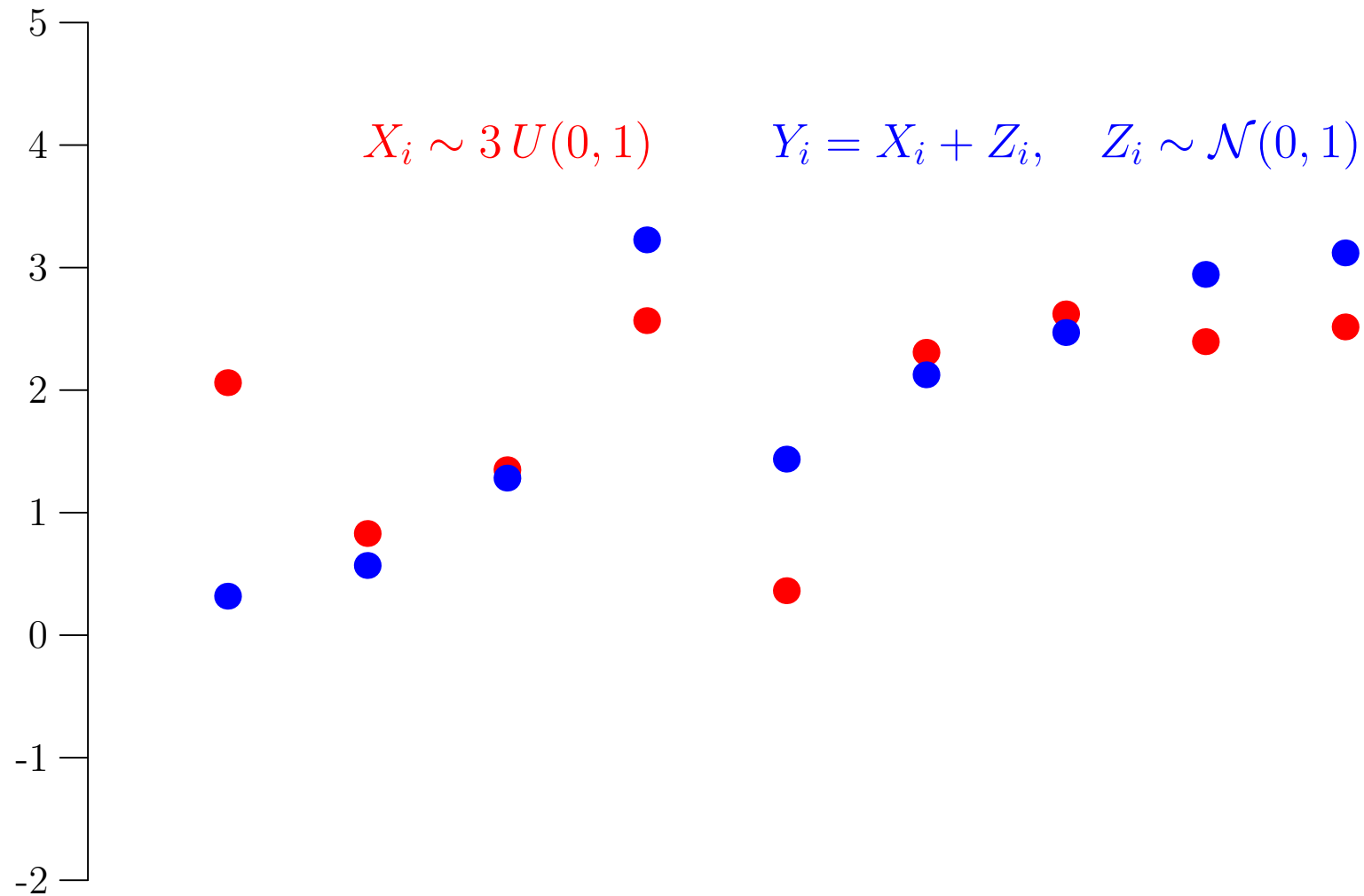
# The Over-fitting Game



# The Over-fitting Game

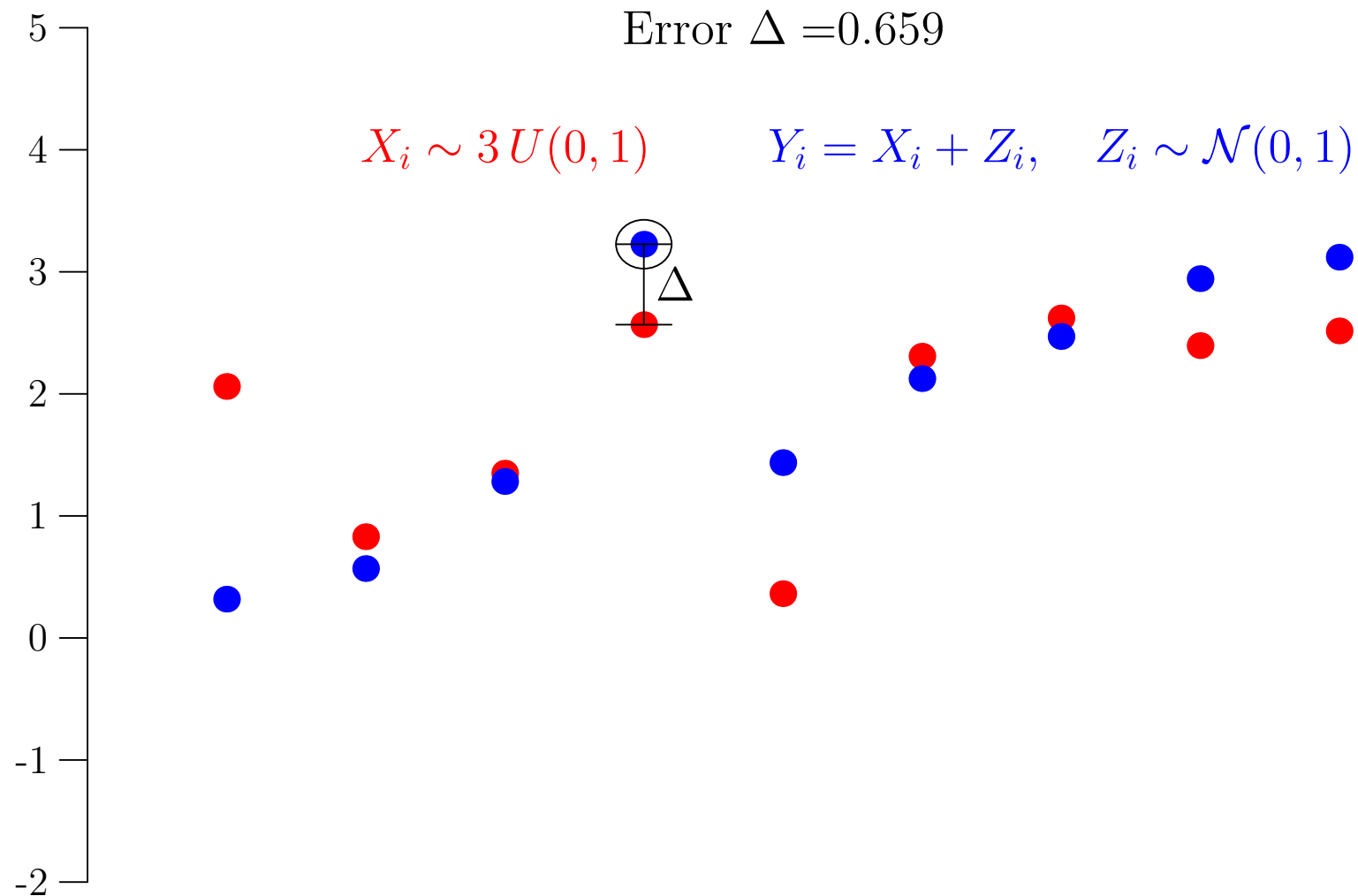


# The Over-fitting Game

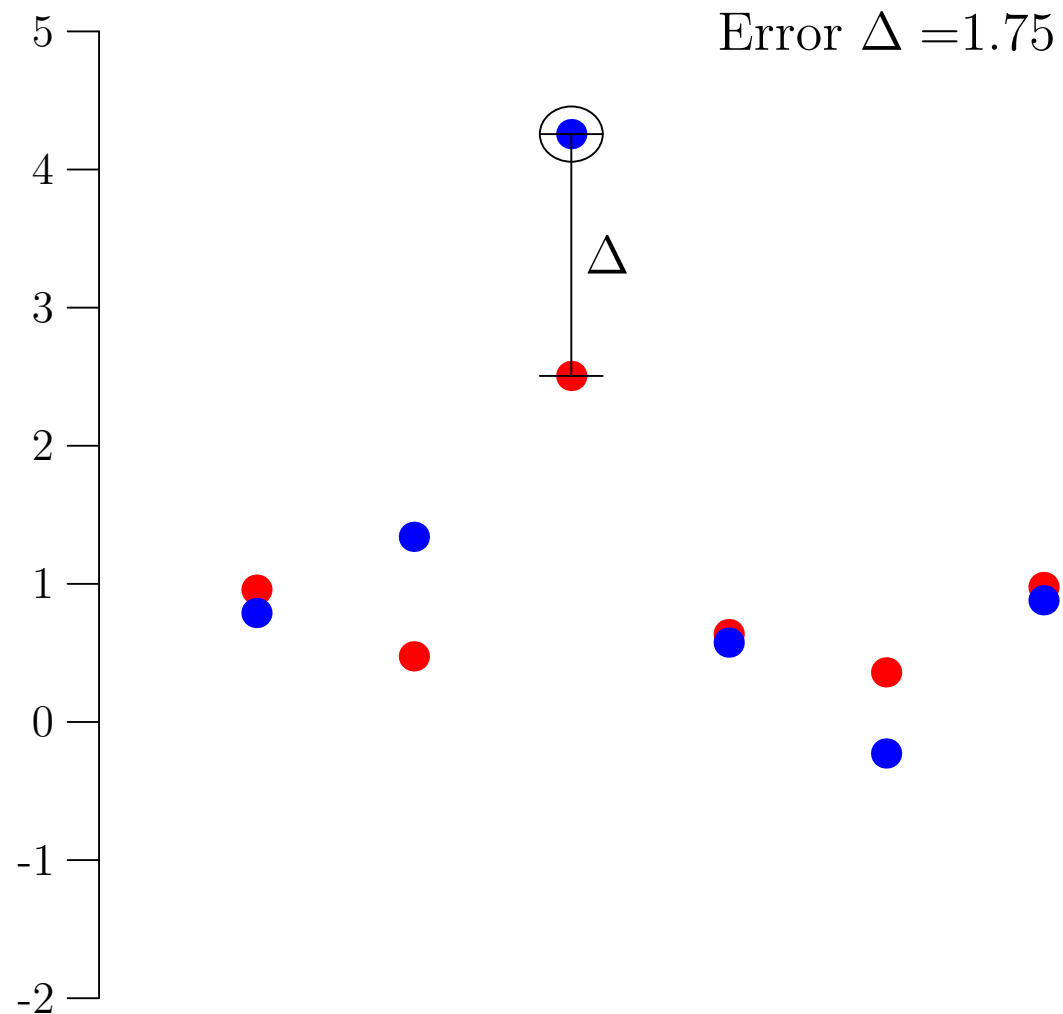




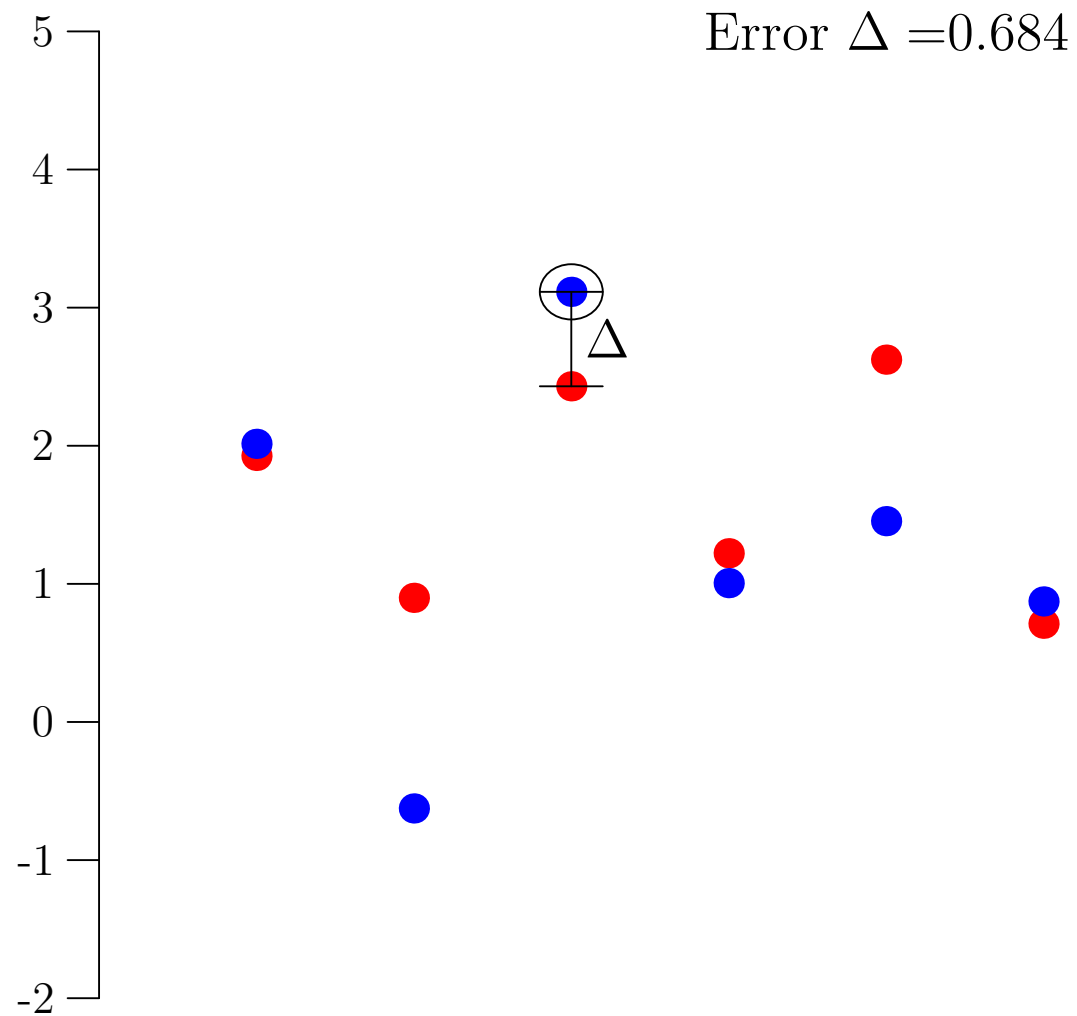
# The Over-fitting Game



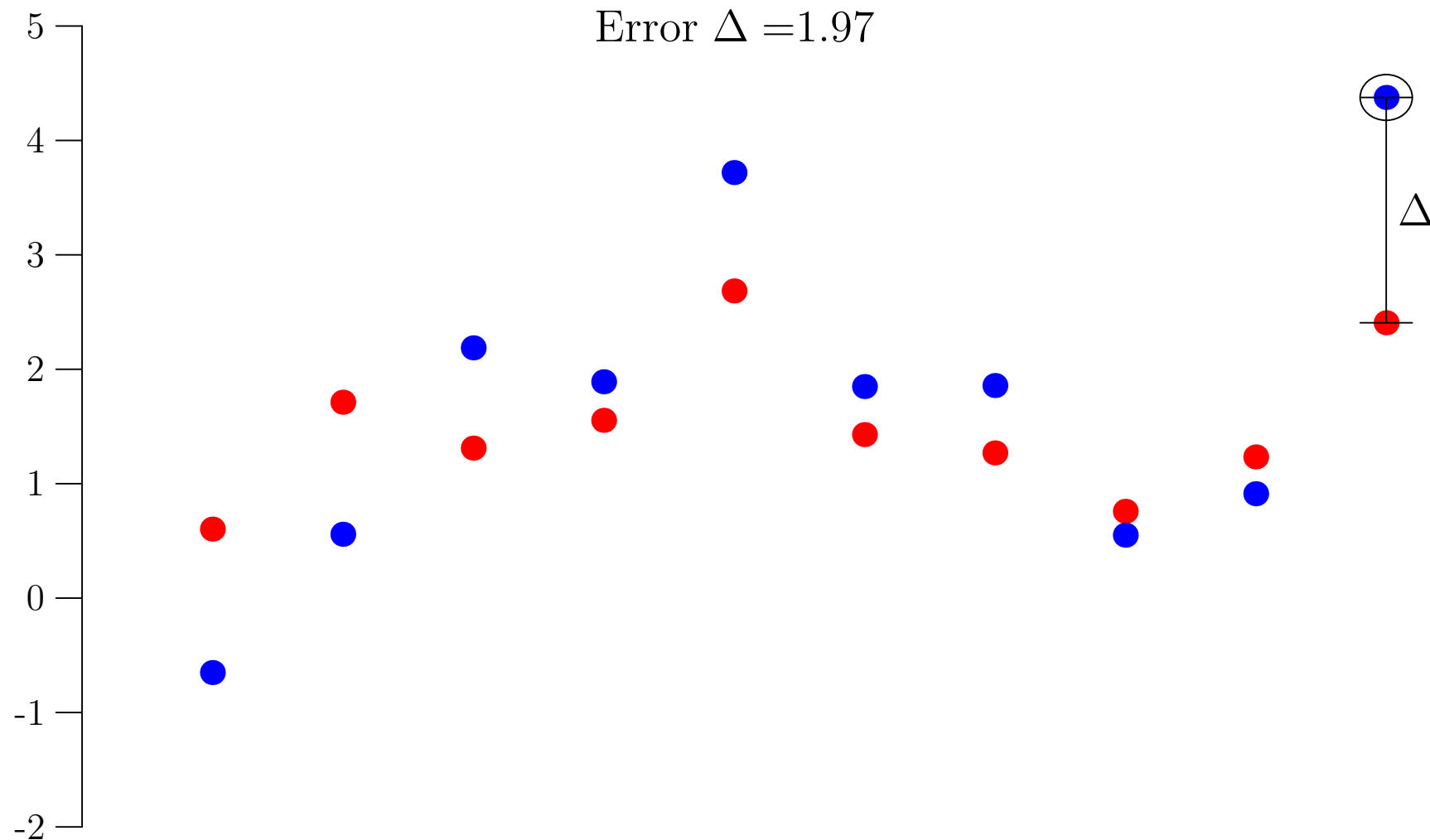
# The Over-fitting Game



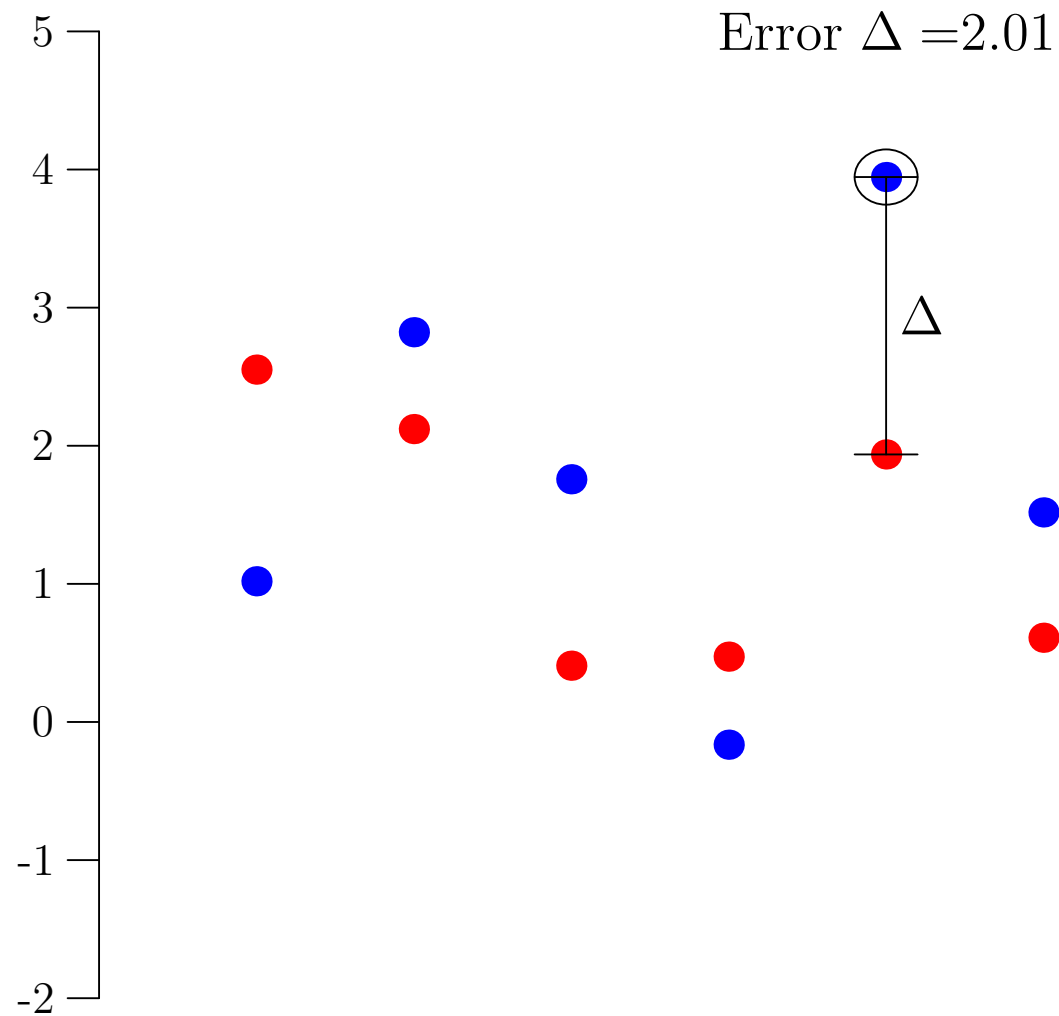
# The Over-fitting Game



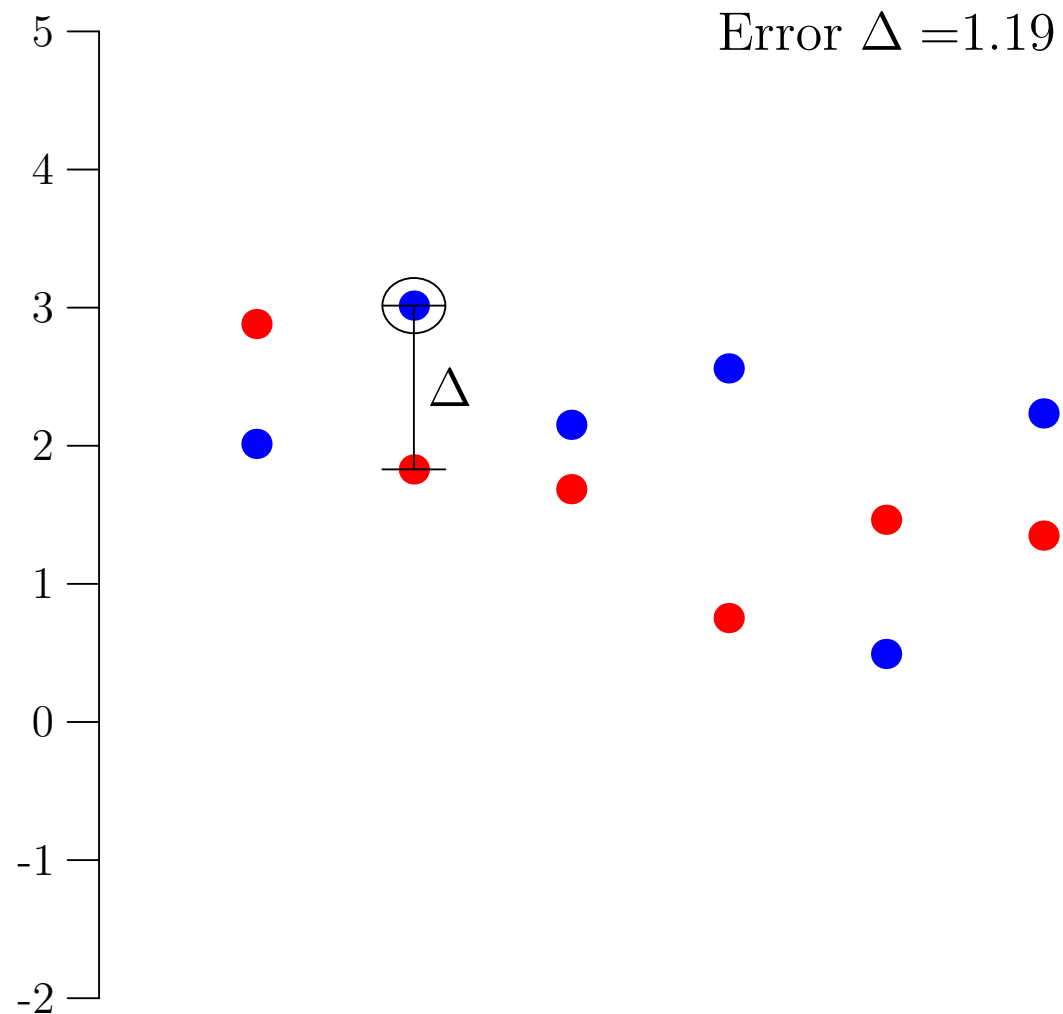
# The Over-fitting Game



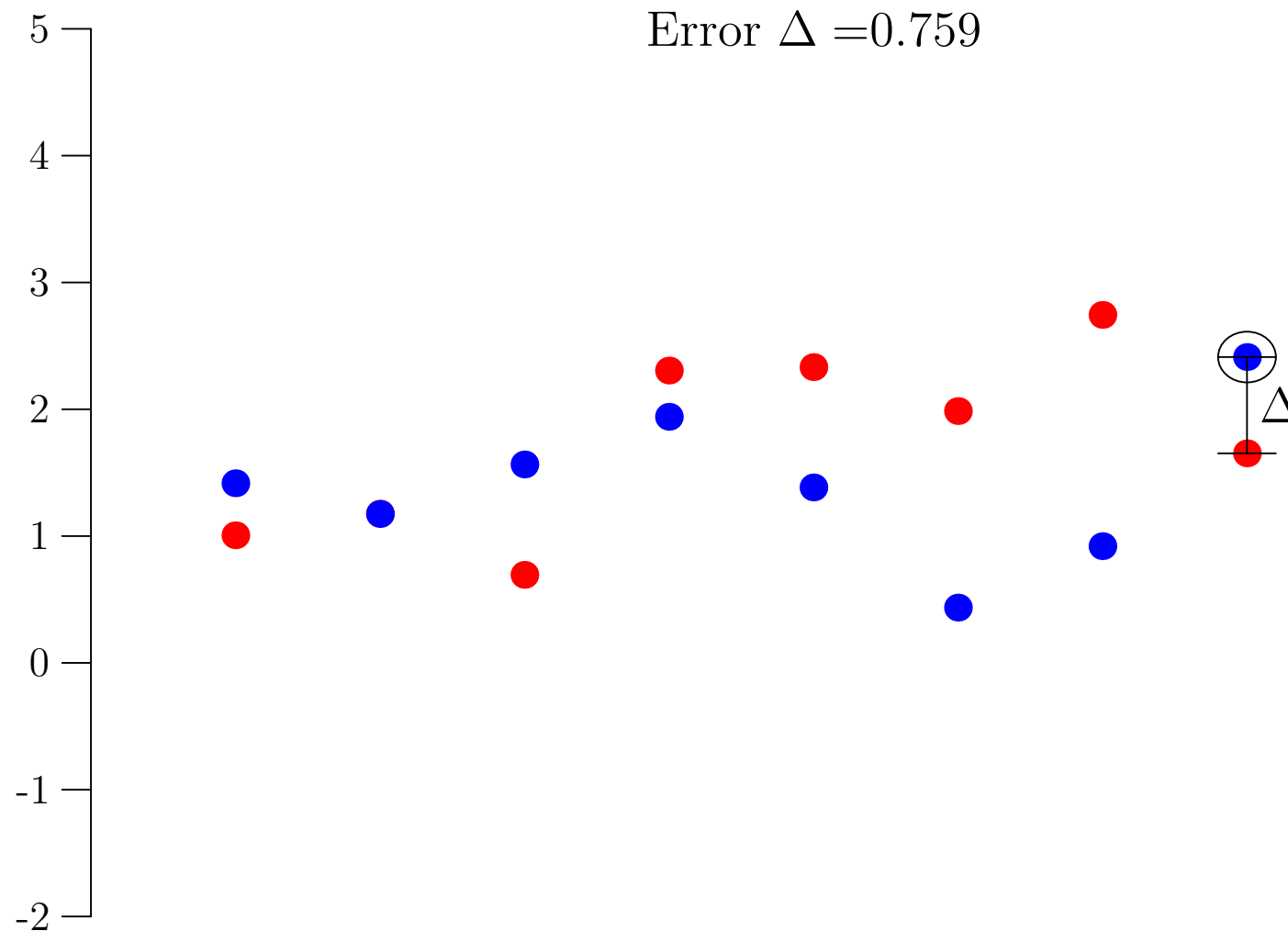
# The Over-fitting Game



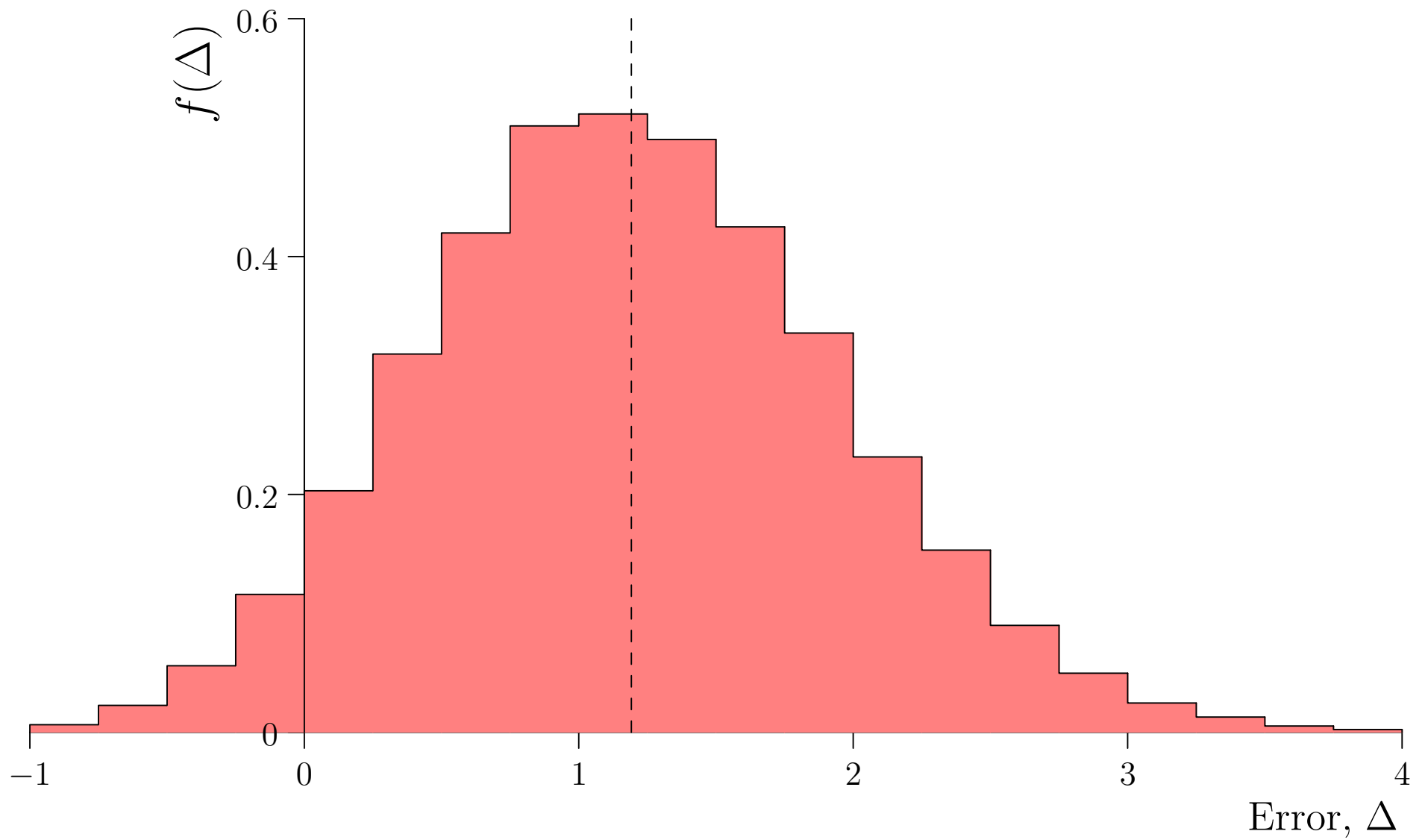
# The Over-fitting Game



# The Over-fitting Game

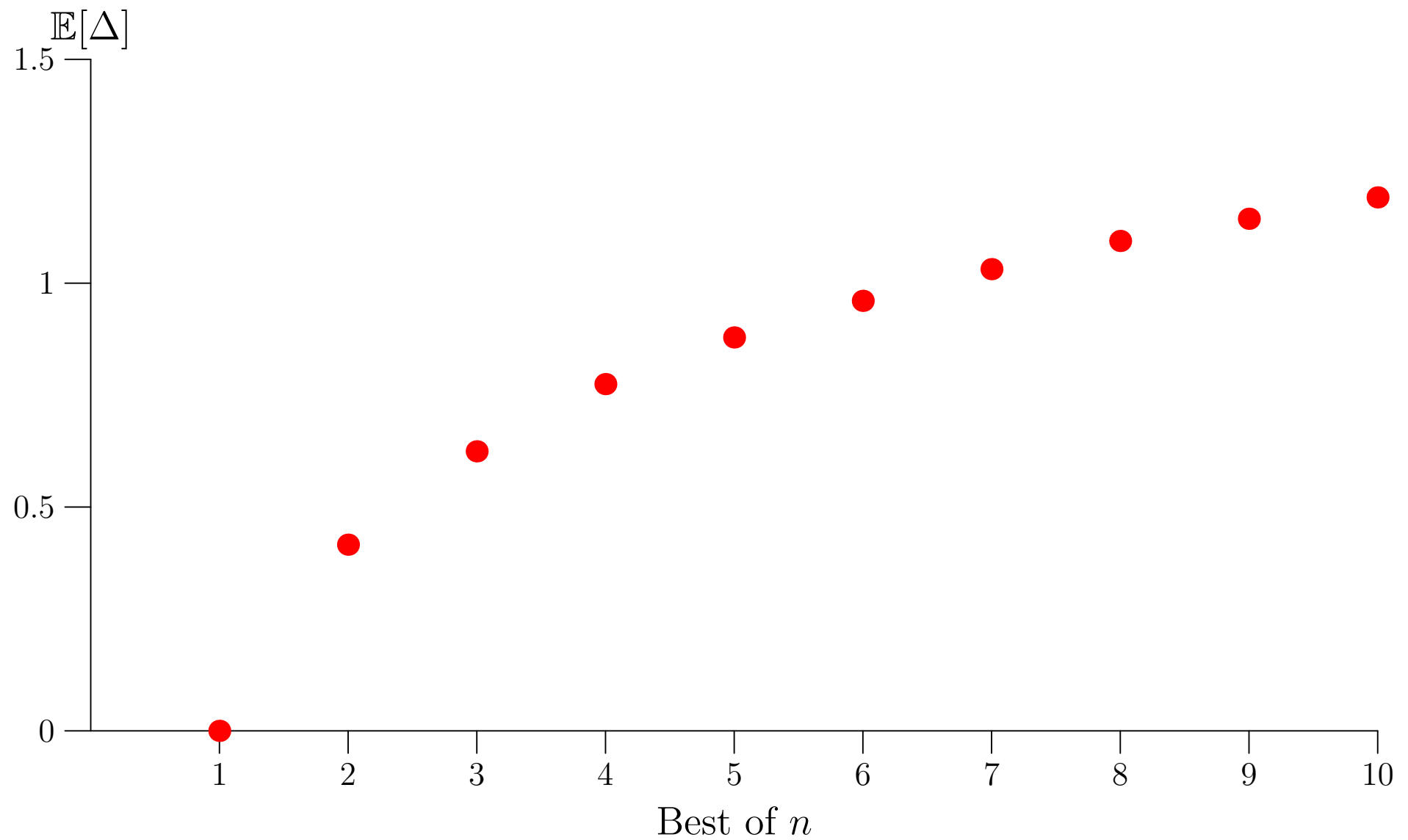


# The Over-fitting Game



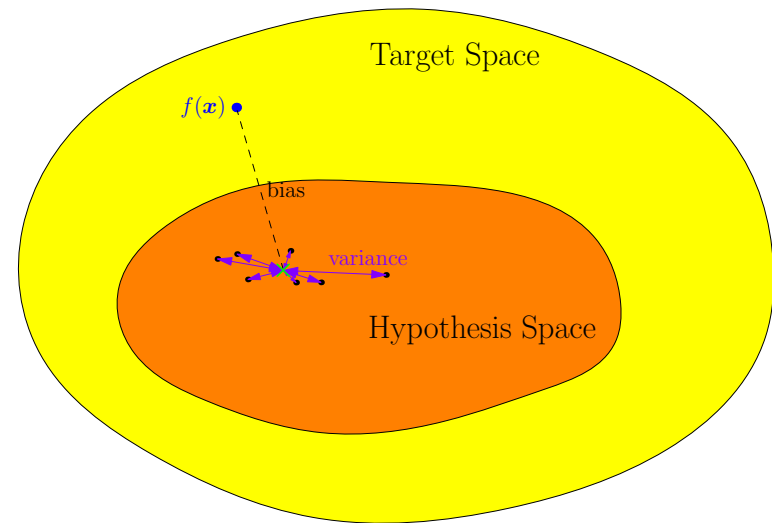


# The Over-fitting Game



# Outline

1. What Makes a Good Learning Machine?
2. **SVMs**
3. Ensemble Methods
4. Bayesian Inference



# Support Vector Machines

- Support vector machines, when used right, often have the best generalisation results
- They are typically used on numerical data, but can and have been adapted to text, sequences, etc.
- Although not as trendy as deep learning, they will often be the method of choice on small data sets
- They subtly regularise themselves, choosing a solution that generalises well from a host of different solutions

# Support Vector Machines

- Support vector machines, when used right, often have the best generalisation results
- They are typically used on numerical data, but can and have been adapted to text, sequences, etc.
- Although not as trendy as deep learning, they will often be the method of choice on small data sets
- They subtly regularise themselves, choosing a solution that generalises well from a host of different solutions

# Support Vector Machines

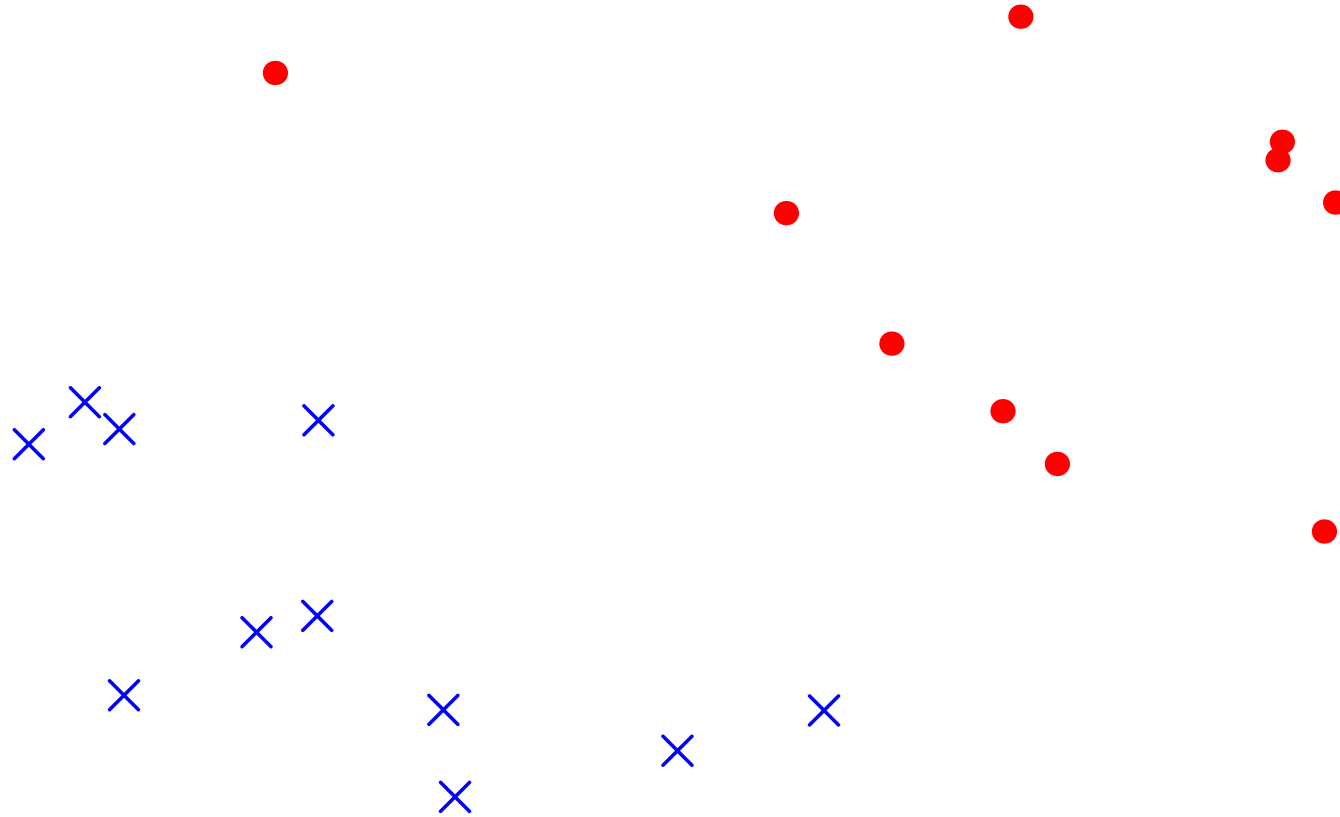
- Support vector machines, when used right, often have the best generalisation results
- They are typically used on numerical data, but can and have been adapted to text, sequences, etc.
- Although not as trendy as deep learning, they will often be the method of choice on small data sets
- They subtly regularise themselves, choosing a solution that generalises well from a host of different solutions

# Support Vector Machines

- Support vector machines, when used right, often have the best generalisation results
- They are typically used on numerical data, but can and have been adapted to text, sequences, etc.
- Although not as trendy as deep learning, they will often be the method of choice on small data sets
- They subtly regularise themselves, choosing a solution that generalises well from a host of different solutions

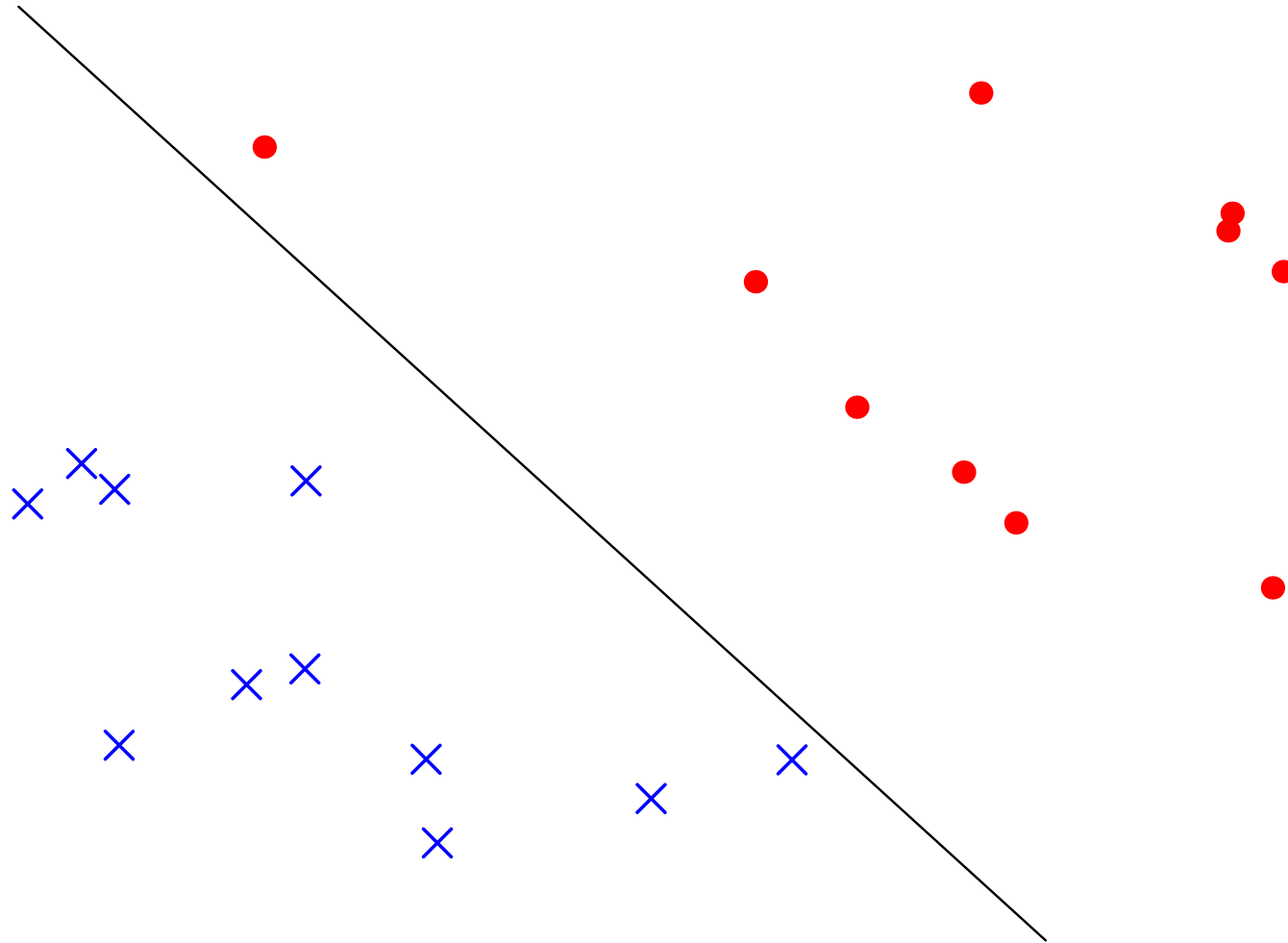
# Linear Separation of Data

- SVMs classify linearly separable data



# Linear Separation of Data

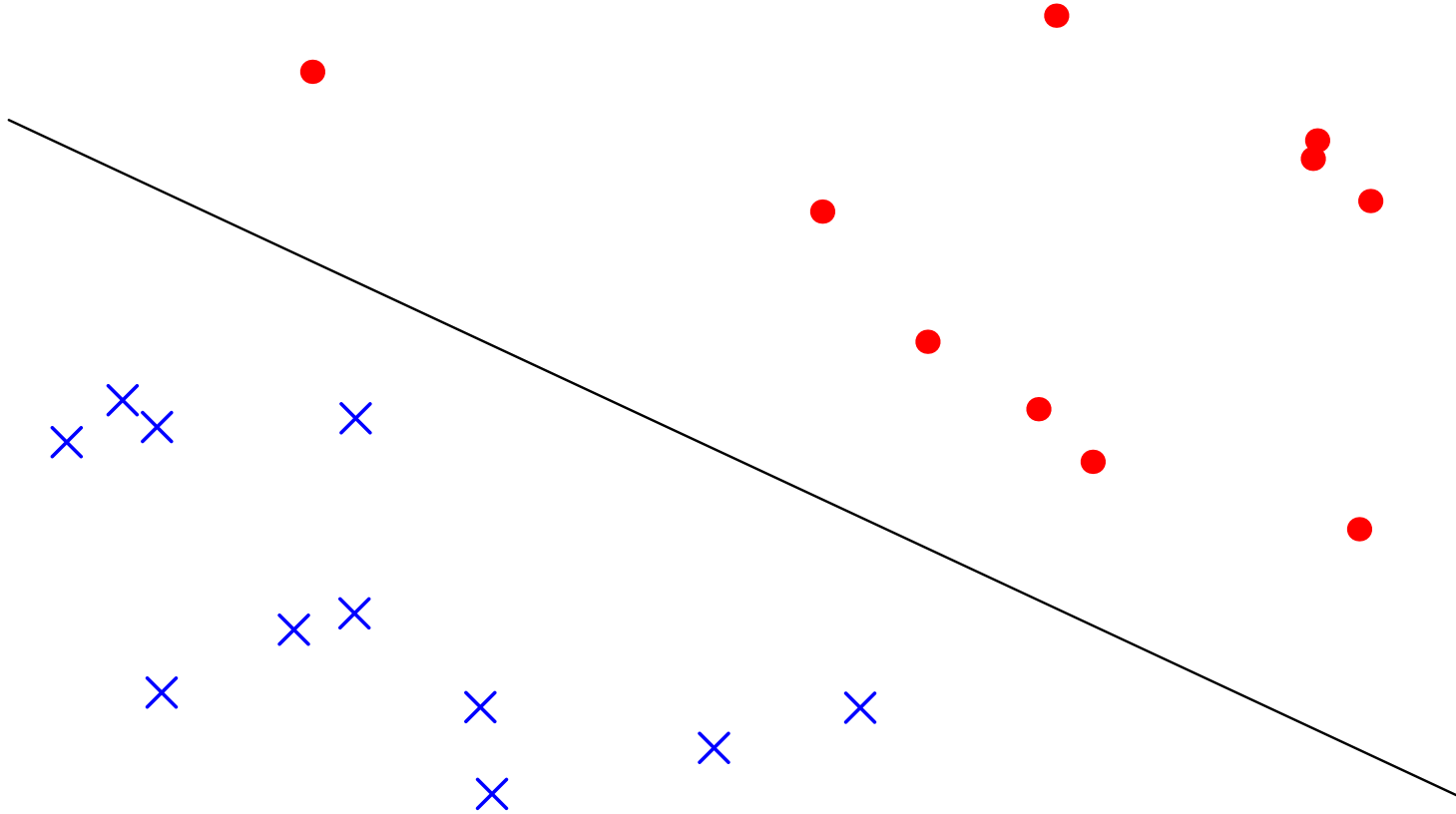
- SVMs classify linearly separable data





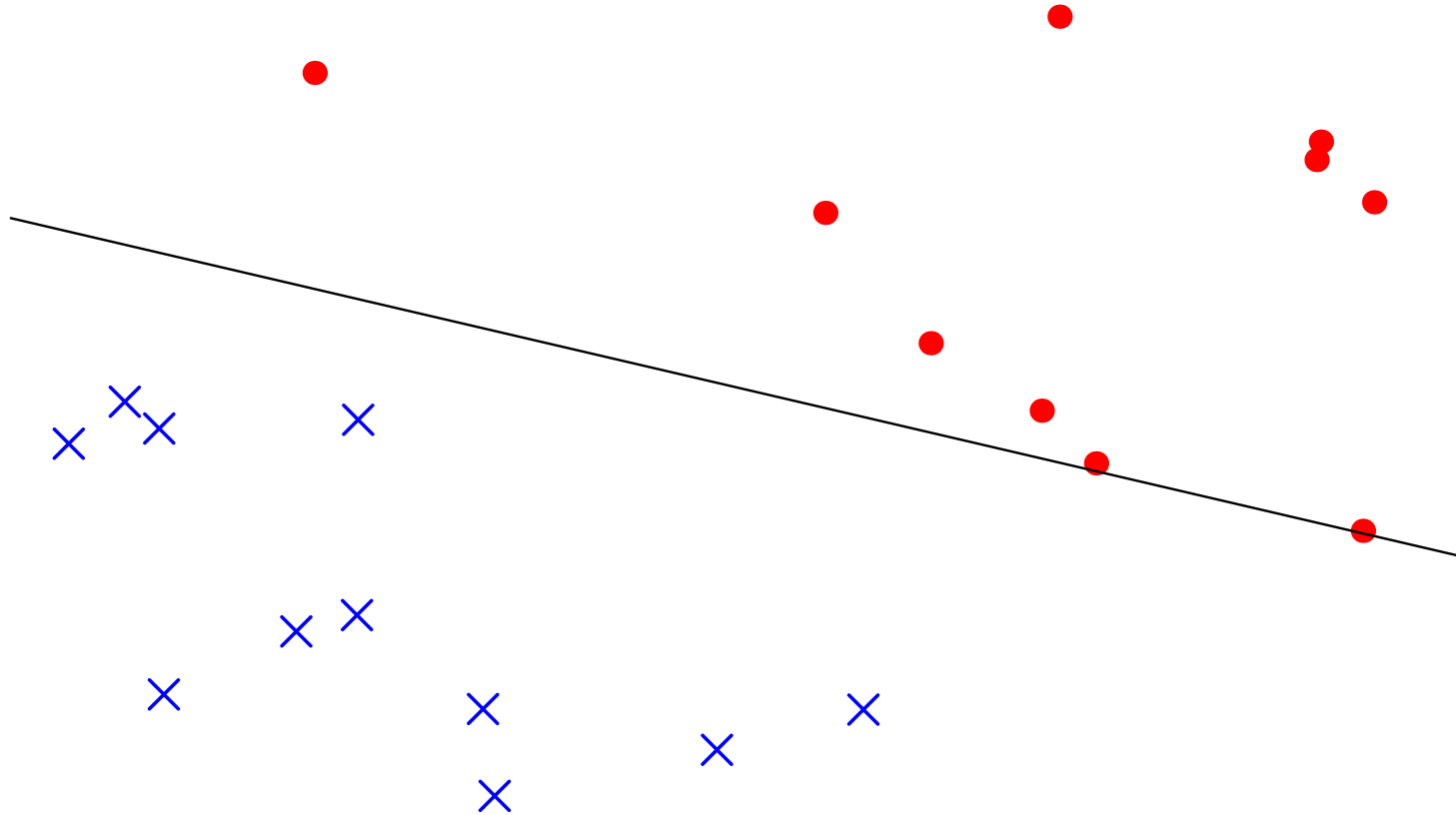
# Linear Separation of Data

- SVMs classify linearly separable data



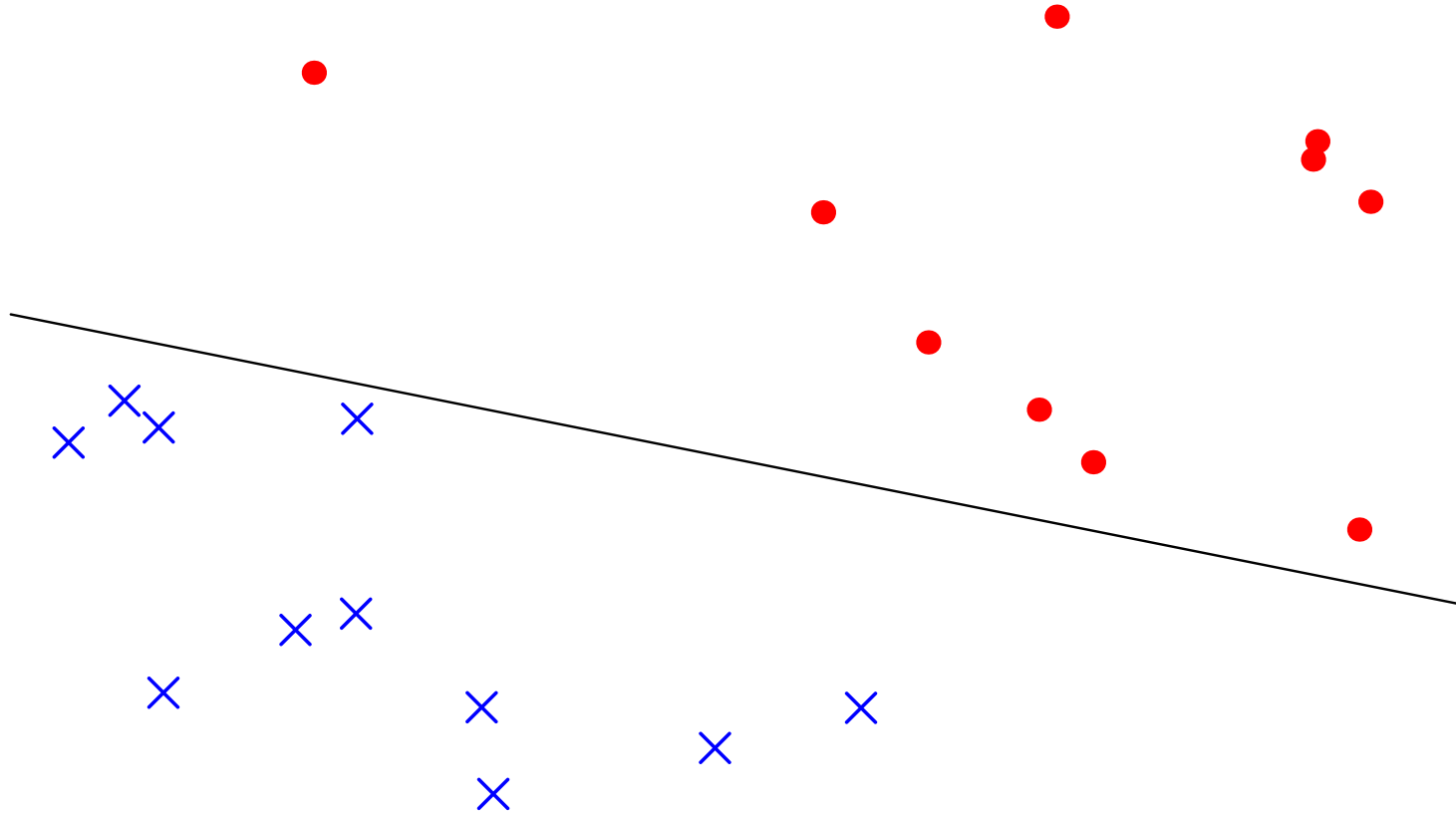
# Linear Separation of Data

- SVMs classify linearly separable data



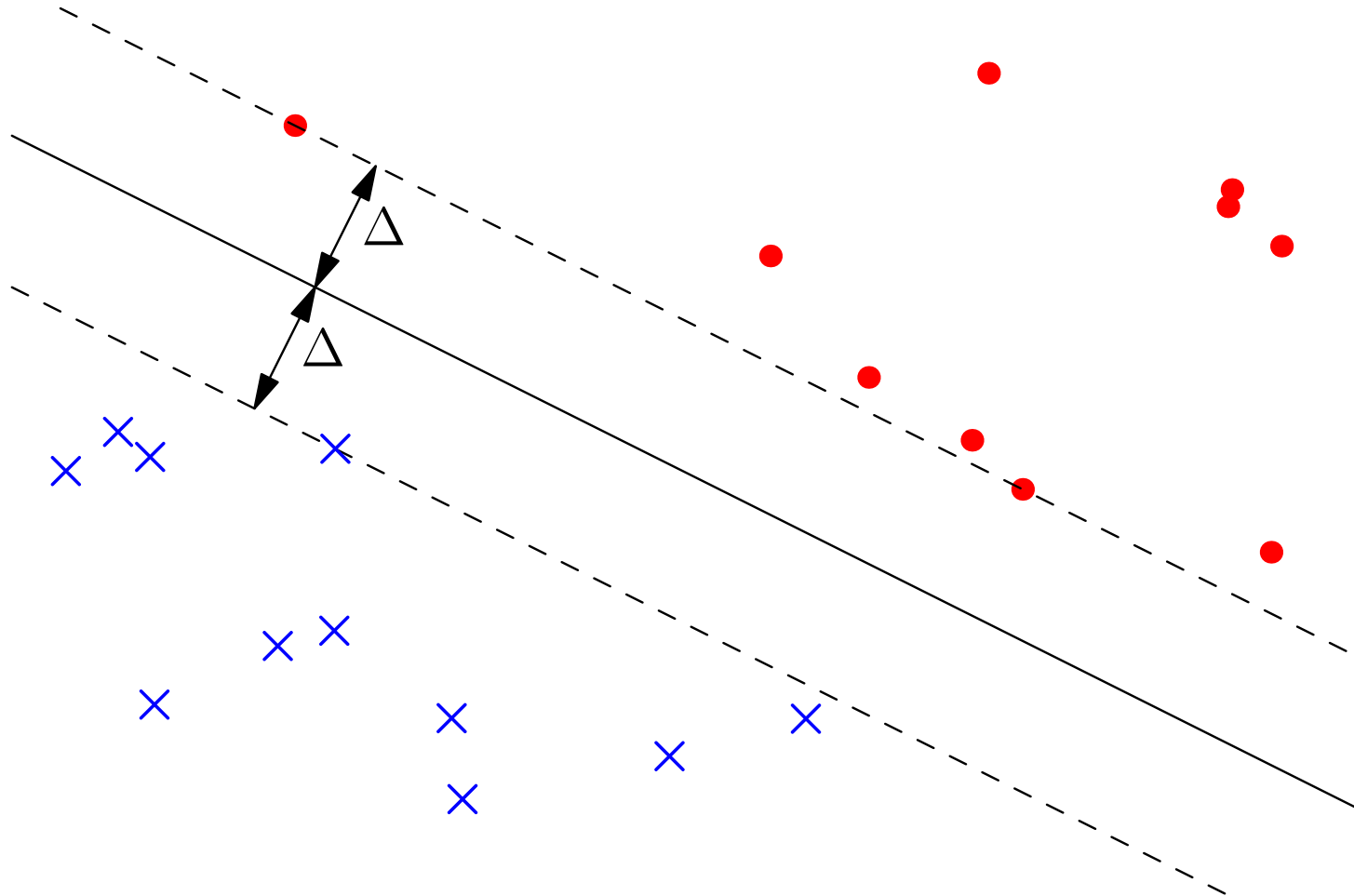
# Linear Separation of Data

- SVMs classify linearly separable data



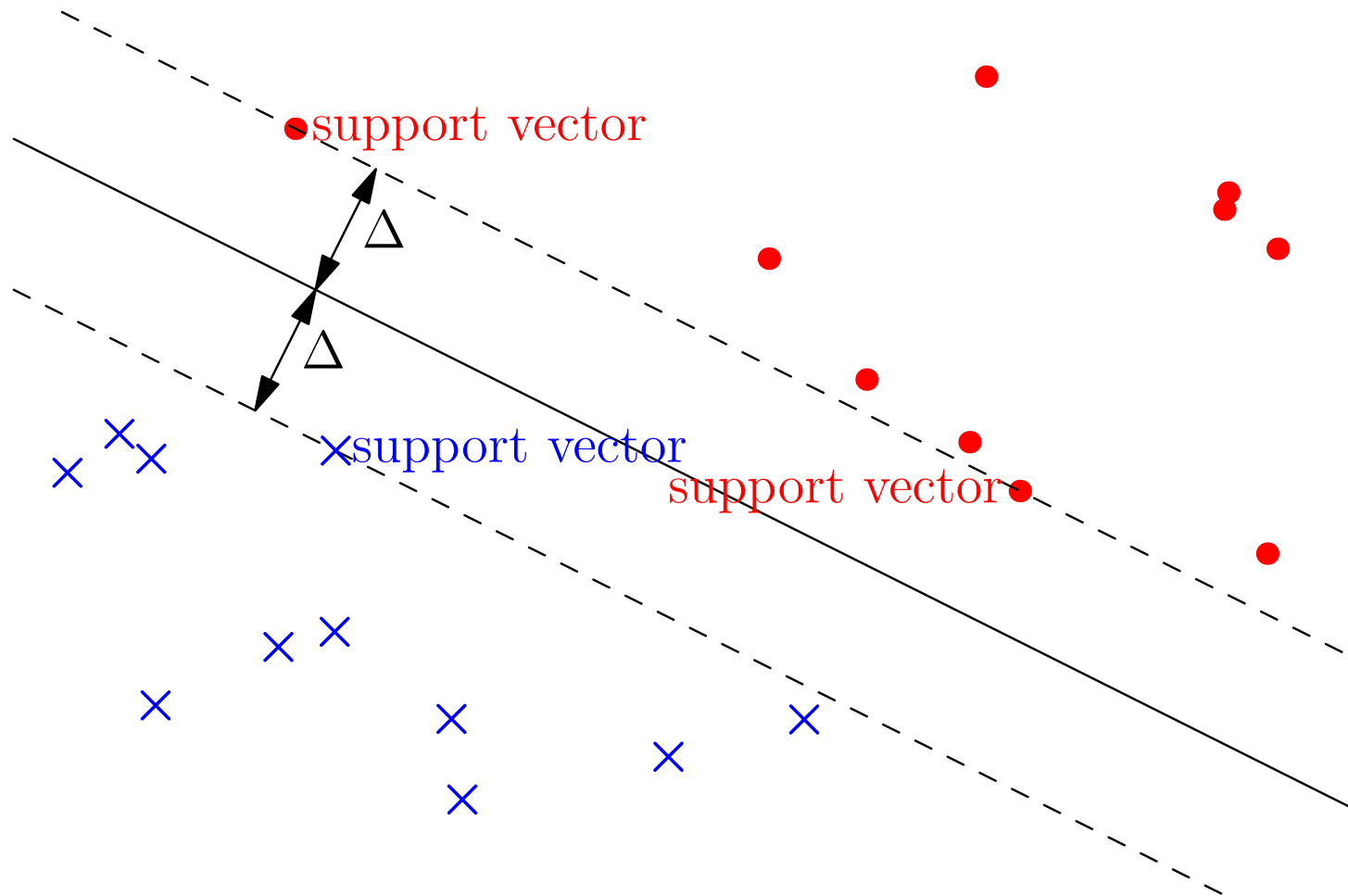
# Linear Separation of Data

- SVMs classify linearly separable data



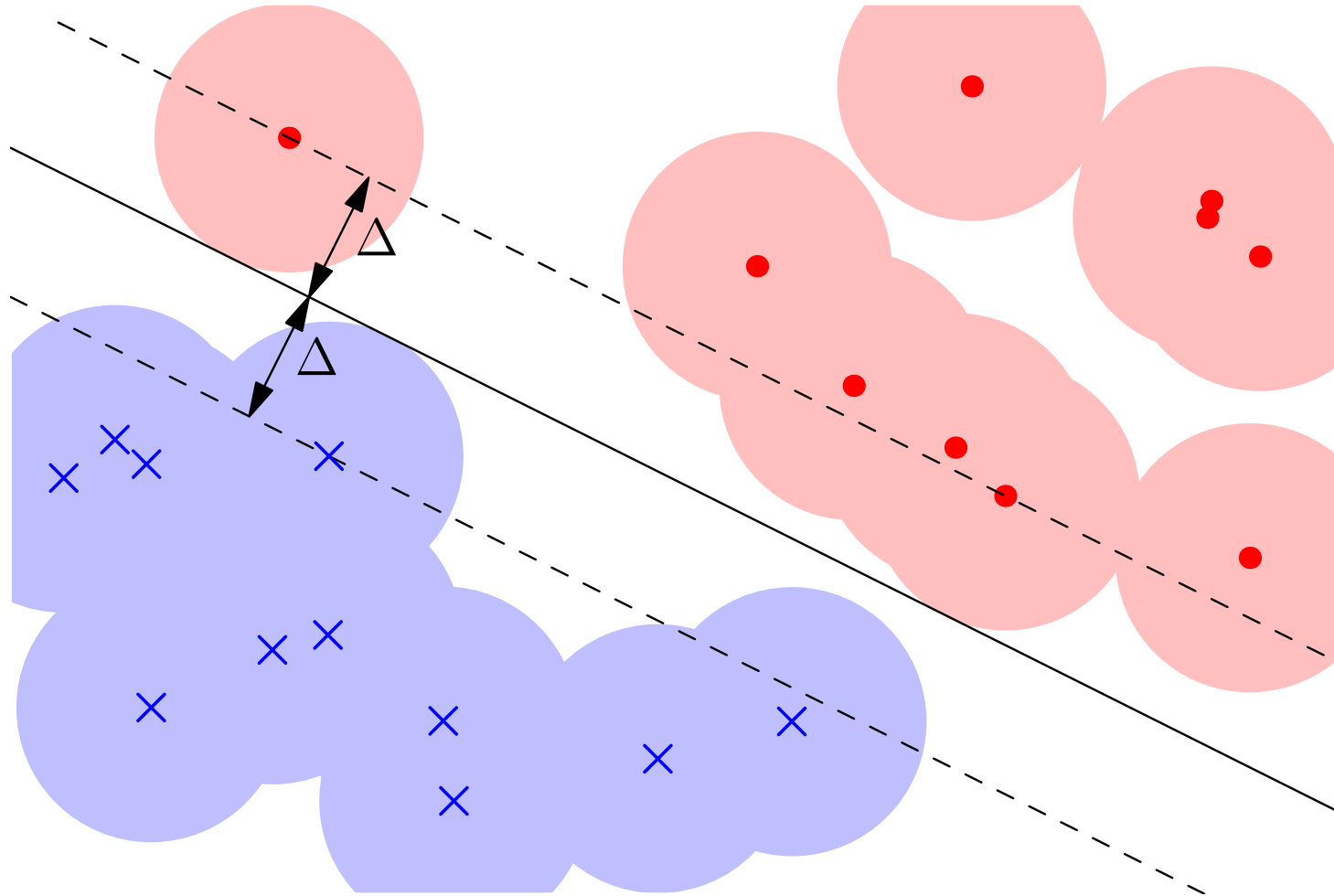
# Linear Separation of Data

- SVMs classify linearly separable data



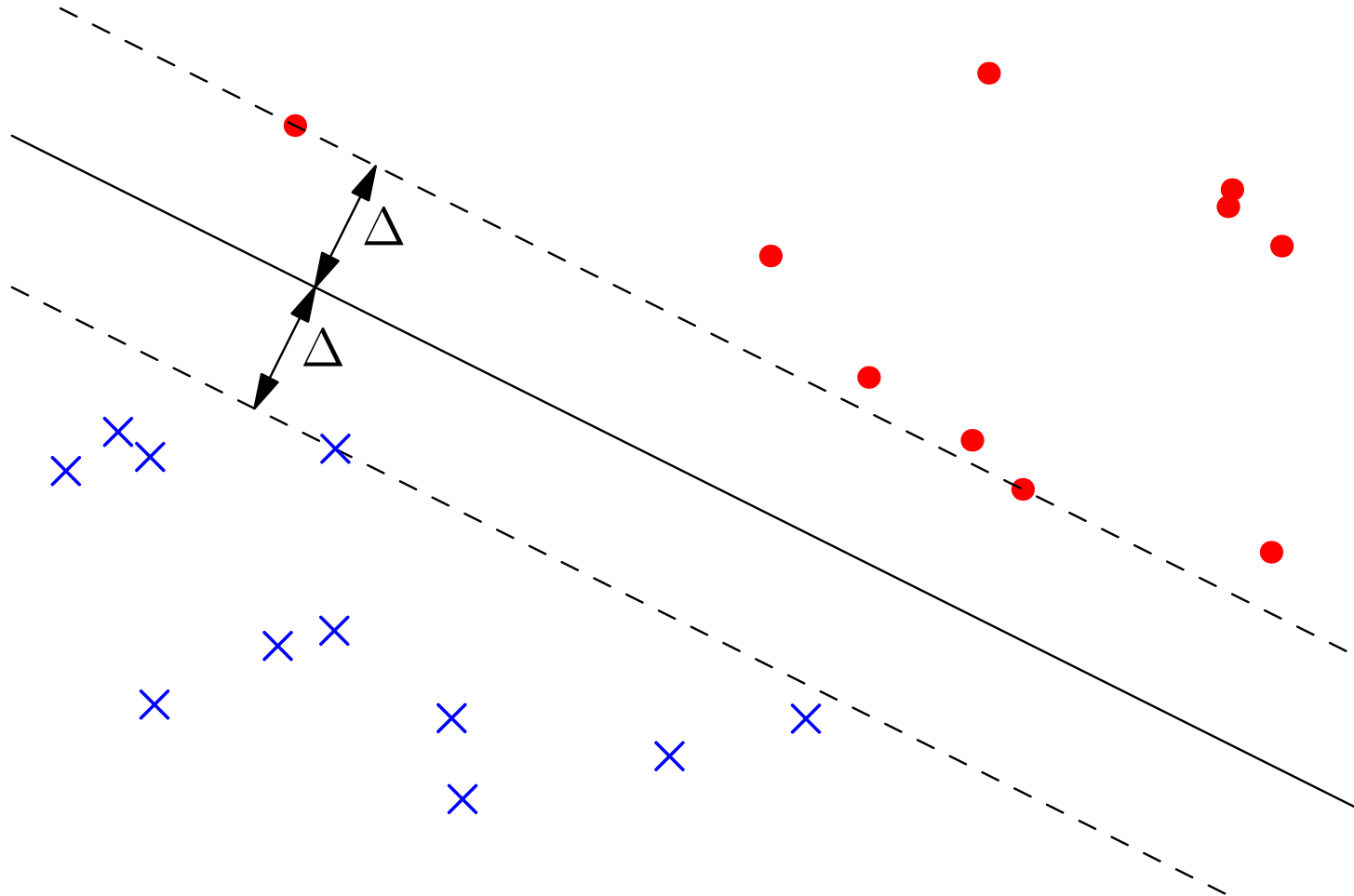
# Linear Separation of Data

- SVMs classify linearly separable data



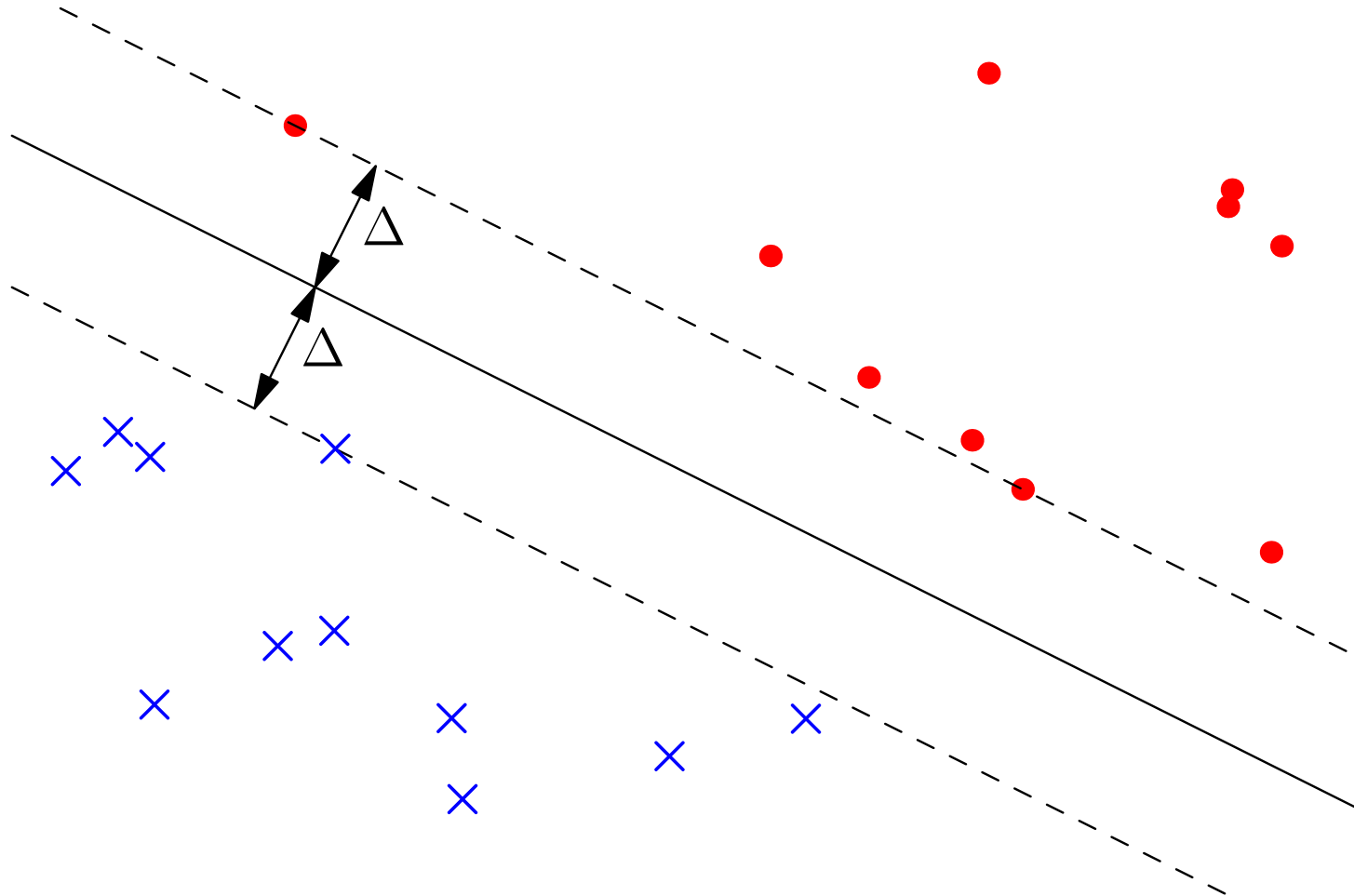
# Linear Separation of Data

- SVMs classify linearly separable data



# Linear Separation of Data

- SVMs classify linearly separable data



- Finds maximum-margin separating plane



# Extended Feature Space

- To increase the likelihood of linear-separability we often use a high-dimensional mapping

$$\mathbf{x} = (x_1, x_2, \dots, x_p) \rightarrow \phi(\mathbf{x}) = (\phi_1(\mathbf{x}), \phi_2(\mathbf{x}), \dots, \phi_{p'}(\mathbf{x}))$$

$$p' \gg p$$

- Finding the maximum margin hyper-plane is time consuming in “primal” form if  $p'$  is large
- We can work in the “dual” space of patterns, then we only need to compute dot products

$$\phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j)$$

# Extended Feature Space

- To increase the likelihood of linear-separability we often use a high-dimensional mapping

$$\mathbf{x} = (x_1, x_2, \dots, x_p) \rightarrow \phi(\mathbf{x}) = (\phi_1(\mathbf{x}), \phi_2(\mathbf{x}), \dots, \phi_{p'}(\mathbf{x}))$$

$$p' \gg p$$

- Finding the maximum margin hyper-plane is time consuming in “primal” form if  $p'$  is large
- We can work in the “dual” space of patterns, then we only need to compute dot products

$$\phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j)$$

# Extended Feature Space

- To increase the likelihood of linear-separability we often use a high-dimensional mapping

$$\mathbf{x} = (x_1, x_2, \dots, x_p) \rightarrow \phi(\mathbf{x}) = (\phi_1(\mathbf{x}), \phi_2(\mathbf{x}), \dots, \phi_{p'}(\mathbf{x}))$$

$$p' \gg p$$

- Finding the maximum margin hyper-plane is time consuming in “primal” form if  $p'$  is large
- We can work in the “dual” space of patterns, then we only need to compute dot products

$$\phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j)$$

# Extended Feature Space

- To increase the likelihood of linear-separability we often use a high-dimensional mapping

$$\mathbf{x} = (x_1, x_2, \dots, x_p) \rightarrow \phi(\mathbf{x}) = (\phi_1(\mathbf{x}), \phi_2(\mathbf{x}), \dots, \phi_{p'}(\mathbf{x}))$$

$$p' \gg p$$

- Finding the maximum margin hyper-plane is time consuming in “primal” form if  $p'$  is large
- We can work in the “dual” space of patterns, then we only need to compute dot products

$$\phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j) = \sum_{k=1}^{p'} \phi_k(\mathbf{x}_i) \phi_k(\mathbf{x}_j)$$

# Kernel Trick

- If we choose a **positive semi-definite** kernel function  $K(\mathbf{x}, \mathbf{y})$  then there exists functions  $\phi_k(\mathbf{x})$ , such that

$$K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j)$$

(like an eigenvector decomposition of a matrix)

- Never need to compute  $\phi_k(\mathbf{x}_i)$  explicitly as we only need the dot-product  $\phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j) = K(\mathbf{x}_i, \mathbf{x}_j)$  to compute maximum margin separating hyper-plane
- Sometimes  $\phi(\mathbf{x}_i)$  is an infinite dimensional vector so its good we don't have to compute it!

# Kernel Trick

- If we choose a **positive semi-definite** kernel function  $K(\mathbf{x}, \mathbf{y})$  then there exists functions  $\phi_k(\mathbf{x})$ , such that

$$K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j)$$

(like an eigenvector decomposition of a matrix)

- Never need to compute  $\phi_k(\mathbf{x}_i)$  explicitly as we only need the dot-product  $\phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j) = K(\mathbf{x}_i, \mathbf{x}_j)$  to compute maximum margin separating hyper-plane
- Sometimes  $\phi(\mathbf{x}_i)$  is an infinite dimensional vector so its good we don't have to compute it!

# Kernel Trick

- If we choose a **positive semi-definite** kernel function  $K(\mathbf{x}, \mathbf{y})$  then there exists functions  $\phi_k(\mathbf{x})$ , such that

$$K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j)$$

(like an eigenvector decomposition of a matrix)

- Never need to compute  $\phi_k(\mathbf{x}_i)$  explicitly as we only need the dot-product  $\phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j) = K(\mathbf{x}_i, \mathbf{x}_j)$  to compute maximum margin separating hyper-plane
- Sometimes  $\phi(\mathbf{x}_i)$  is an infinite dimensional vector so its good we don't have to compute it!

# Kernel Trick

- If we choose a **positive semi-definite** kernel function  $K(\mathbf{x}, \mathbf{y})$  then there exists functions  $\phi_k(\mathbf{x})$ , such that

$$K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j)$$

(like an eigenvector decomposition of a matrix)

- Never need to compute  $\phi_k(\mathbf{x}_i)$  explicitly as we only need the dot-product  $\phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j) = K(\mathbf{x}_i, \mathbf{x}_j)$  to compute maximum margin separating hyper-plane
- Sometimes  $\phi(\mathbf{x}_i)$  is an infinite dimensional vector so its good we don't have to compute it!



# Kernel Functions

- Kernel functions are symmetric functions of two variable
- Strong restriction: *positive semi-definite*
- Examples

Quadratic kernel:  $K(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1^\top \mathbf{x}_2)^2$

Gaussian (RBF) kernel:  $K(\mathbf{x}_1, \mathbf{x}_2) = e^{-\gamma \|\mathbf{x}_1 - \mathbf{x}_2\|^2}$

- Consider the mapping

$$\mathbf{x}_i = \begin{pmatrix} x_i \\ y_i \end{pmatrix} \rightarrow \phi(\mathbf{x}_i) = \begin{pmatrix} x_i^2 \\ y_i^2 \\ \sqrt{2} x_i y_i \end{pmatrix}$$

# Kernel Functions

- Kernel functions are symmetric functions of two variable
- Strong restriction: *positive semi-definite*
- Examples

Quadratic kernel:  $K(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1^\top \mathbf{x}_2)^2$

Gaussian (RBF) kernel:  $K(\mathbf{x}_1, \mathbf{x}_2) = e^{-\gamma \|\mathbf{x}_1 - \mathbf{x}_2\|^2}$

- Consider the mapping

$$\mathbf{x}_i = \begin{pmatrix} x_i \\ y_i \end{pmatrix} \rightarrow \phi(\mathbf{x}_i) = \begin{pmatrix} x_i^2 \\ y_i^2 \\ \sqrt{2} x_i y_i \end{pmatrix}$$

# Kernel Functions

- Kernel functions are symmetric functions of two variable
- Strong restriction: *positive semi-definite*
- Examples

Quadratic kernel:  $K(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1^\top \mathbf{x}_2)^2$

Gaussian (RBF) kernel:  $K(\mathbf{x}_1, \mathbf{x}_2) = e^{-\gamma \|\mathbf{x}_1 - \mathbf{x}_2\|^2}$

- Consider the mapping

$$\mathbf{x}_i = \begin{pmatrix} x_i \\ y_i \end{pmatrix} \rightarrow \phi(\mathbf{x}_i) = \begin{pmatrix} x_i^2 \\ y_i^2 \\ \sqrt{2} x_i y_i \end{pmatrix}$$

# Kernel Functions

- Kernel functions are symmetric functions of two variable
- Strong restriction: *positive semi-definite*
- Examples

Quadratic kernel:  $K(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1^\top \mathbf{x}_2)^2$

Gaussian (RBF) kernel:  $K(\mathbf{x}_1, \mathbf{x}_2) = e^{-\gamma \|\mathbf{x}_1 - \mathbf{x}_2\|^2}$

- Consider the mapping

$$\mathbf{x}_i = \begin{pmatrix} x_i \\ y_i \end{pmatrix} \rightarrow \phi(\mathbf{x}_i) = \begin{pmatrix} x_i^2 \\ y_i^2 \\ \sqrt{2} x_i y_i \end{pmatrix}$$

# Non-linearly Separation of Data

$$K(\boldsymbol{x}_1, \boldsymbol{x}_2) = \boldsymbol{\phi}^\top(\boldsymbol{x}_1)\boldsymbol{\phi}(\boldsymbol{x}_2)$$

# Non-linearly Separation of Data

$$K(\mathbf{x}_1, \mathbf{x}_2) = \phi^\top(\mathbf{x}_1)\phi(\mathbf{x}_2) = \begin{pmatrix} x_1^2 & y_1^2 & \sqrt{2} x_1 y_1 \end{pmatrix} \begin{pmatrix} x_2^2 \\ y_2^2 \\ \sqrt{2} x_2 y_2 \end{pmatrix}$$

# Non-linearly Separation of Data

$$\begin{aligned} K(\mathbf{x}_1, \mathbf{x}_2) &= \boldsymbol{\phi}^\top(\mathbf{x}_1)\boldsymbol{\phi}(\mathbf{x}_2) = \begin{pmatrix} x_1^2 & y_1^2 & \sqrt{2} x_1 y_1 \end{pmatrix} \begin{pmatrix} x_2^2 \\ y_2^2 \\ \sqrt{2} x_2 y_2 \end{pmatrix} \\ &= x_1^2 x_2^2 + y_1^2 y_2^2 + 2 x_1 y_1 x_2 y_2 \end{aligned}$$

# Non-linearly Separation of Data

$$\begin{aligned} K(\mathbf{x}_1, \mathbf{x}_2) &= \boldsymbol{\phi}^\top(\mathbf{x}_1)\boldsymbol{\phi}(\mathbf{x}_2) = \begin{pmatrix} x_1^2 & y_1^2 & \sqrt{2} x_1 y_1 \end{pmatrix} \begin{pmatrix} x_2^2 \\ y_2^2 \\ \sqrt{2} x_2 y_2 \end{pmatrix} \\ &= x_1^2 x_2^2 + y_1^2 y_2^2 + 2 x_1 y_1 x_2 y_2 = (x_1 x_2 + y_1 y_2)^2 \end{aligned}$$

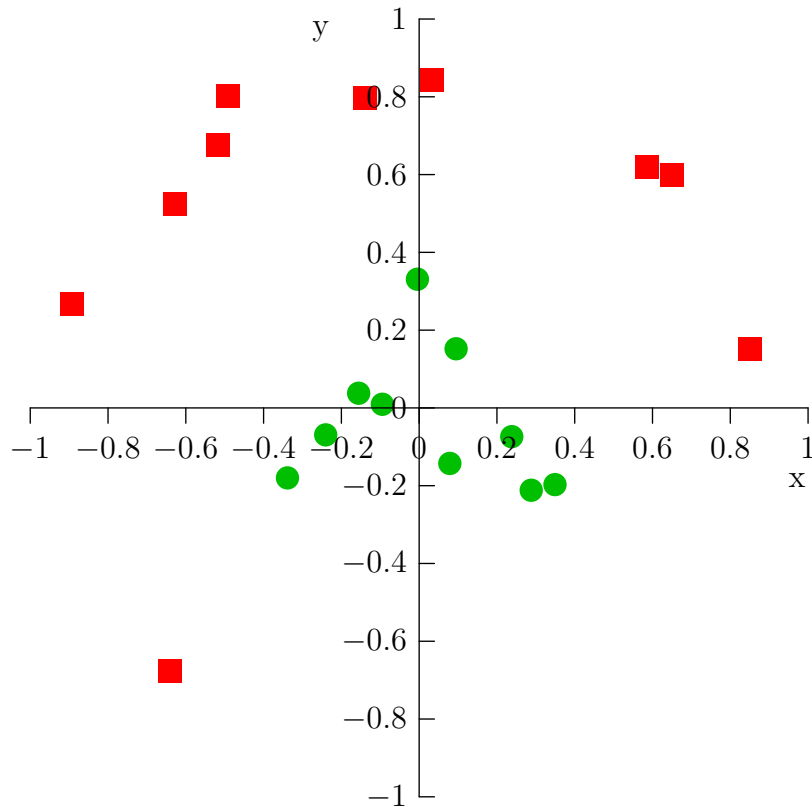


# Non-linearly Separation of Data

$$\begin{aligned} K(\mathbf{x}_1, \mathbf{x}_2) &= \boldsymbol{\phi}^\top(\mathbf{x}_1)\boldsymbol{\phi}(\mathbf{x}_2) = \begin{pmatrix} x_1^2 & y_1^2 & \sqrt{2} x_1 y_1 \end{pmatrix} \begin{pmatrix} x_2^2 \\ y_2^2 \\ \sqrt{2} x_2 y_2 \end{pmatrix} \\ &= x_1^2 x_2^2 + y_1^2 y_2^2 + 2 x_1 y_1 x_2 y_2 = (x_1 x_2 + y_1 y_2)^2 = (\mathbf{x}_1^T \mathbf{x}_2)^2 \end{aligned}$$

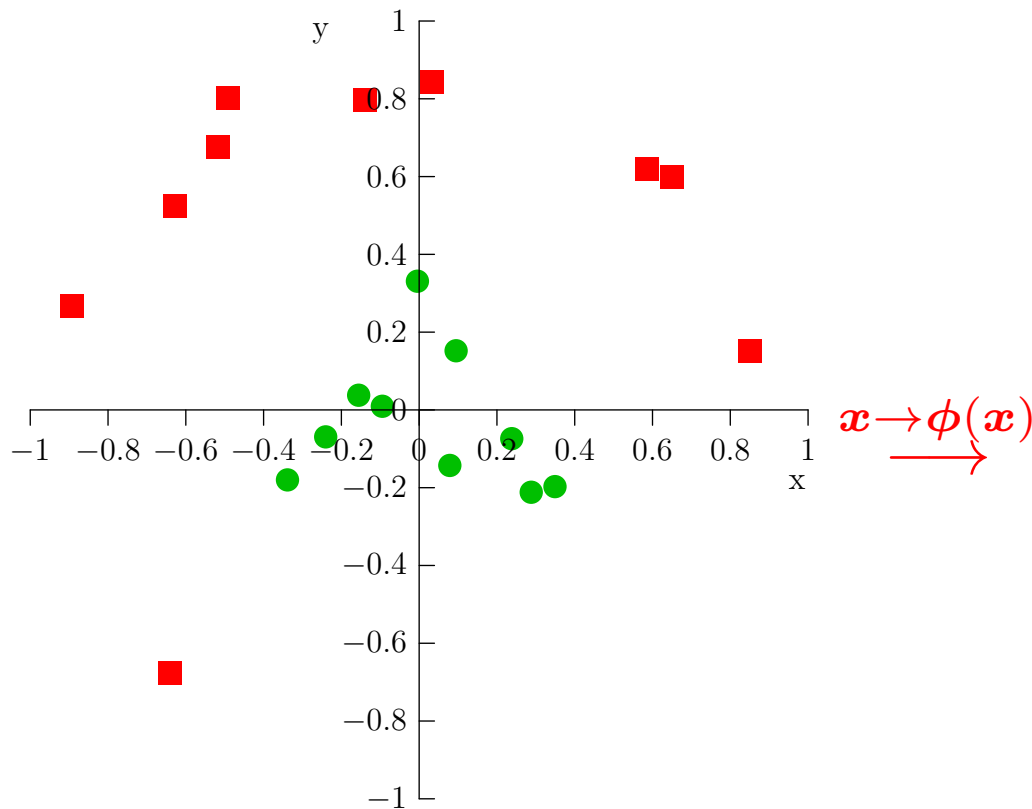
# Non-linearly Separable Data

$$\begin{aligned} K(\mathbf{x}_1, \mathbf{x}_2) &= \phi^\top(\mathbf{x}_1)\phi(\mathbf{x}_2) = \begin{pmatrix} x_1^2 & y_1^2 & \sqrt{2} x_1 y_1 \end{pmatrix} \begin{pmatrix} x_2^2 \\ y_2^2 \\ \sqrt{2} x_2 y_2 \end{pmatrix} \\ &= x_1^2 x_2^2 + y_1^2 y_2^2 + 2 x_1 y_1 x_2 y_2 = (x_1 x_2 + y_1 y_2)^2 = (\mathbf{x}_1^T \mathbf{x}_2)^2 \end{aligned}$$



# Non-linearly Separation of Data

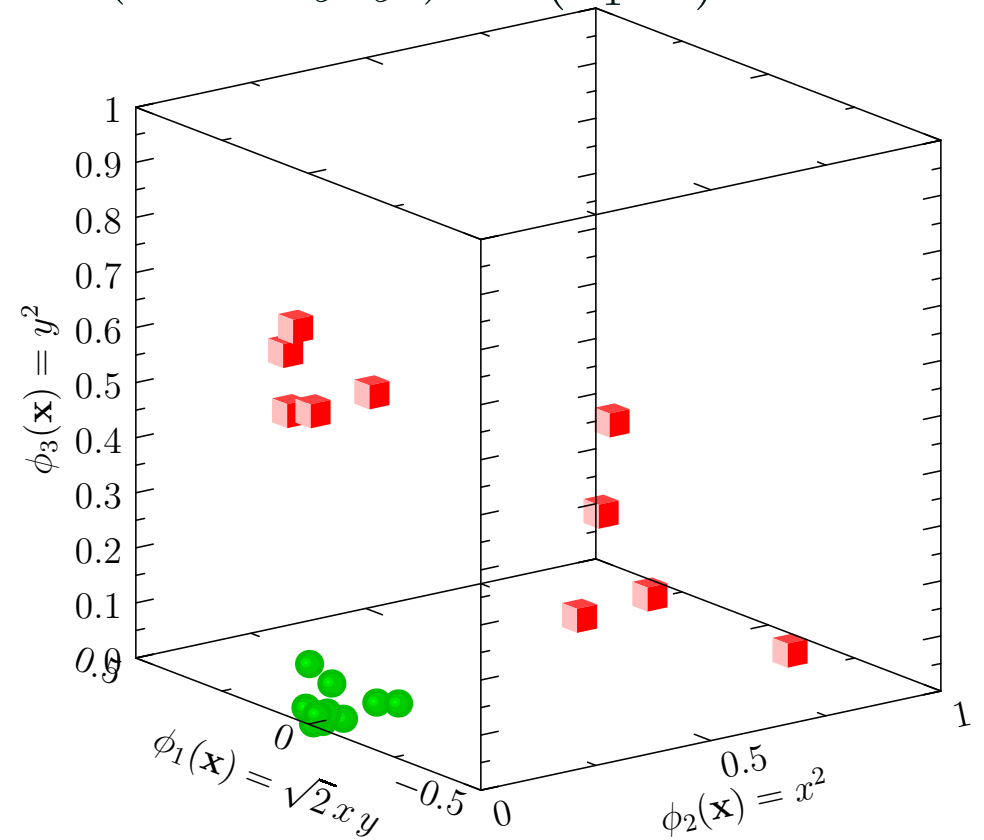
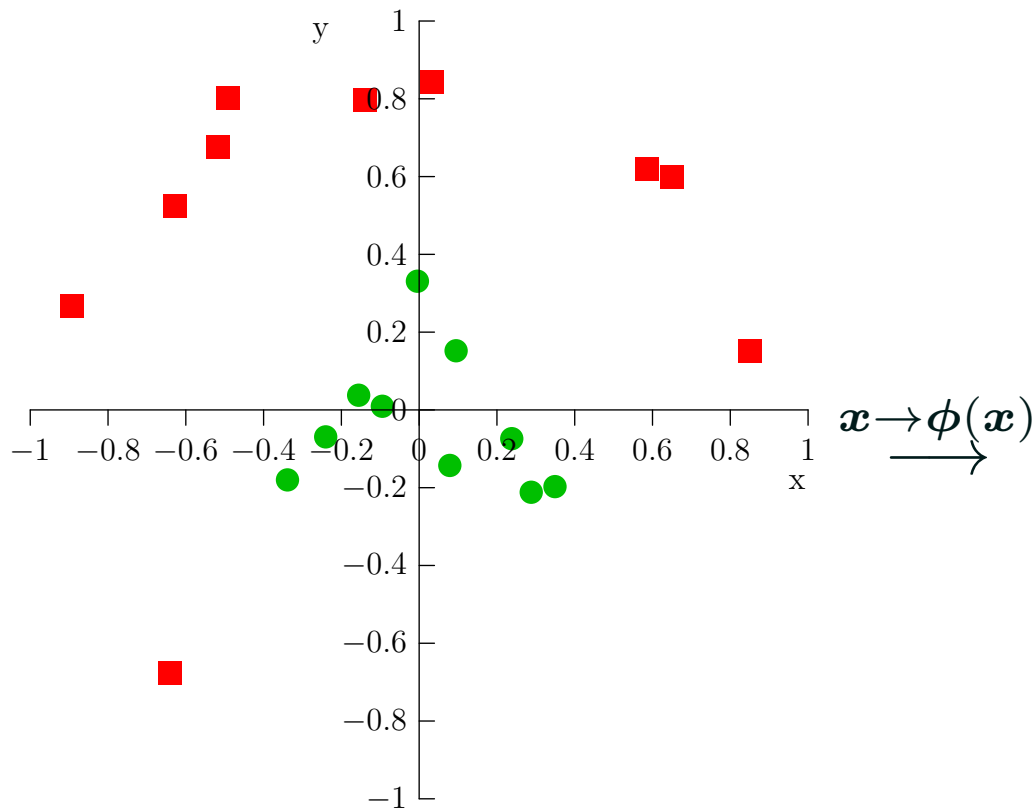
$$\begin{aligned} K(\mathbf{x}_1, \mathbf{x}_2) &= \phi^\top(\mathbf{x}_1)\phi(\mathbf{x}_2) = \begin{pmatrix} x_1^2 & y_1^2 & \sqrt{2} x_1 y_1 \end{pmatrix} \begin{pmatrix} x_2^2 \\ y_2^2 \\ \sqrt{2} x_2 y_2 \end{pmatrix} \\ &= x_1^2 x_2^2 + y_1^2 y_2^2 + 2 x_1 y_1 x_2 y_2 = (x_1 x_2 + y_1 y_2)^2 = (\mathbf{x}_1^T \mathbf{x}_2)^2 \end{aligned}$$



# Non-linearly Separation of Data

$$K(\mathbf{x}_1, \mathbf{x}_2) = \phi^\top(\mathbf{x}_1)\phi(\mathbf{x}_2) = \begin{pmatrix} x_1^2 & y_1^2 & \sqrt{2} x_1 y_1 \end{pmatrix} \begin{pmatrix} x_2^2 \\ y_2^2 \\ \sqrt{2} x_2 y_2 \end{pmatrix}$$

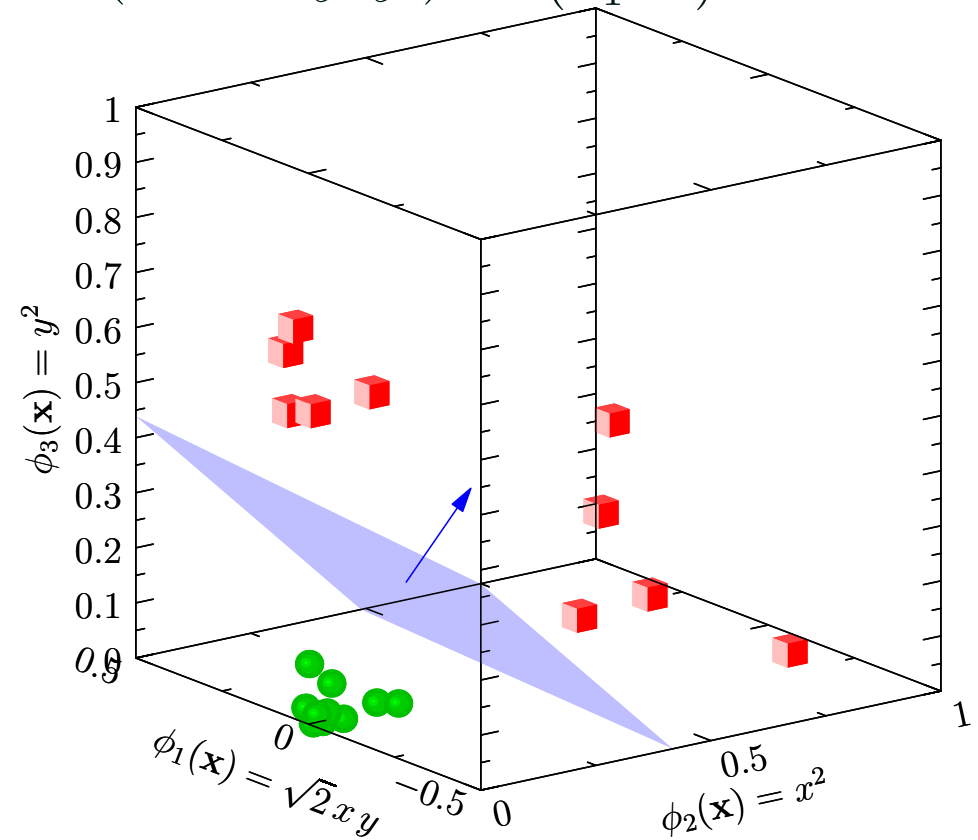
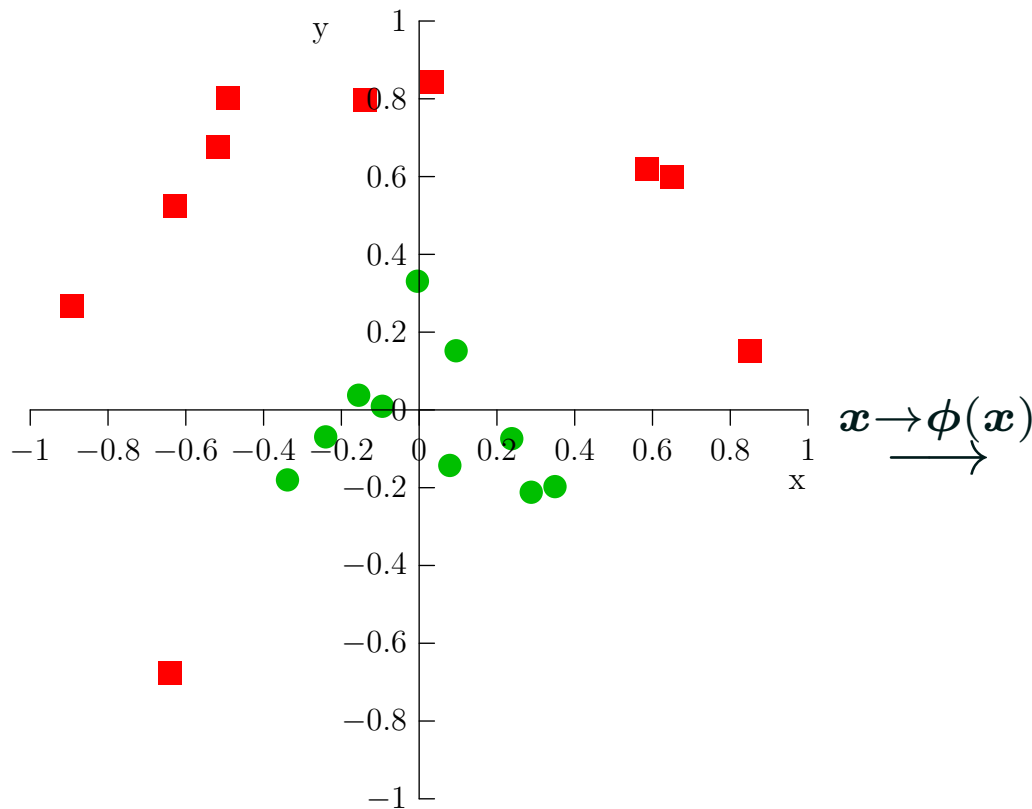
$$= x_1^2 x_2^2 + y_1^2 y_2^2 + 2 x_1 y_1 x_2 y_2 = (x_1 x_2 + y_1 y_2)^2 = (\mathbf{x}_1^T \mathbf{x}_2)^2$$



# Non-linearly Separation of Data

$$K(\mathbf{x}_1, \mathbf{x}_2) = \phi^\top(\mathbf{x}_1)\phi(\mathbf{x}_2) = \begin{pmatrix} x_1^2 & y_1^2 & \sqrt{2} x_1 y_1 \end{pmatrix} \begin{pmatrix} x_2^2 \\ y_2^2 \\ \sqrt{2} x_2 y_2 \end{pmatrix}$$

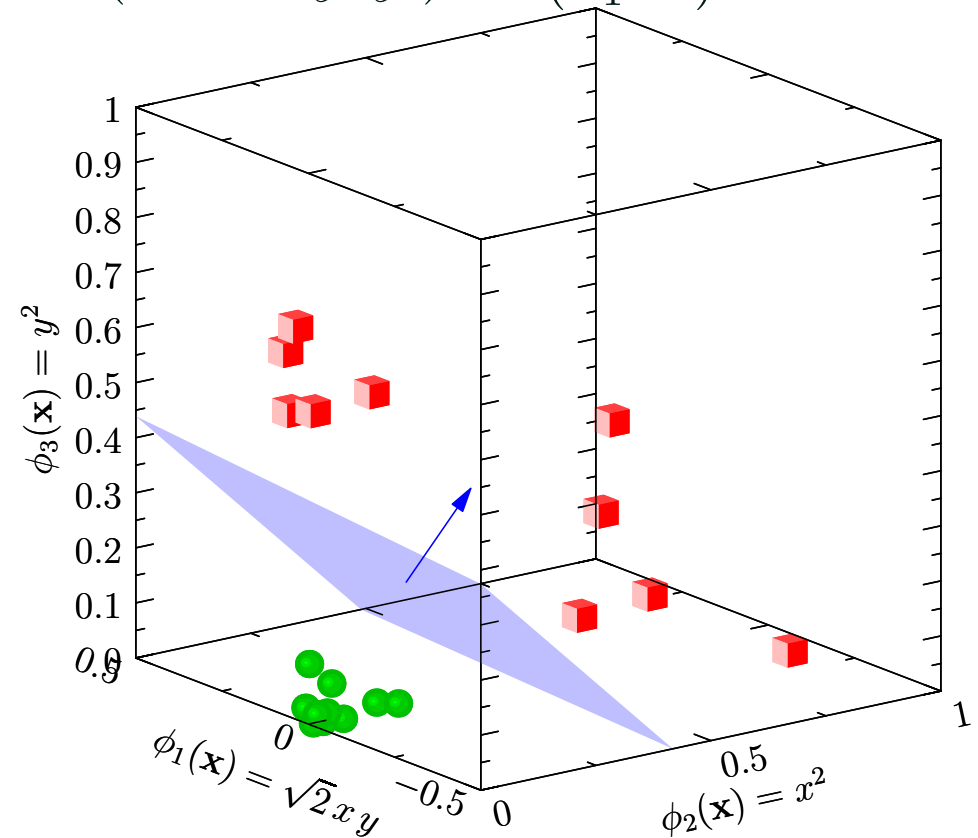
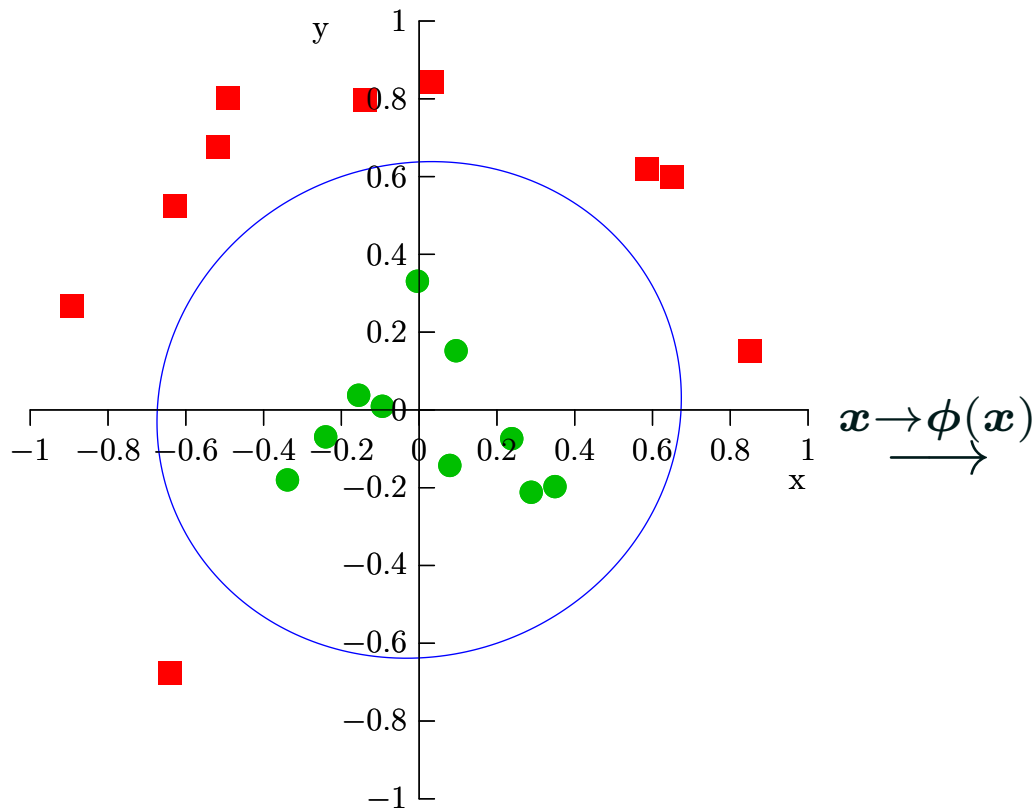
$$= x_1^2 x_2^2 + y_1^2 y_2^2 + 2 x_1 y_1 x_2 y_2 = (x_1 x_2 + y_1 y_2)^2 = (\mathbf{x}_1^T \mathbf{x}_2)^2$$



# Non-linearly Separation of Data

$$K(\mathbf{x}_1, \mathbf{x}_2) = \phi^\top(\mathbf{x}_1)\phi(\mathbf{x}_2) = \begin{pmatrix} x_1^2 & y_1^2 & \sqrt{2}x_1y_1 \end{pmatrix} \begin{pmatrix} x_2^2 \\ y_2^2 \\ \sqrt{2}x_2y_2 \end{pmatrix}$$

$$= x_1^2 x_2^2 + y_1^2 y_2^2 + 2x_1 y_1 x_2 y_2 = (x_1 x_2 + y_1 y_2)^2 = (\mathbf{x}_1^T \mathbf{x}_2)^2$$



# Computing the Maximum-Margin Hyper-plane

- Although it is not hugely difficult to derive the equations for finding the maximum-margin hyper-plane it is not that illuminating (although very elegant)
- Never need to compute  $\phi_i(\mathbf{x})$  only need to compute  $\phi^\top(\mathbf{x}_i)\phi(\mathbf{x}_j)$
- We end up with a quadratic programming problem, which requires some sophisticated algorithms to solve
- When we use the kernel trick the time to compute the solution to the quadratic programming problem is  $pm^3$  where  $m$  is the number of training examples and  $p$  is the number of features

# Computing the Maximum-Margin Hyper-plane

- Although it is not hugely difficult to derive the equations for finding the maximum-margin hyper-plane it is not that illuminating (although very elegant)
- Never need to compute  $\phi_i(\mathbf{x})$  only need to compute  $\phi^\top(\mathbf{x}_i)\phi(\mathbf{x}_j)$
- We end up with a quadratic programming problem, which requires some sophisticated algorithms to solve
- When we use the kernel trick the time to compute the solution to the quadratic programming problem is  $pm^3$  where  $m$  is the number of training examples and  $p$  is the number of features



# Computing the Maximum-Margin Hyper-plane

- Although it is not hugely difficult to derive the equations for finding the maximum-margin hyper-plane it is not that illuminating (although very elegant)
- Never need to compute  $\phi_i(\mathbf{x})$  only need to compute  $\phi^\top(\mathbf{x}_i)\phi(\mathbf{x}_j)$
- We end up with a quadratic programming problem, which requires some sophisticated algorithms to solve
- When we use the kernel trick the time to compute the solution to the quadratic programming problem is  $pm^3$  where  $m$  is the number of training examples and  $p$  is the number of features

# Computing the Maximum-Margin Hyper-plane

- Although it is not hugely difficult to derive the equations for finding the maximum-margin hyper-plane it is not that illuminating (although very elegant)
- Never need to compute  $\phi_i(\mathbf{x})$  only need to compute  $\phi^\top(\mathbf{x}_i)\phi(\mathbf{x}_j)$
- We end up with a quadratic programming problem, which requires some sophisticated algorithms to solve
- When we use the kernel trick the time to compute the solution to the quadratic programming problem is  $pm^3$  where  $m$  is the number of training examples and  $p$  is the number of features

# Getting SVMs to Work Well

- SVMs rely on distances between data points
- These will change relative to each other if we rescale some features but not other—giving different maximum-margin hyper-planes
- If we don't know what features are important (most often the case), then it is worth scaling each feature (for example, so their range is between 0 and 1 or their variance is 1)

# Getting SVMs to Work Well

- SVMs rely on distances between data points
- These will change relative to each other if we rescale some features but not other—giving different maximum-margin hyper-planes
- If we don't know what features are important (most often the case), then it is worth scaling each feature (for example, so their range is between 0 and 1 or their variance is 1)

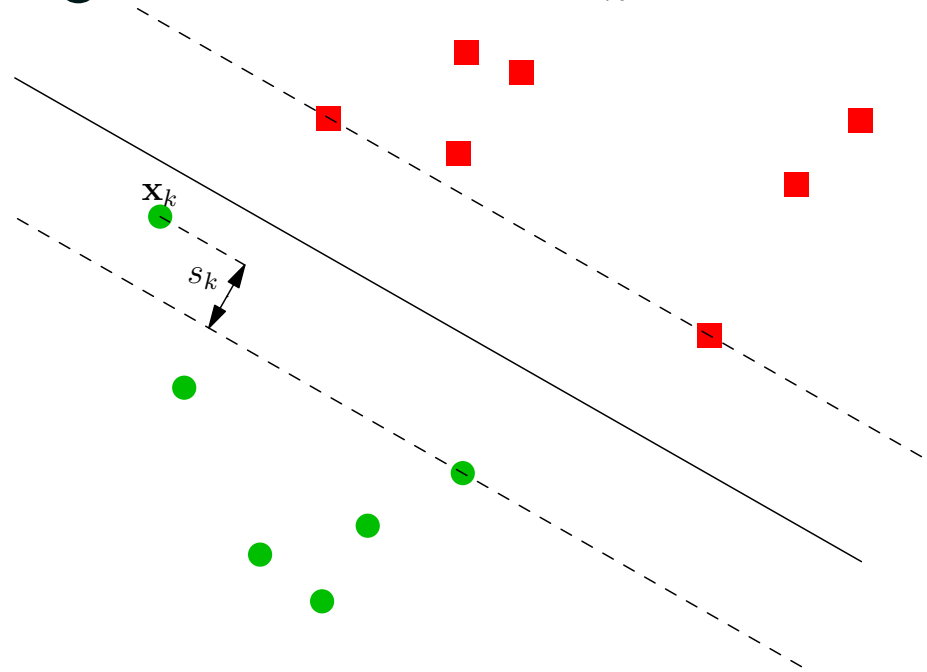
# Getting SVMs to Work Well

- SVMs rely on distances between data points
- These will change relative to each other if we rescale some features but not other—giving different maximum-margin hyper-planes
- If we don't know what features are important (most often the case), then it is worth scaling each feature (for example, so their range is between 0 and 1 or their variance is 1)

# Soft Margins

- Sometimes the margin constraint is too severe
- Relax constraints by introducing *slack variables*,  $s_k \geq 0$

$$y_k(\mathbf{x}_k^\top \mathbf{w} - b) \geq 1 - s_k$$

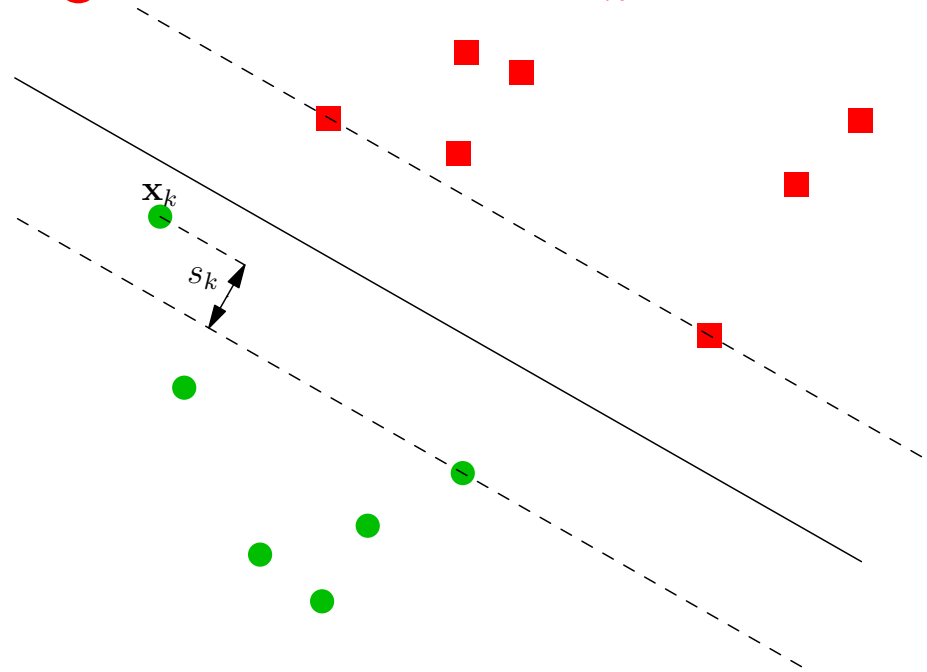


- Minimise  $\frac{\|\mathbf{w}\|^2}{2} + C \sum_{k=1}^n s_k$  subject to constraints
- Large  $C$  punishes slack variables

# Soft Margins

- Sometimes the margin constraint is too severe
- Relax constraints by introducing *slack variables*,  $s_k \geq 0$

$$y_k(\mathbf{x}_k^T \mathbf{w} - b) \geq 1 - s_k$$

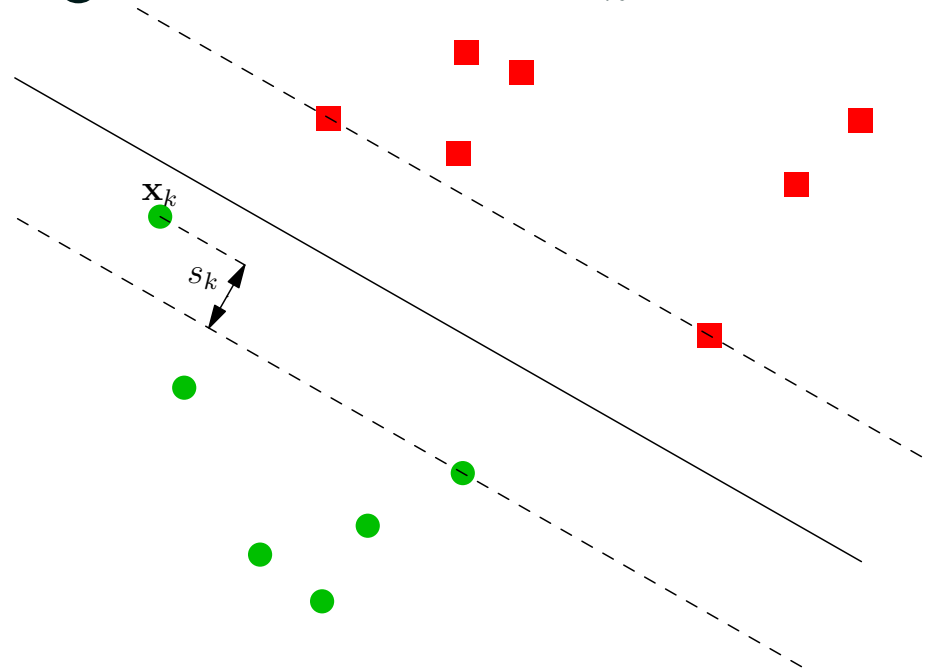


- Minimise  $\frac{\|\mathbf{w}\|^2}{2} + C \sum_{k=1}^n s_k$  subject to constraints
- Large  $C$  punishes slack variables

# Soft Margins

- Sometimes the margin constraint is too severe
- Relax constraints by introducing *slack variables*,  $s_k \geq 0$

$$y_k(\mathbf{x}_k^\top \mathbf{w} - b) \geq 1 - s_k$$



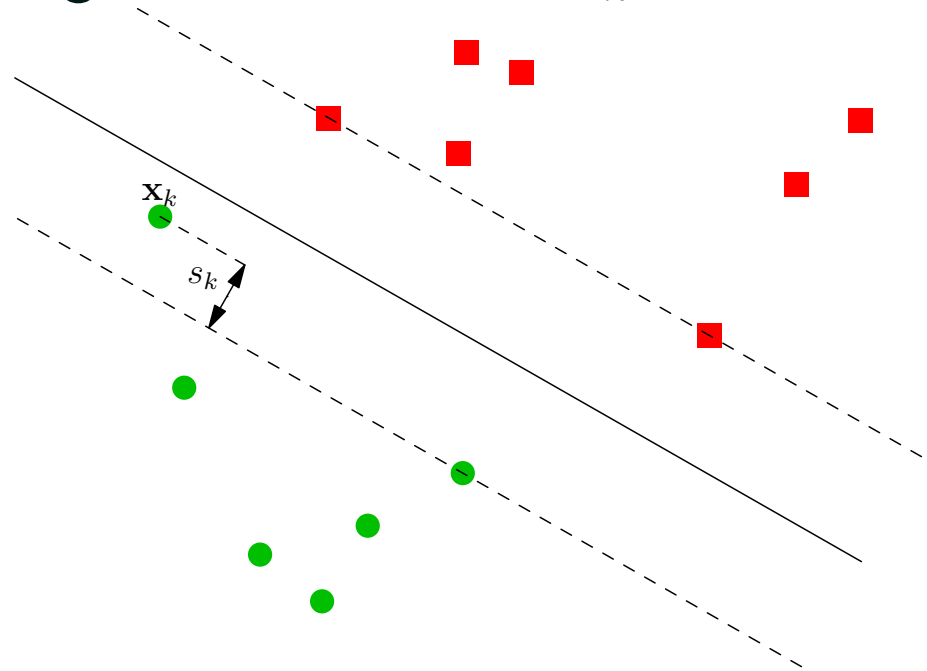
- Minimise  $\frac{\|\mathbf{w}\|^2}{2} + C \sum_{k=1}^n s_k$  subject to constraints
- Large  $C$  punishes slack variables



# Soft Margins

- Sometimes the margin constraint is too severe
- Relax constraints by introducing *slack variables*,  $s_k \geq 0$

$$y_k(\mathbf{x}_k^\top \mathbf{w} - b) \geq 1 - s_k$$



- Minimise  $\frac{\|\mathbf{w}\|^2}{2} + C \sum_{k=1}^n s_k$  subject to constraints
- Large  $C$  punishes slack variables

# Optimising $C$

- In practice it can make a huge difference to the performance if we change  $C$
- Optimal  $C$  values changes by many orders of magnitude e.g.  $2^{-5}$ – $2^{15}$
- Typically optimised by a grid search (start from  $2^{-5}$  say and double until you reach  $2^{15}$ )

# Optimising C

- In practice it can make a huge difference to the performance if we change  $C$
- Optimal  $C$  values changes by many orders of magnitude e.g.  $2^{-5}$ – $2^{15}$
- Typically optimised by a grid search (start from  $2^{-5}$  say and double until you reach  $2^{15}$ )

# Optimising $C$

- In practice it can make a huge difference to the performance if we change  $C$
- Optimal  $C$  values changes by many orders of magnitude e.g.  $2^{-5}$ – $2^{15}$
- Typically optimised by a grid search (start from  $2^{-5}$  say and double until you reach  $2^{15}$ )

# Choosing the Right Kernel Function

- There are kernels designed for particular data types (e.g. string kernels for text or biological sequences)
- For numerical data people tend to look at using no kernel (linear SVM), a radial basis function (Gaussian) kernel or polynomial kernels
- Kernel's often come with parameters, e.g. the popular radial basis function kernel

$$K(\mathbf{x}, \mathbf{y}) = e^{-\gamma \|\mathbf{x} - \mathbf{y}\|^2}$$

- Optimal  $\gamma$  values range over  $2^{-15}$ – $2^3$

# Choosing the Right Kernel Function

- There are kernels designed for particular data types (e.g. string kernels for text or biological sequences)
- For numerical data people tend to look at using no kernel (linear SVM), a radial basis function (Gaussian) kernel or polynomial kernels
- Kernel's often come with parameters, e.g. the popular radial basis function kernel

$$K(\mathbf{x}, \mathbf{y}) = e^{-\gamma \|\mathbf{x} - \mathbf{y}\|^2}$$

- Optimal  $\gamma$  values range over  $2^{-15}$ – $2^3$

# Choosing the Right Kernel Function

- There are kernels designed for particular data types (e.g. string kernels for text or biological sequences)
- For numerical data people tend to look at using no kernel (linear SVM), a radial basis function (Gaussian) kernel or polynomial kernels
- Kernel's often come with parameters, e.g. the popular radial basis function kernel

$$K(\mathbf{x}, \mathbf{y}) = e^{-\gamma \|\mathbf{x} - \mathbf{y}\|^2}$$

- Optimal  $\gamma$  values range over  $2^{-15}$ – $2^3$

# Choosing the Right Kernel Function

- There are kernels designed for particular data types (e.g. string kernels for text or biological sequences)
- For numerical data people tend to look at using no kernel (linear SVM), a radial basis function (Gaussian) kernel or polynomial kernels
- Kernel's often come with parameters, e.g. the popular radial basis function kernel

$$K(\mathbf{x}, \mathbf{y}) = e^{-\gamma \|\mathbf{x} - \mathbf{y}\|^2}$$

- Optimal  $\gamma$  values range over  $2^{-15}$ – $2^3$



# SVM Libraries

- Although SVMs have unique solutions, they require very well written optimisers
- If you have a large data set they can be very slow
- There are good libraries out there, svmlib, svm-lite, etc.
- These will often automate normalisation of data and grid search for parameters

# SVM Libraries

- Although SVMs have unique solutions, they require very well written optimisers
- If you have a large data set they can be very slow
- There are good libraries out there, svmlib, svm-lite, etc.
- These will often automate normalisation of data and grid search for parameters

# SVM Libraries

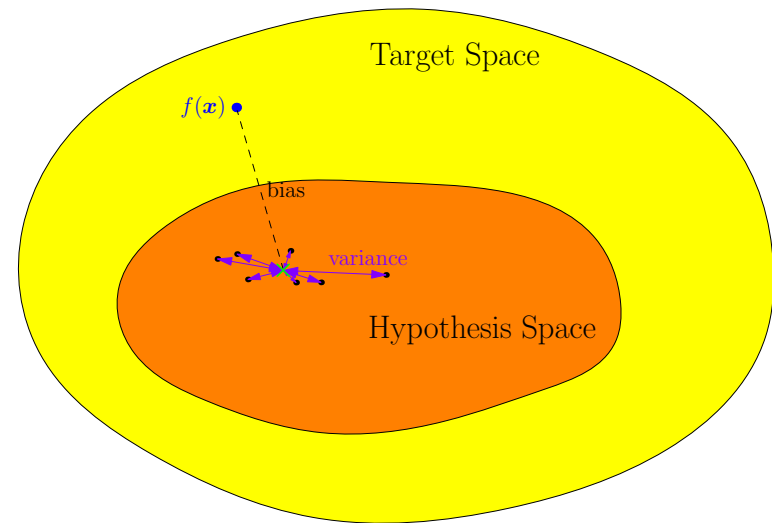
- Although SVMs have unique solutions, they require very well written optimisers
- If you have a large data set they can be very slow
- There are good libraries out there, svmlib, svm-lite, etc.
- These will often automate normalisation of data and grid search for parameters

# SVM Libraries

- Although SVMs have unique solutions, they require very well written optimisers
- If you have a large data set they can be very slow
- There are good libraries out there, svmlib, svm-lite, etc.
- These will often automate normalisation of data and grid search for parameters

# Outline

1. What Makes a Good Learning Machine?
2. SVMs
3. **Ensemble Methods**
4. Bayesian Inference



# Removing Variance By Averaging

- We can reduce the variance and hence improve our generalisation error by averaging over different learning machines
- There are a number of different techniques for doing this that go by the name of **ensemble methods** or **ensemble learning**
- This trick can be used with many different learning machines, but is clearly most practical for machine that can be trained quickly (nevertheless, even for deep learning taking the average response of many machines is usually done to win competitions)

# Removing Variance By Averaging

- We can reduce the variance and hence improve our generalisation error by averaging over different learning machines
- There are a number of different techniques for doing this that go by the name of **ensemble methods** or **ensemble learning**
- This trick can be used with many different learning machines, but is clearly most practical for machine that can be trained quickly (nevertheless, even for deep learning taking the average response of many machines is usually done to win competitions)

# Removing Variance By Averaging

- We can reduce the variance and hence improve our generalisation error by averaging over different learning machines
- There are a number of different techniques for doing this that go by the name of **ensemble methods** or **ensemble learning**
- This trick can be used with many different learning machines, but is clearly most practical for machine that can be trained quickly (nevertheless, even for deep learning taking the average response of many machines is usually done to win competitions)



# Ensembling of Decision Trees

- One set of algorithms where ensembling are common place are decision trees
- These are particularly good for handling messy data
  - ★ categorical data
  - ★ mixture of data types
  - ★ missing data
  - ★ large data sets
  - ★ multiclass
- In many competitions ensembled trees, particularly *random forests* and *gradient boosting* beats all other techniques

# Ensembling of Decision Trees

- One set of algorithms where ensembling are common place are decision trees
- These are particularly good for handling messy data
  - ★ categorical data
  - ★ mixture of data types
  - ★ missing data
  - ★ large data sets
  - ★ multiclass
- In many competitions ensembled trees, particularly *random forests* and *gradient boosting* beats all other techniques

# Ensembling of Decision Trees

- One set of algorithms where ensembling are common place are decision trees
- These are particularly good for handling messy data
  - ★ categorical data
  - ★ mixture of data types
  - ★ missing data
  - ★ large data sets
  - ★ multiclass
- In many competitions ensembled trees, particularly *random forests* and *gradient boosting* beats all other techniques

# Ensembling of Decision Trees

- One set of algorithms where ensembling are common place are decision trees
- These are particularly good for handling messy data
  - ★ categorical data
  - ★ mixture of data types
  - ★ missing data
  - ★ large data sets
  - ★ multiclass
- In many competitions ensembled trees, particularly *random forests* and *gradient boosting* beats all other techniques

# Ensembling of Decision Trees

- One set of algorithms where ensembling are common place are decision trees
- These are particularly good for handling messy data
  - ★ categorical data
  - ★ mixture of data types
  - ★ missing data
  - ★ large data sets
  - ★ multiclass
- In many competitions ensembled trees, particularly *random forests* and *gradient boosting* beats all other techniques

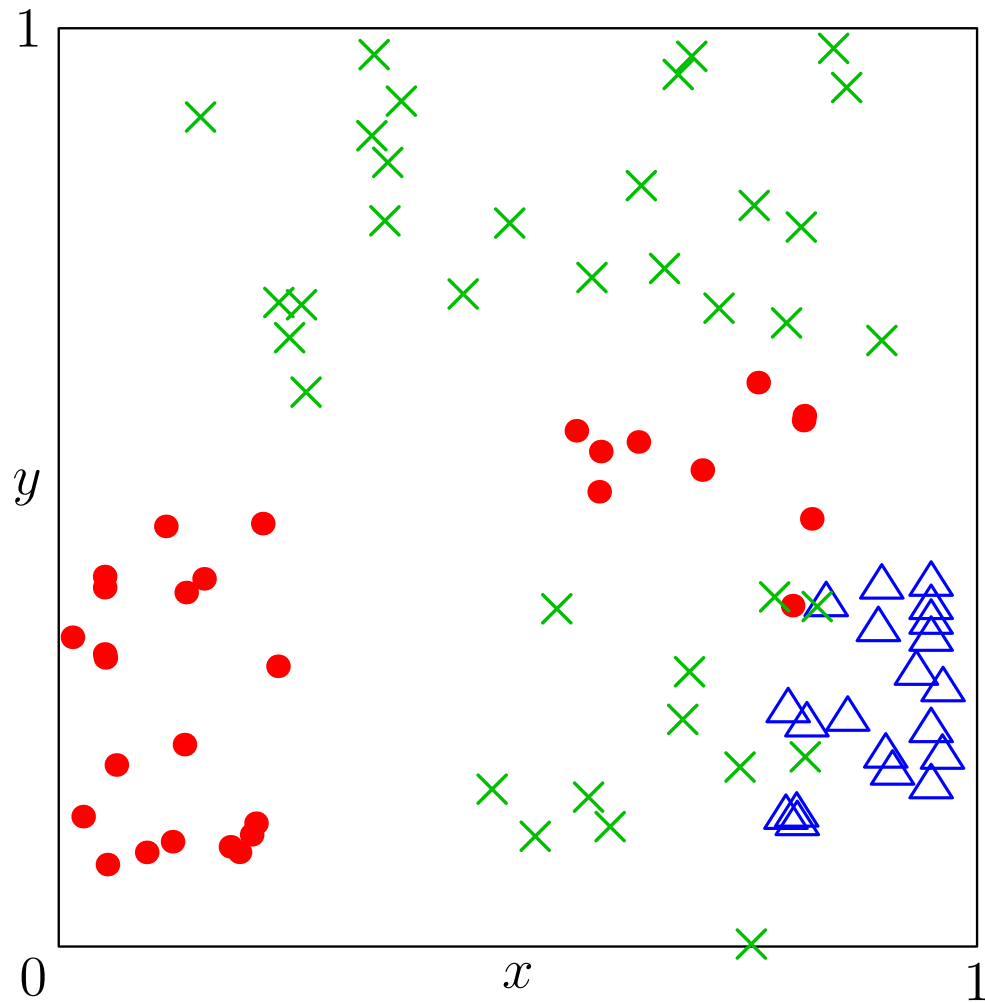
# Ensembling of Decision Trees

- One set of algorithms where ensembling are common place are decision trees
- These are particularly good for handling messy data
  - ★ categorical data
  - ★ mixture of data types
  - ★ missing data
  - ★ large data sets
  - ★ **multiclass**
- In many competitions ensembled trees, particularly *random forests* and *gradient boosting* beats all other techniques

# Ensembling of Decision Trees

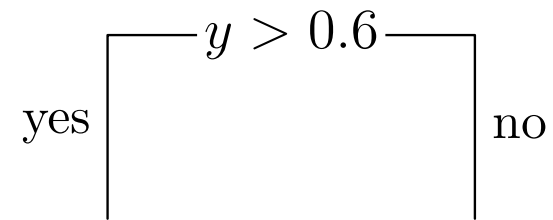
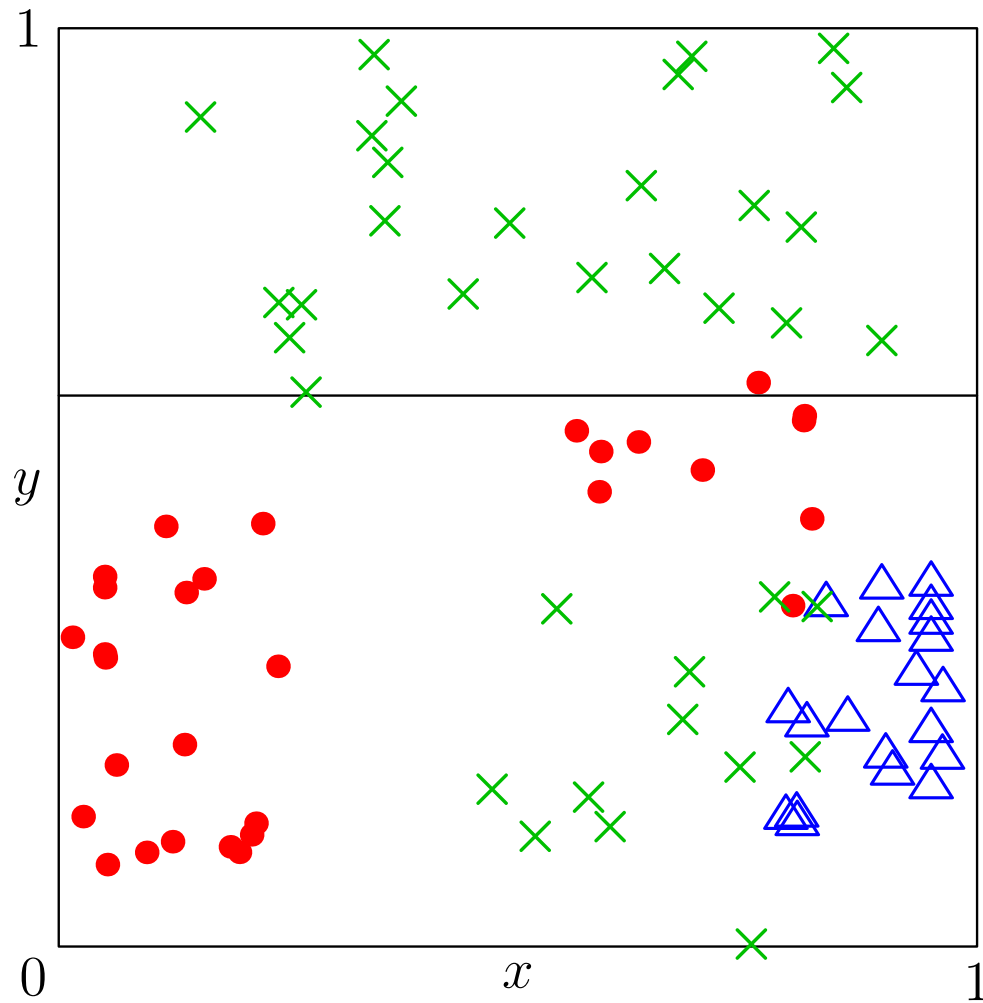
- One set of algorithms where ensembling are common place are decision trees
- These are particularly good for handling messy data
  - ★ categorical data
  - ★ mixture of data types
  - ★ missing data
  - ★ large data sets
  - ★ multiclass
- In many competitions ensembled trees, particularly *random forests* and *gradient boosting* beats all other techniques

# Building Decision Trees

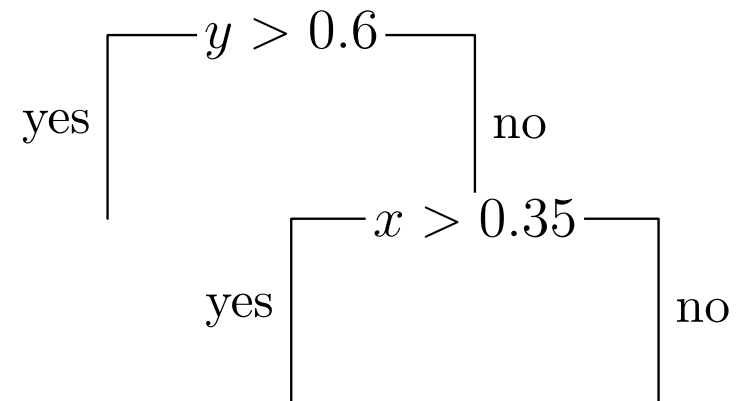
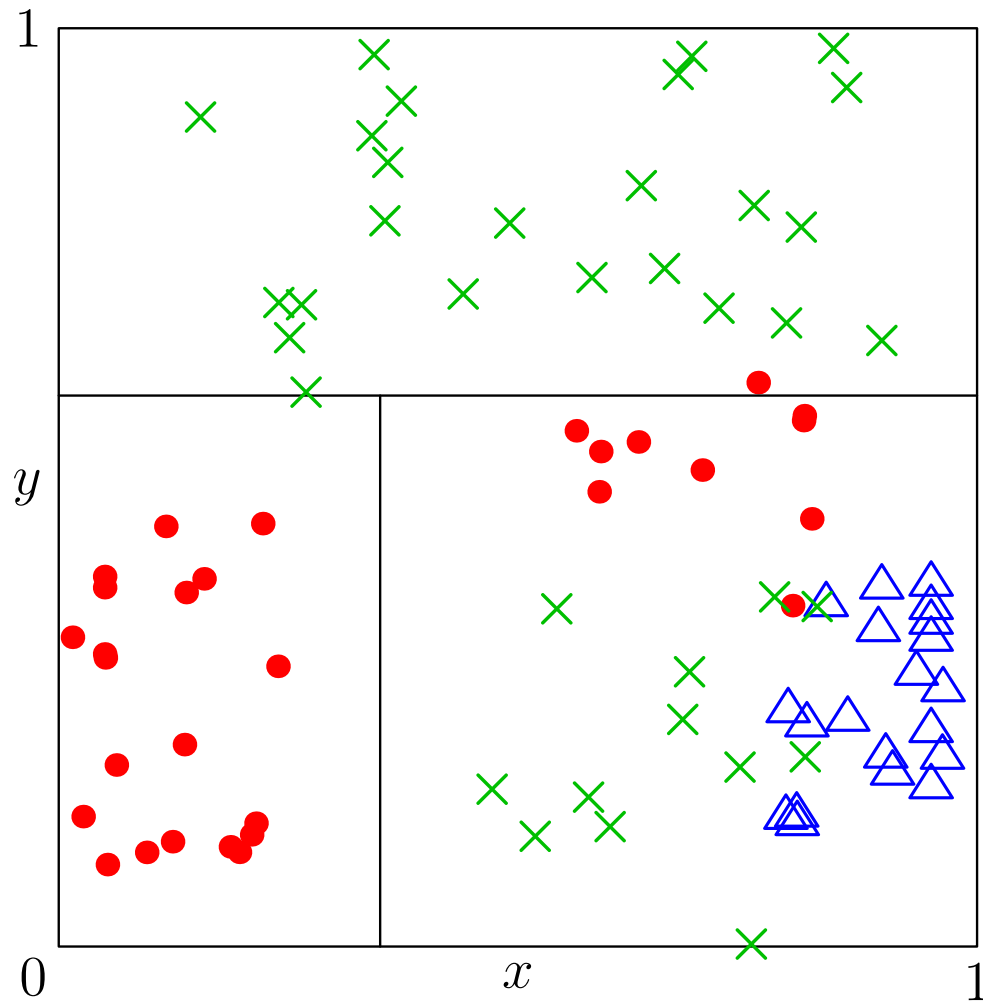




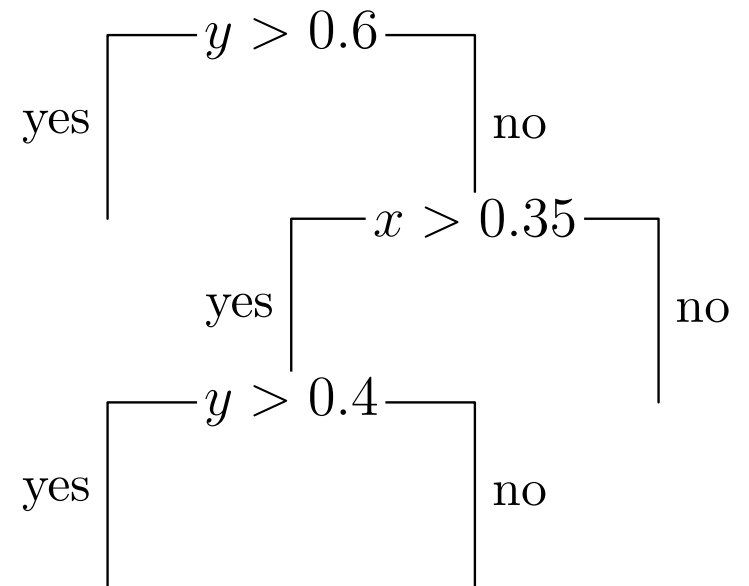
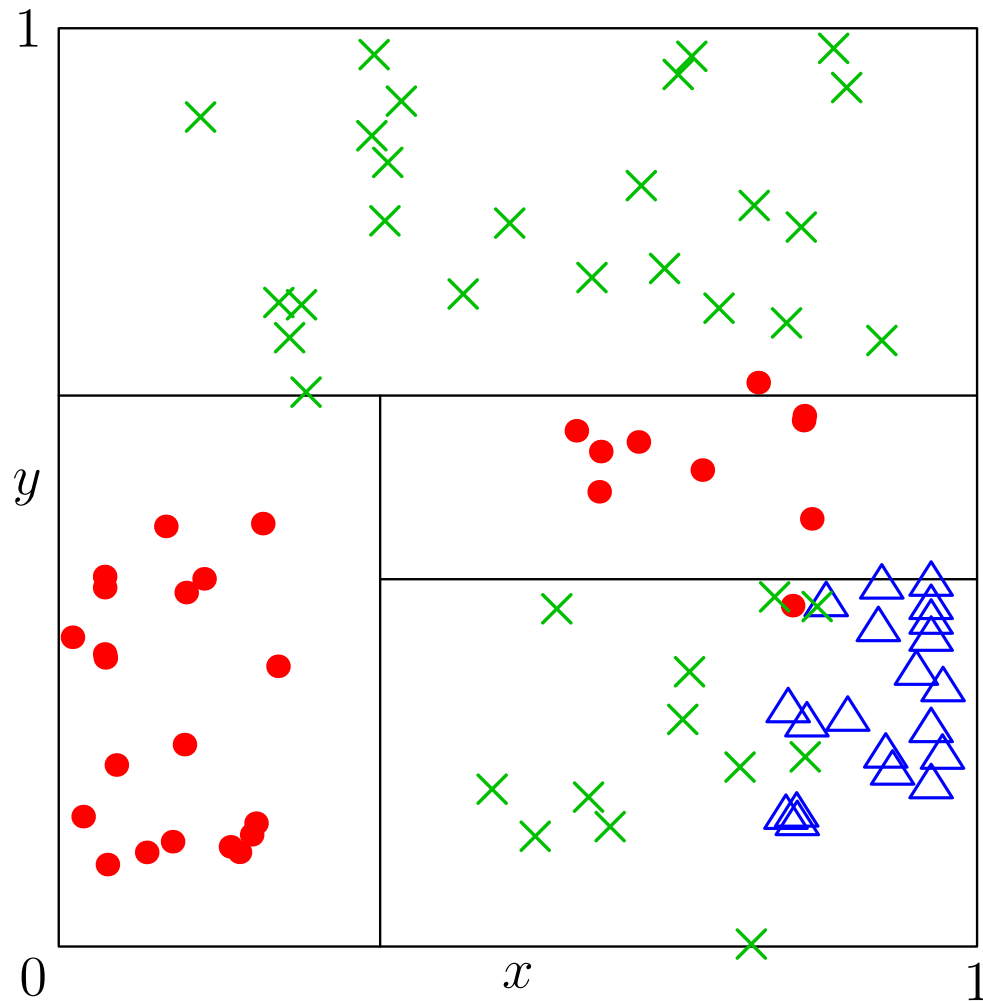
# Building Decision Trees



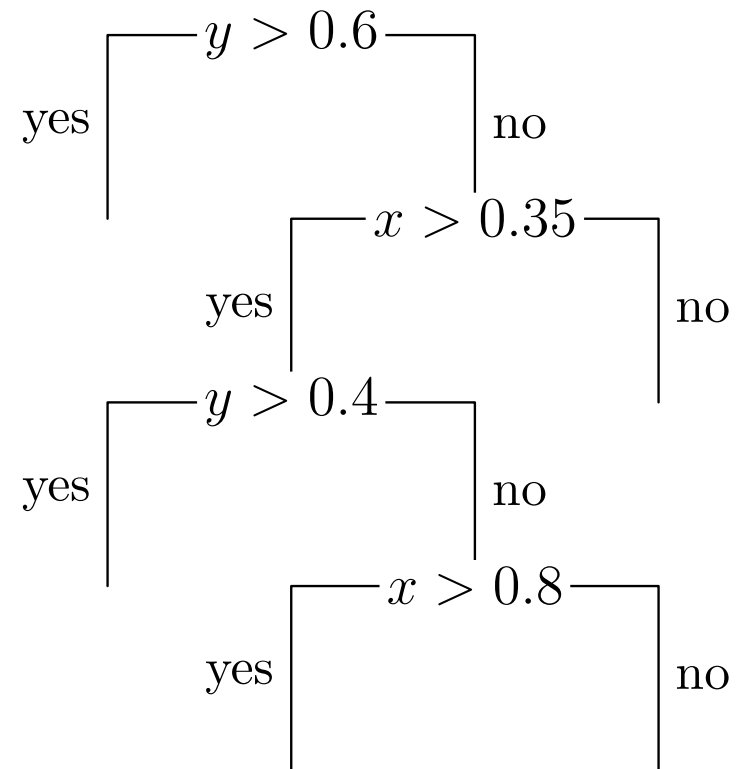
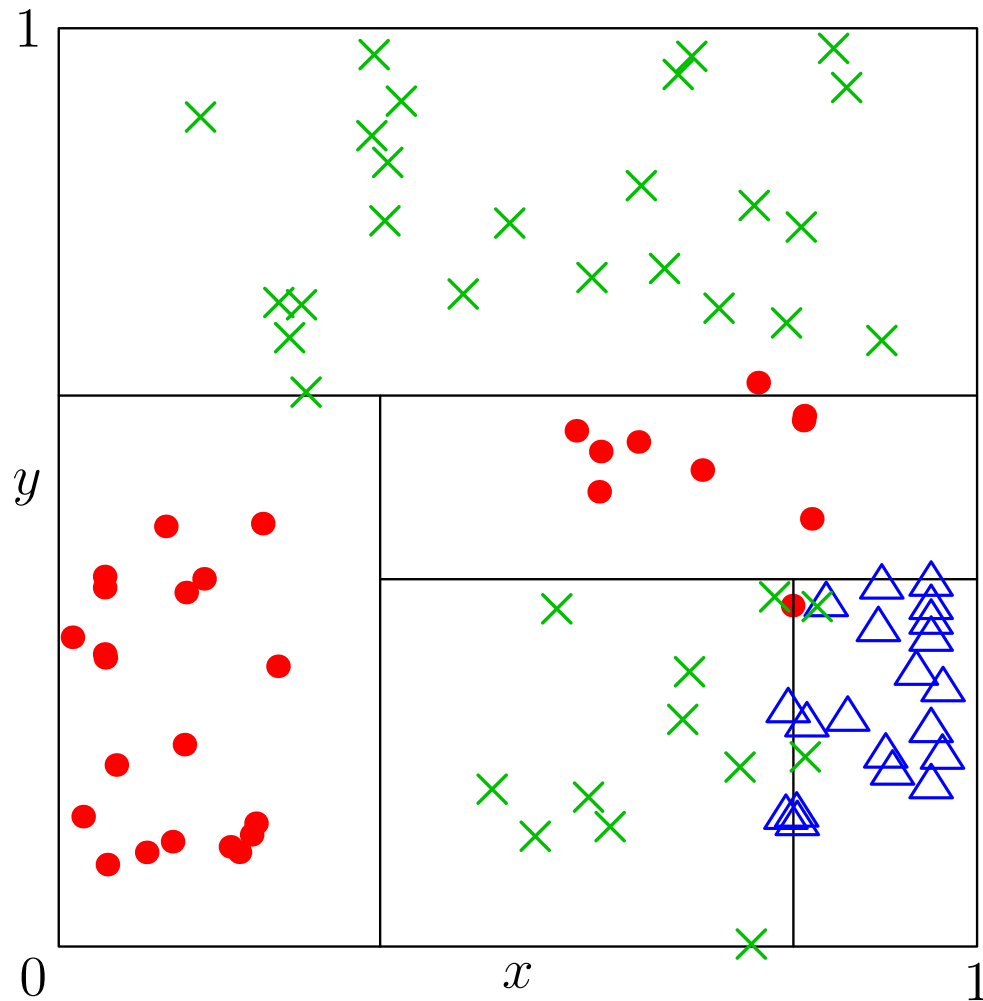
# Building Decision Trees



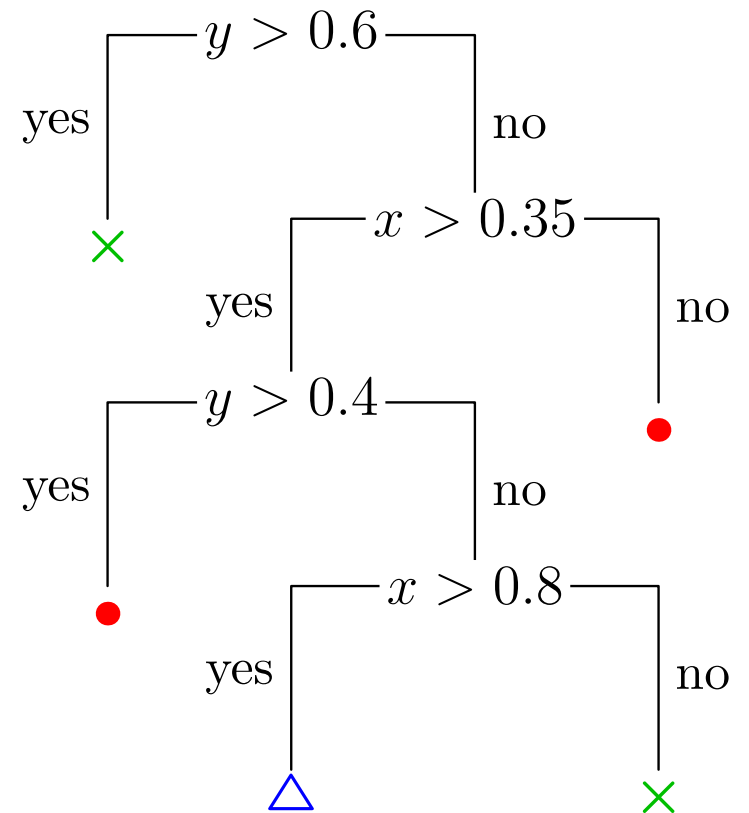
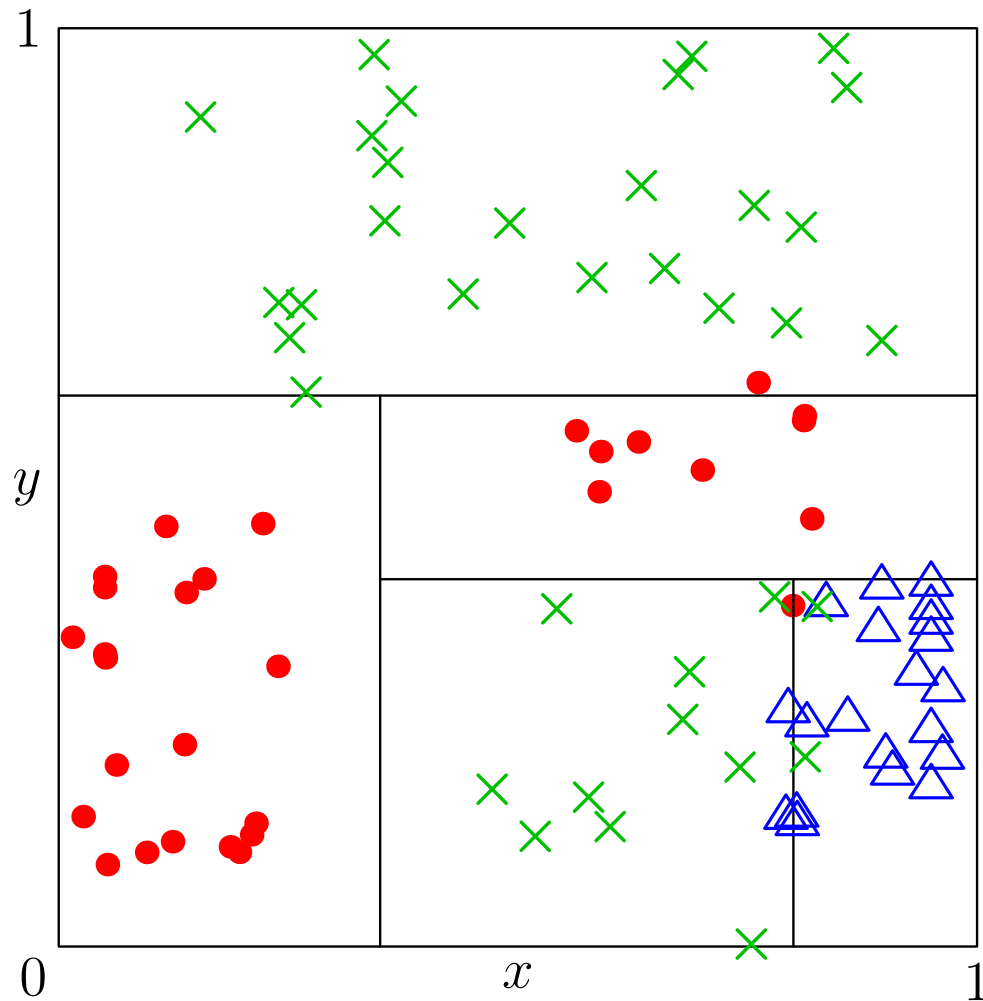
# Building Decision Trees



# Building Decision Trees



# Building Decision Trees



# Using Decision Trees

- Decision trees are very dependent on the first few splits of data
- As a result they have high variance and usually not great generalisation performance
- Nevertheless they can be very useful for quickly exploring data
- They are readily interpretable and often reveal what are important parameters
- However, when ensembled they become state-of-the-art

# Using Decision Trees

- Decision trees are very dependent on the first few splits of data
- As a result they have high variance and usually not great generalisation performance
- Nevertheless they can be very useful for quickly exploring data
- They are readily interpretable and often reveal what are important parameters
- However, when ensembled they become state-of-the-art

# Using Decision Trees

- Decision trees are very dependent on the first few splits of data
- As a result they have high variance and usually not great generalisation performance
- Nevertheless they can be very useful for quickly exploring data
- They are readily interpretable and often reveal what are important parameters
- However, when ensembled they become state-of-the-art



# Using Decision Trees

- Decision trees are very dependent on the first few splits of data
- As a result they have high variance and usually not great generalisation performance
- Nevertheless they can be very useful for quickly exploring data
- They are readily interpretable and often reveal what are important parameters
- However, when ensembled they become state-of-the-art

# Using Decision Trees

- Decision trees are very dependent on the first few splits of data
- As a result they have high variance and usually not great generalisation performance
- Nevertheless they can be very useful for quickly exploring data
- They are readily interpretable and often reveal what are important parameters
- However, when ensembled they become state-of-the-art

# Bootstrap Aggregation (Bagging)

- In order to average machine they must learn something different
- We only have one data set, but we can resample from the data set to make them look a bit different—this is known as **bootstrapping**

# Bootstrap Aggregation (Bagging)

- In order to average machine they must learn something different
- We only have one data set, but we can resample from the data set to make them look a bit different—this is known as **bootstrapping**

# Bootstrap Aggregation (Bagging)

- In order to average machine they must learn something different
- We only have one data set, but we can resample from the data set to make them look a bit different—this is known a **bootstrapping**

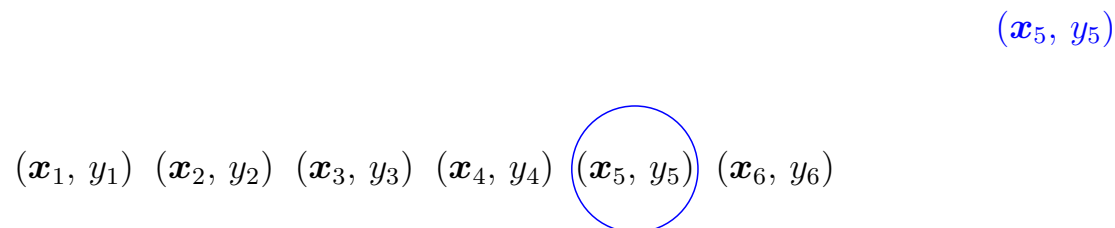
# Bootstrap Aggregation (Bagging)

- In order to average machine they must learn something different
- We only have one data set, but we can resample from the data set to make them look a bit different—this is known as **bootstrapping**

$(\mathbf{x}_1, y_1) (\mathbf{x}_2, y_2) (\mathbf{x}_3, y_3) (\mathbf{x}_4, y_4) (\mathbf{x}_5, y_5) (\mathbf{x}_6, y_6)$

# Bootstrap Aggregation (Bagging)

- In order to average machine they must learn something different
- We only have one data set, but we can resample from the data set to make them look a bit different—this is known as **bootstrapping**



# Bootstrap Aggregation (Bagging)

- In order to average machine they must learn something different
- We only have one data set, but we can resample from the data set to make them look a bit different—this is known as **bootstrapping**

$(\mathbf{x}_5, y_5)$   $(\mathbf{x}_4, y_4)$

$(\mathbf{x}_1, y_1)$   $(\mathbf{x}_2, y_2)$   $(\mathbf{x}_3, y_3)$   $(\mathbf{x}_4, y_4)$   $(\mathbf{x}_5, y_5)$   $(\mathbf{x}_6, y_6)$



# Bootstrap Aggregation (Bagging)

- In order to average machine they must learn something different
- We only have one data set, but we can resample from the data set to make them look a bit different—this is known as **bootstrapping**

$(\mathbf{x}_5, y_5)$   $(\mathbf{x}_4, y_4)$   $(\mathbf{x}_5, y_5)$

$(\mathbf{x}_1, y_1)$   $(\mathbf{x}_2, y_2)$   $(\mathbf{x}_3, y_3)$   $(\mathbf{x}_4, y_4)$   $(\mathbf{x}_5, y_5)$   $(\mathbf{x}_6, y_6)$

# Bootstrap Aggregation (Bagging)

- In order to average machine they must learn something different
- We only have one data set, but we can resample from the data set to make them look a bit different—this is known as **bootstrapping**

$(\mathbf{x}_5, y_5)$   $(\mathbf{x}_4, y_4)$   $(\mathbf{x}_5, y_5)$   $(\mathbf{x}_6, y_6)$

$(\mathbf{x}_1, y_1)$   $(\mathbf{x}_2, y_2)$   $(\mathbf{x}_3, y_3)$   $(\mathbf{x}_4, y_4)$   $(\mathbf{x}_5, y_5)$   $(\mathbf{x}_6, y_6)$

# Bootstrap Aggregation (Bagging)

- In order to average machine they must learn something different
- We only have one data set, but we can resample from the data set to make them look a bit different—this is known as **bootstrapping**

$(\mathbf{x}_5, y_5)$   $(\mathbf{x}_4, y_4)$   $(\mathbf{x}_5, y_5)$   $(\mathbf{x}_6, y_6)$   $(\mathbf{x}_2, y_2)$

$(\mathbf{x}_1, y_1)$   $(\mathbf{x}_2, y_2)$   $(\mathbf{x}_3, y_3)$   $(\mathbf{x}_4, y_4)$   $(\mathbf{x}_5, y_5)$   $(\mathbf{x}_6, y_6)$

# Bootstrap Aggregation (Bagging)

- In order to average machine they must learn something different
- We only have one data set, but we can resample from the data set to make them look a bit different—this is known as **bootstrapping**

$(\mathbf{x}_5, y_5)$   $(\mathbf{x}_4, y_4)$   $(\mathbf{x}_5, y_5)$   $(\mathbf{x}_6, y_6)$   $(\mathbf{x}_2, y_2)$   $(\mathbf{x}_1, y_1)$

$(\mathbf{x}_1, y_1)$   $(\mathbf{x}_2, y_2)$   $(\mathbf{x}_3, y_3)$   $(\mathbf{x}_4, y_4)$   $(\mathbf{x}_5, y_5)$   $(\mathbf{x}_6, y_6)$

# Bootstrap Aggregation (Bagging)

- In order to average machine they must learn something different
- We only have one data set, but we can resample from the data set to make them look a bit different—this is known as **bootstrapping**

$(\mathbf{x}_5, y_5)$   $(\mathbf{x}_4, y_4)$   $(\mathbf{x}_5, y_5)$   $(\mathbf{x}_6, y_6)$   $(\mathbf{x}_2, y_2)$   $(\mathbf{x}_1, y_1)$   
 $(\mathbf{x}_1, y_1)$   $(\mathbf{x}_5, y_5)$   $(\mathbf{x}_3, y_3)$   $(\mathbf{x}_2, y_2)$   $(\mathbf{x}_5, y_5)$   $(\mathbf{x}_2, y_2)$

$(\mathbf{x}_1, y_1)$   $(\mathbf{x}_2, y_2)$   $(\mathbf{x}_3, y_3)$   $(\mathbf{x}_4, y_4)$   $(\mathbf{x}_5, y_5)$   $(\mathbf{x}_6, y_6)$

# Bootstrap Aggregation (Bagging)

- In order to average machine they must learn something different
- We only have one data set, but we can resample from the data set to make them look a bit different—this is known as **bootstrapping**

$(\mathbf{x}_1, y_1)$   $(\mathbf{x}_2, y_2)$   $(\mathbf{x}_3, y_3)$   $(\mathbf{x}_4, y_4)$   $(\mathbf{x}_5, y_5)$   $(\mathbf{x}_6, y_6)$

$(\mathbf{x}_5, y_5)$   $(\mathbf{x}_4, y_4)$   $(\mathbf{x}_5, y_5)$   $(\mathbf{x}_6, y_6)$   $(\mathbf{x}_2, y_2)$   $(\mathbf{x}_1, y_1)$

$(\mathbf{x}_1, y_1)$   $(\mathbf{x}_5, y_5)$   $(\mathbf{x}_3, y_3)$   $(\mathbf{x}_2, y_2)$   $(\mathbf{x}_5, y_5)$   $(\mathbf{x}_2, y_2)$

$(\mathbf{x}_6, y_6)$   $(\mathbf{x}_1, y_1)$   $(\mathbf{x}_3, y_3)$   $(\mathbf{x}_6, y_6)$   $(\mathbf{x}_3, y_3)$   $(\mathbf{x}_4, y_4)$

# Bootstrap Aggregation (Bagging)

- In order to average machine they must learn something different
- We only have one data set, but we can resample from the data set to make them look a bit different—this is known as **bootstrapping**

$(\mathbf{x}_1, y_1)$   $(\mathbf{x}_2, y_2)$   $(\mathbf{x}_3, y_3)$   $(\mathbf{x}_4, y_4)$   $(\mathbf{x}_5, y_5)$   $(\mathbf{x}_6, y_6)$

$(\mathbf{x}_5, y_5)$   $(\mathbf{x}_4, y_4)$   $(\mathbf{x}_5, y_5)$   $(\mathbf{x}_6, y_6)$   $(\mathbf{x}_2, y_2)$   $(\mathbf{x}_1, y_1)$

$(\mathbf{x}_1, y_1)$   $(\mathbf{x}_5, y_5)$   $(\mathbf{x}_3, y_3)$   $(\mathbf{x}_2, y_2)$   $(\mathbf{x}_5, y_5)$   $(\mathbf{x}_2, y_2)$

$(\mathbf{x}_6, y_6)$   $(\mathbf{x}_1, y_1)$   $(\mathbf{x}_3, y_3)$   $(\mathbf{x}_6, y_6)$   $(\mathbf{x}_3, y_3)$   $(\mathbf{x}_4, y_4)$

$(\mathbf{x}_4, y_4)$   $(\mathbf{x}_3, y_3)$   $(\mathbf{x}_1, y_1)$   $(\mathbf{x}_6, y_6)$   $(\mathbf{x}_4, y_4)$   $(\mathbf{x}_6, y_6)$

# Bootstrap Aggregation (Bagging)

- In order to average machine they must learn something different
- We only have one data set, but we can resample from the data set to make them look a bit different—this is known as **bootstrapping**

$(\mathbf{x}_1, y_1)$   $(\mathbf{x}_2, y_2)$   $(\mathbf{x}_3, y_3)$   $(\mathbf{x}_4, y_4)$   $(\mathbf{x}_5, y_5)$   $(\mathbf{x}_6, y_6)$

$(\mathbf{x}_5, y_5)$   $(\mathbf{x}_4, y_4)$   $(\mathbf{x}_5, y_5)$   $(\mathbf{x}_6, y_6)$   $(\mathbf{x}_2, y_2)$   $(\mathbf{x}_1, y_1)$

$(\mathbf{x}_1, y_1)$   $(\mathbf{x}_5, y_5)$   $(\mathbf{x}_3, y_3)$   $(\mathbf{x}_2, y_2)$   $(\mathbf{x}_5, y_5)$   $(\mathbf{x}_2, y_2)$

$(\mathbf{x}_6, y_6)$   $(\mathbf{x}_1, y_1)$   $(\mathbf{x}_3, y_3)$   $(\mathbf{x}_6, y_6)$   $(\mathbf{x}_3, y_3)$   $(\mathbf{x}_4, y_4)$

$(\mathbf{x}_4, y_4)$   $(\mathbf{x}_3, y_3)$   $(\mathbf{x}_1, y_1)$   $(\mathbf{x}_6, y_6)$   $(\mathbf{x}_4, y_4)$   $(\mathbf{x}_6, y_6)$

$(\mathbf{x}_3, y_3)$   $(\mathbf{x}_2, y_2)$   $(\mathbf{x}_3, y_3)$   $(\mathbf{x}_4, y_4)$   $(\mathbf{x}_6, y_6)$   $(\mathbf{x}_2, y_2)$



# Random Forest

- To make decision trees more distinct we can choose a subset of variables on which to split the tree
- Bagging with trees is known as **random forest** and is a very powerful technique
- In many instances it has been superseded by boosting algorithms

# Random Forest

- To make decision trees more distinct we can choose a subset of variables on which to split the tree
- Bagging with trees is known as **random forest** and is a very powerful technique
- In many instances it has been superseded by boosting algorithms

# Random Forest

- To make decision trees more distinct we can choose a subset of variables on which to split the tree
- Bagging with trees is known as **random forest** and is a very powerful technique
- In many instances it has been superseded by boosting algorithms

# Boosting

- In boosting we make a **strong learner** by using a weighted sum of **weak learners**

$$\hat{f}(\mathbf{x}) = \sum_{i=1}^n w_i \hat{h}_i(\mathbf{x})$$

- Weak learners ( $\hat{h}_i(\mathbf{x})$ ) are learning machine that do a little better than chance
- The trick is to choose the weights  $w_i$

# Boosting

- In boosting we make a **strong learner** by using a weighted sum of **weak learners**

$$\hat{f}(\mathbf{x}) = \sum_{i=1}^n w_i \hat{h}_i(\mathbf{x})$$

- Weak learners ( $\hat{h}_i(\mathbf{x})$ ) are learning machine that do a little better than chance
- The trick is to choose the weights  $w_i$

# Boosting

- In boosting we make a **strong learner** by using a weighted sum of **weak learners**

$$\hat{f}(\mathbf{x}) = \sum_{i=1}^n w_i \hat{h}_i(\mathbf{x})$$

- Weak learners ( $\hat{h}_i(\mathbf{x})$ ) are learning machine that do a little better than chance
- The trick is to choose the weights  $w_i$

# Shallow Trees

- One of the most effective type of weak learner are very shallow trees, (e.g.  $\sqrt{p}$ )
- Use at most only small number of random variables, (e.g.  $\sqrt{p}$ )
- There are different algorithms for choosing the weights
  - ★ adaboost
  - ★ gradient boosting
- Boosting can overfit, but because we use weak learners it tends to do this rather slowly

# Shallow Trees

- One of the most effective type of weak learner are very shallow trees, (e.g.  $\sqrt{p}$ )
- Use at most only small number of random variables, (e.g.  $\sqrt{p}$ )
- There are different algorithms for choosing the weights
  - ★ adaboost
  - ★ gradient boosting
- Boosting can overfit, but because we use weak learners it tends to do this rather slowly



# Shallow Trees

- One of the most effective type of weak learner are very shallow trees, (e.g.  $\sqrt{p}$ )
- Use at most only small number of random variables, (e.g.  $\sqrt{p}$ )
- There are different algorithms for choosing the weights
  - ★ adaboost
  - ★ gradient boosting
- Boosting can overfit, but because we use weak learners it tends to do this rather slowly

# Shallow Trees

- One of the most effective type of weak learner are very shallow trees, (e.g.  $\sqrt{p}$ )
- Use at most only small number of random variables, (e.g.  $\sqrt{p}$ )
- There are different algorithms for choosing the weights
  - ★ adaboost—a classic algorithm for binary classification problems
  - ★ gradient boosting
- Boosting can overfit, but because we use weak learners it tends to do this rather slowly

# Shallow Trees

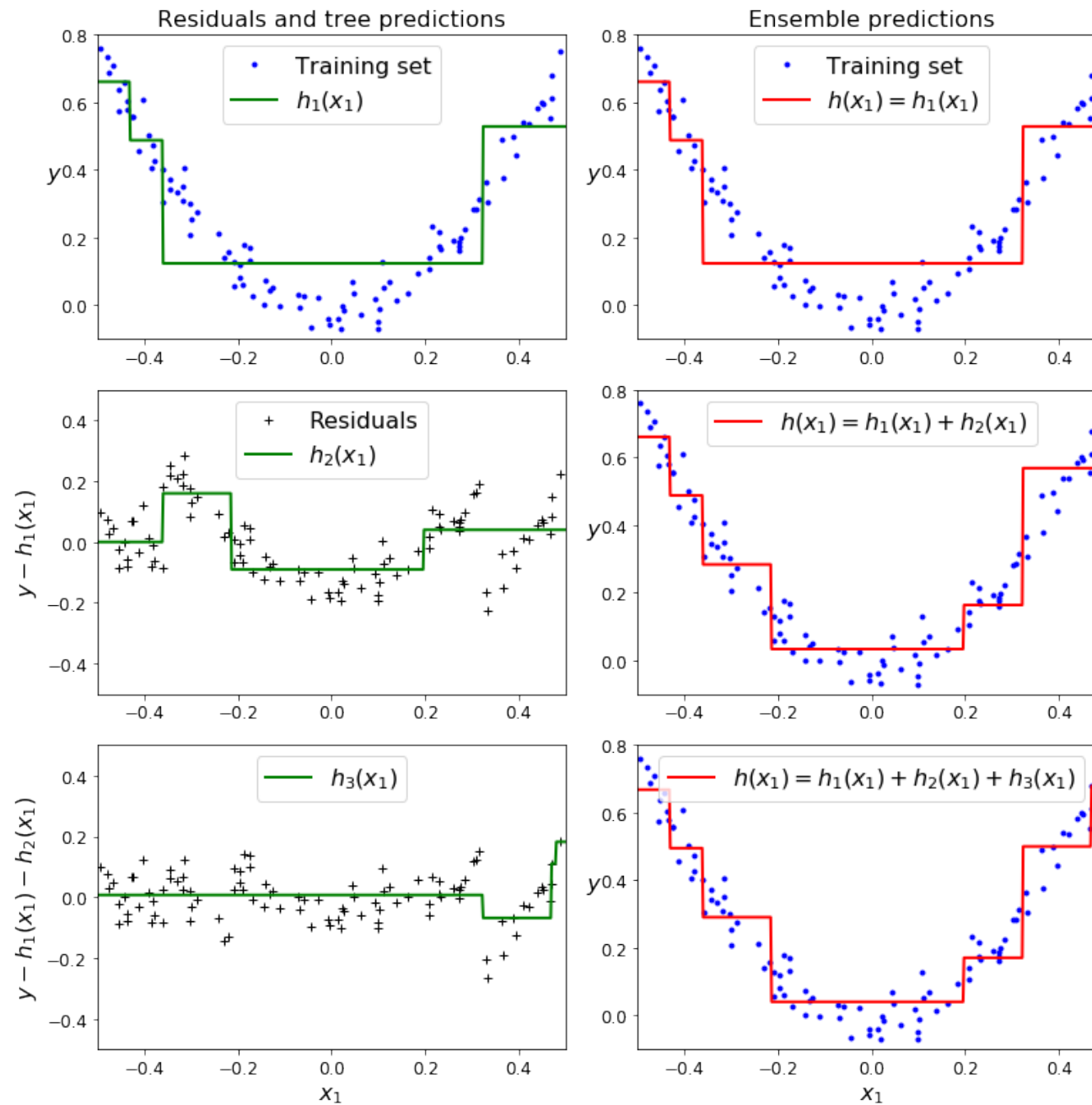
- One of the most effective type of weak learner are very shallow trees, (e.g.  $\sqrt{p}$ )
- Use at most only small number of random variables, (e.g.  $\sqrt{p}$ )
- There are different algorithms for choosing the weights
  - ★ adaboost—a classic algorithm for binary classification problems
  - ★ gradient boosting
- Boosting can overfit, but because we use weak learners it tends to do this rather slowly

# Shallow Trees

- One of the most effective type of weak learner are very shallow trees, (e.g.  $\sqrt{p}$ )
- Use at most only small number of random variables, (e.g.  $\sqrt{p}$ )
- There are different algorithms for choosing the weights
  - ★ adaboost—a classic algorithm for binary classification problems
  - ★ gradient boosting—regression algorithm, train a classifier on the residual errors
- Boosting can overfit, but because we use weak learners it tends to do this rather slowly

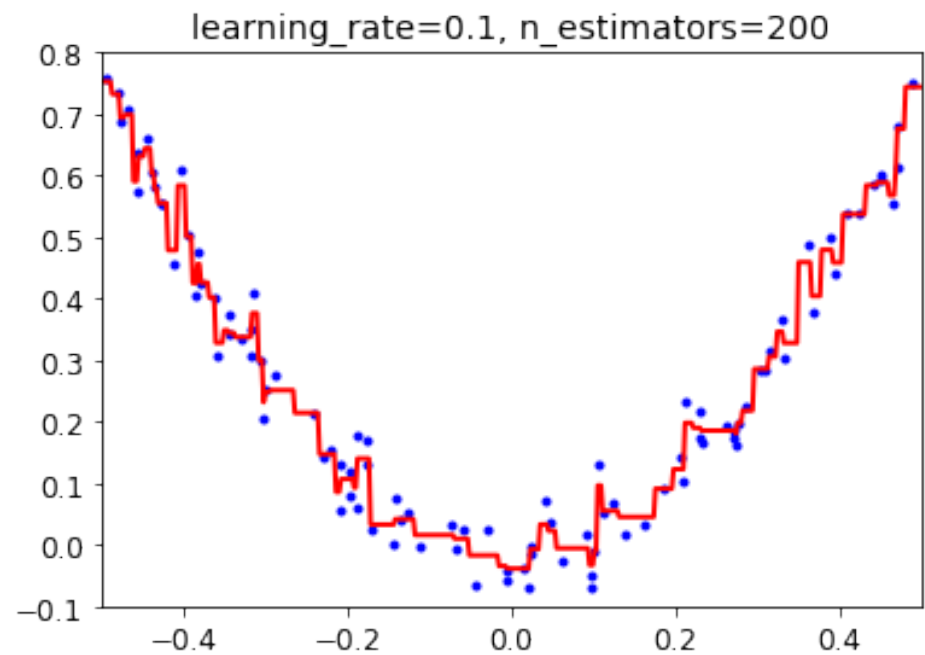
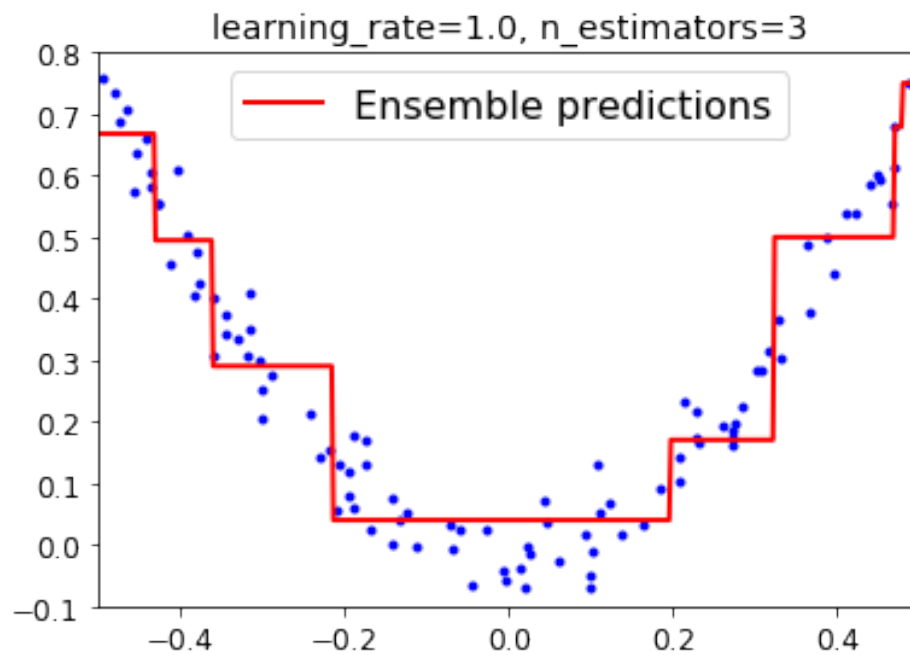
# Shallow Trees

- One of the most effective type of weak learner are very shallow trees, (e.g.  $\sqrt{p}$ )
- Use at most only small number of random variables, (e.g.  $\sqrt{p}$ )
- There are different algorithms for choosing the weights
  - ★ adaboost—a classic algorithm for binary classification problems
  - ★ gradient boosting—regression algorithm, train a classifier on the residual errors
- Boosting can overfit, but because we use weak learners it tends to do this rather slowly



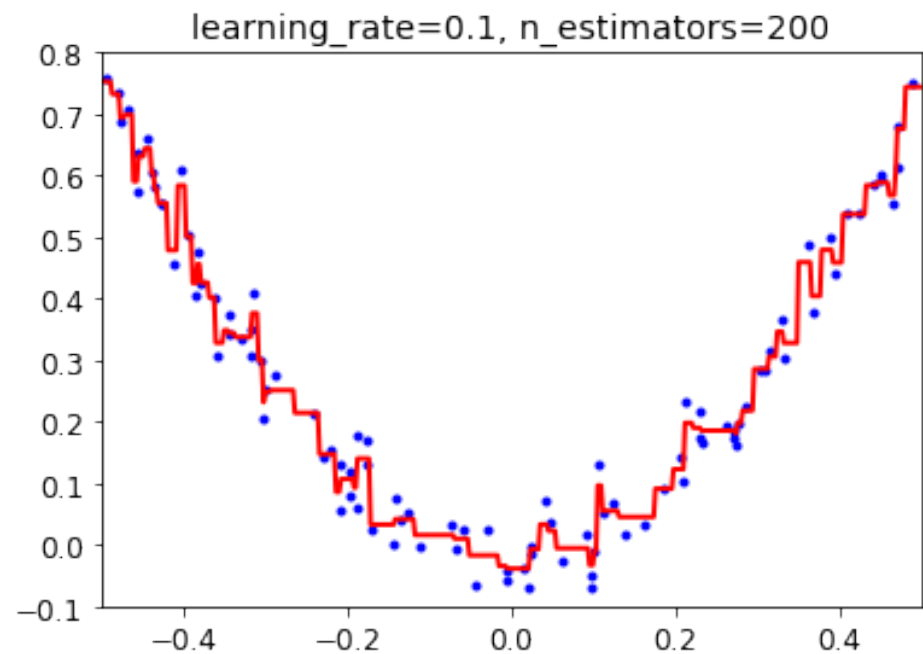
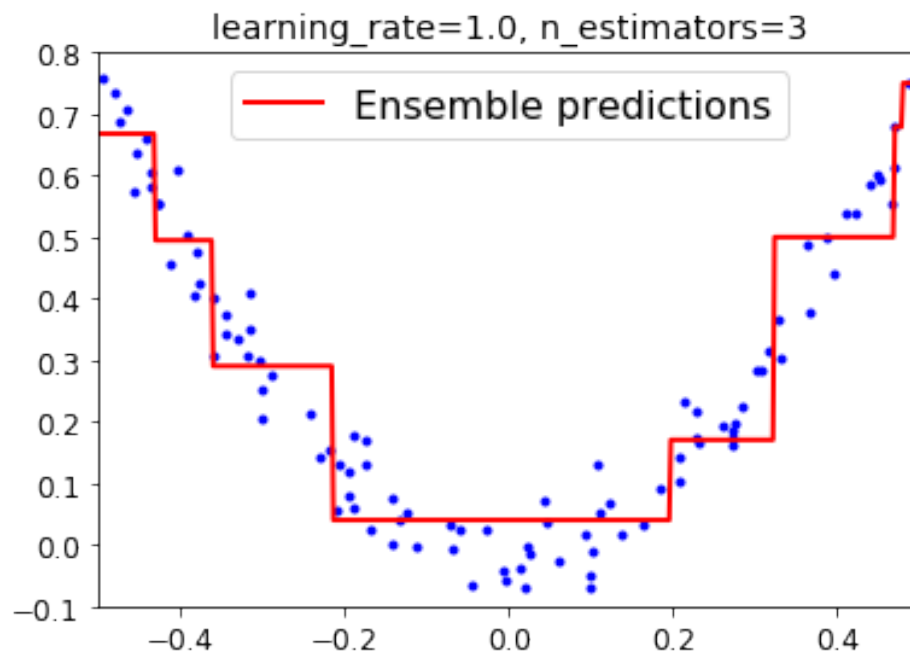
# Keep On Going

- We can keep on going



# Keep On Going

- We can keep on going

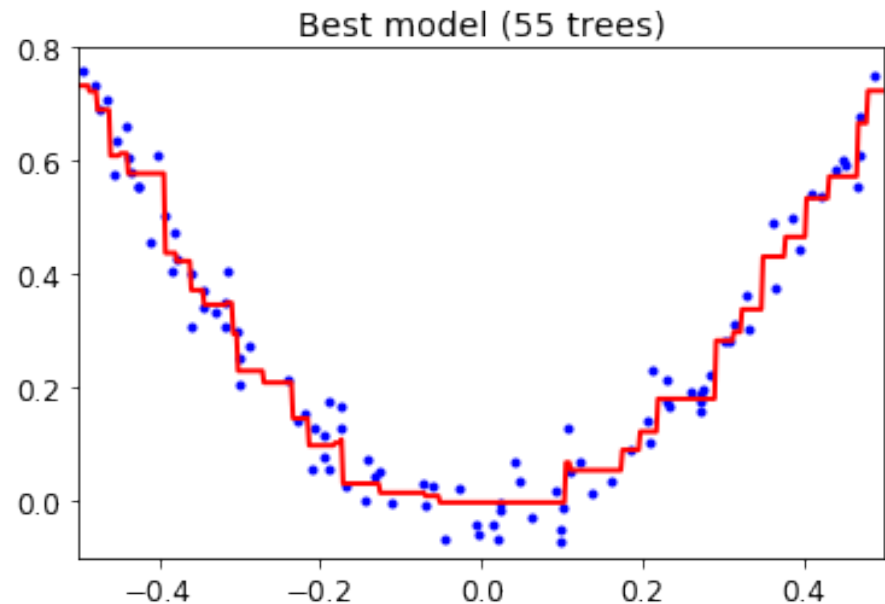
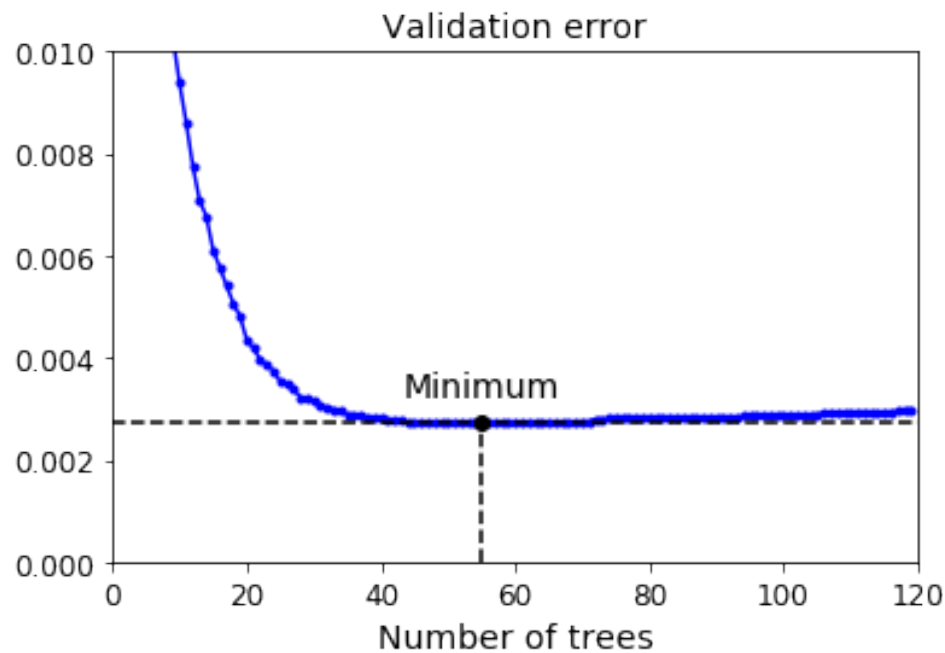


- But we will over-fit eventually



# Early Stopping

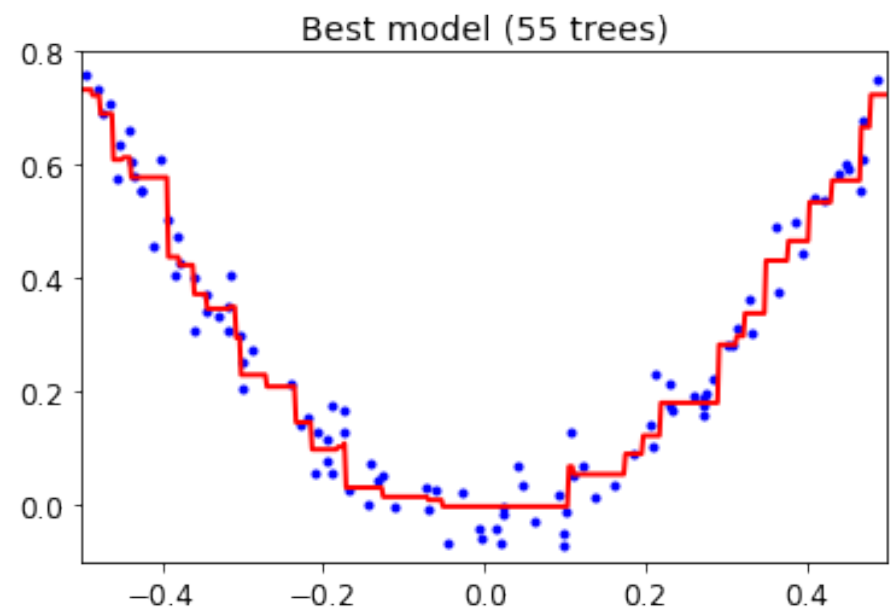
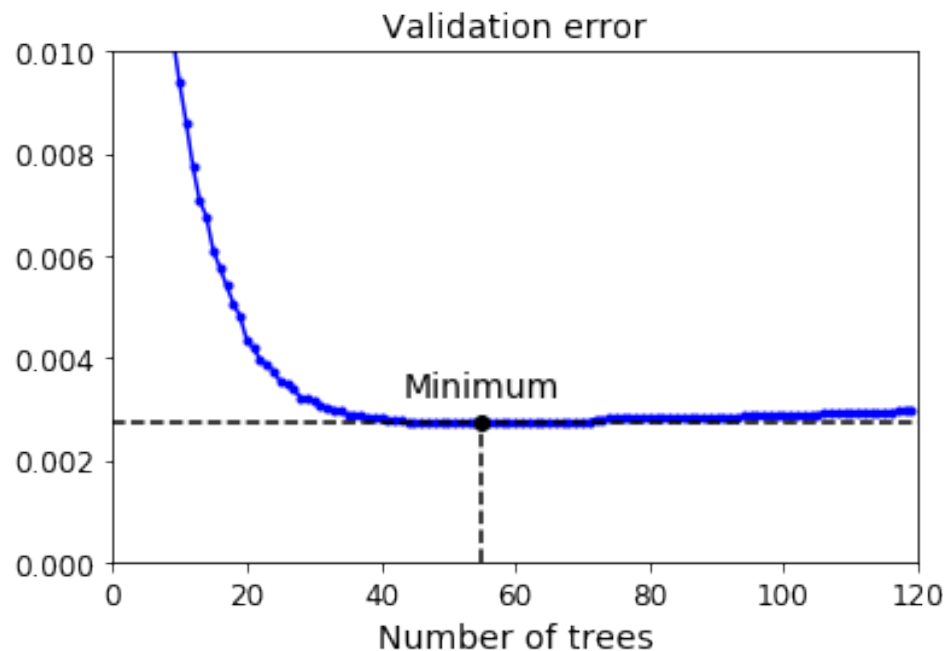
- Like many algorithms we often get better results by early stopping



- Use cross-validation against a validation set to decide when to stop

# Early Stopping

- Like many algorithms we often get better results by early stopping



- Use cross-validation against a validation set to decide when to stop

# XGBoost

- XGBoost is an implementation of gradient boosting that won the Higg's Boson challenge and regularly wins Kaggle competitions
- XGBoost stands for eXtreme Gradient Boosting
- It is much faster than most gradient boosting algorithms and scales to billions of training data points
- It uses a cleverly chosen regularisation term to favour simple trees
- Finds a clever way to approximately minimise error plus regulariser very fast

# XGBoost

- XGBoost is an implementation of gradient boosting that won the Higg's Boson challenge and regularly wins Kaggle competitions
- XGBoost stands for eXtreme Gradient Boosting
- It is much faster than most gradient boosting algorithms and scales to billions of training data points
- It uses a cleverly chosen regularisation term to favour simple trees
- Finds a clever way to approximately minimise error plus regulariser very fast

# XGBoost

- XGBoost is an implementation of gradient boosting that won the Higg's Boson challenge and regularly wins Kaggle competitions
- XGBoost stands for eXtreme Gradient Boosting
- It is much faster than most gradient boosting algorithms and scales to billions of training data points
- It uses a cleverly chosen regularisation term to favour simple trees
- Finds a clever way to approximately minimise error plus regulariser very fast

# XGBoost

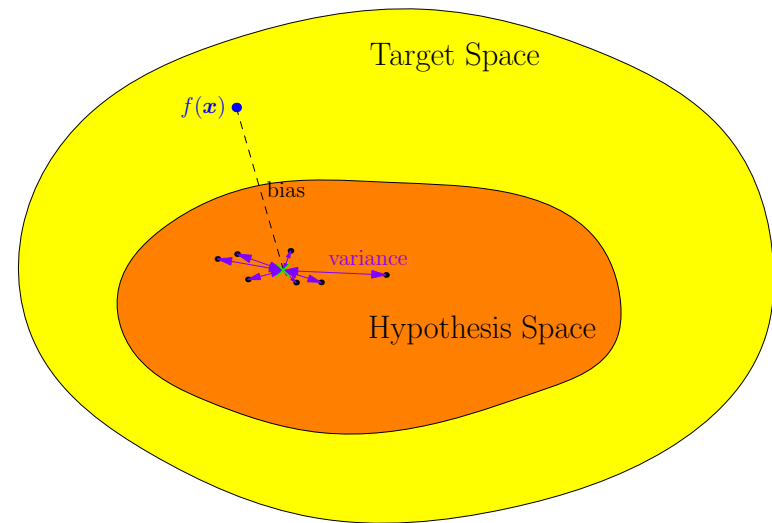
- XGBoost is an implementation of gradient boosting that won the Higg's Boson challenge and regularly wins Kaggle competitions
- XGBoost stands for eXtreme Gradient Boosting
- It is much faster than most gradient boosting algorithms and scales to billions of training data points
- It uses a cleverly chosen regularisation term to favour simple trees
- Finds a clever way to approximately minimise error plus regulariser very fast

# XGBoost

- XGBoost is an implementation of gradient boosting that won the Higg's Boson challenge and regularly wins Kaggle competitions
- XGBoost stands for eXtreme Gradient Boosting
- It is much faster than most gradient boosting algorithms and scales to billions of training data points
- It uses a cleverly chosen regularisation term to favour simple trees
- Finds a clever way to approximately minimise error plus regulariser very fast

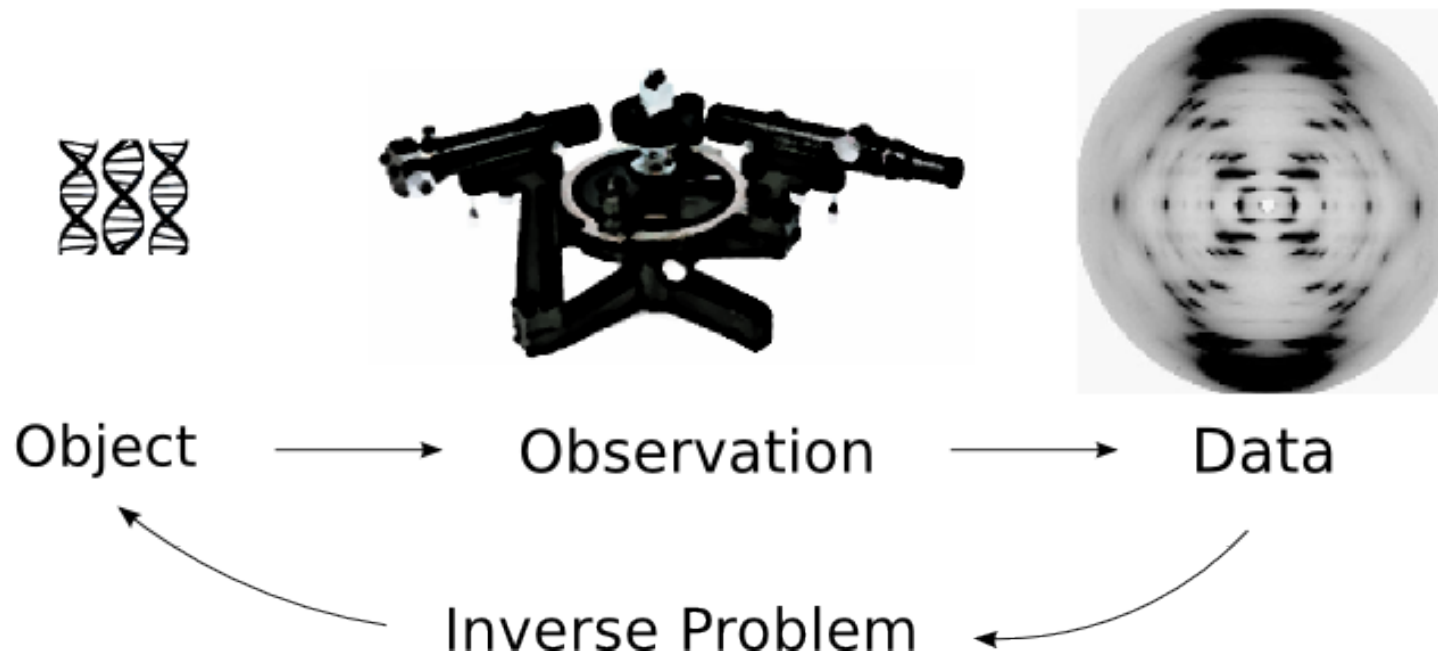
# Outline

1. What Makes a Good Learning Machine?
2. SVMs
3. Ensemble Methods
4. **Bayesian Inference**



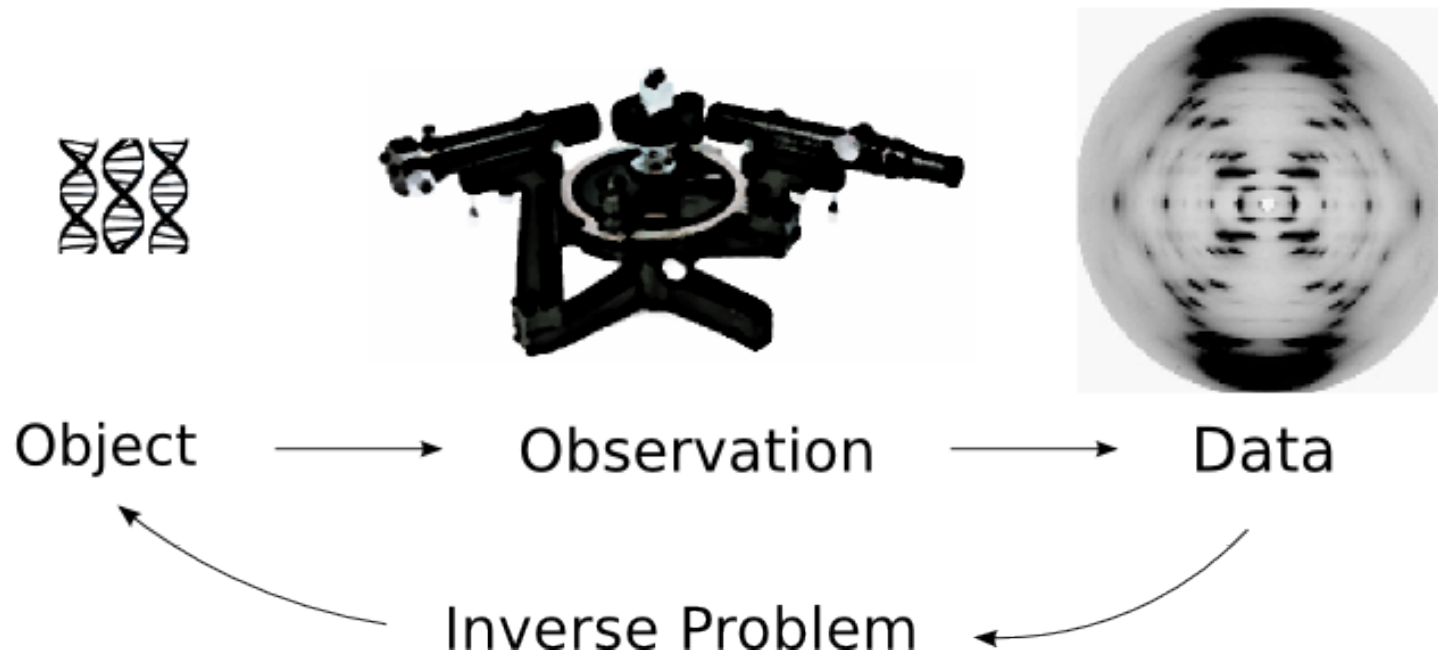


# Inverse Problems



- Machine learning can often be seen as solving inverse problems
- We are given some data generated by the world and we want to infer something about the world

# Inverse Problems



- Machine learning can often be seen as solving inverse problems
- We are given some data generated by the world and we want to infer something about the world

# Bayes' Rule

- A trivial identity in probability known as Bayes' rules tells you how to solve inverse problems

$$\mathbb{P}(W|\mathcal{D}) = \frac{\mathbb{P}(\mathcal{D}|W) \mathbb{P}(W)}{\mathbb{P}(D)}$$

- What we want is to know the probability of the world,  $W$ , given the data,  $\mathcal{D}$  we have observed—this is known as the **posteriori** probability
- This depends on the **likelihood** of the data given the world  $\mathbb{P}(\mathcal{D}|W)$
- Multiplied by the **prior probability** of the world,  $\mathbb{P}(W)$

# Bayes' Rule

- A trivial identity in probability known as Bayes' rules tells you how to solve inverse problems

$$\mathbb{P}(W|\mathcal{D}) = \frac{\mathbb{P}(\mathcal{D}|W) \mathbb{P}(W)}{\mathbb{P}(D)}$$

- What we want is to know the probability of the world,  $W$ , given the data,  $\mathcal{D}$  we have observed—this is known as the **posteriori** probability
- This depends on the **likelihood** of the data given the world  $\mathbb{P}(\mathcal{D}|W)$
- Multiplied by the **prior probability** of the world,  $\mathbb{P}(W)$

# Bayes' Rule

- A trivial identity in probability known as Bayes' rules tells you how to solve inverse problems

$$\mathbb{P}(W|\mathcal{D}) = \frac{\mathbb{P}(\mathcal{D}|W) \mathbb{P}(W)}{\mathbb{P}(D)}$$

- What we want is to know the probability of the world,  $W$ , given the data,  $\mathcal{D}$  we have observed—this is known as the **posteriori probability**
- This depends on the **likelihood** of the data given the world  $\mathbb{P}(\mathcal{D}|W)$
- Multiplied by the **prior probability** of the world,  $\mathbb{P}(W)$

# Bayes' Rule

- A trivial identity in probability known as Bayes' rules tells you how to solve inverse problems

$$\mathbb{P}(W|\mathcal{D}) = \frac{\mathbb{P}(\mathcal{D}|W) \mathbb{P}(W)}{\mathbb{P}(D)}$$

- What we want is to know the probability of the world,  $W$ , given the data,  $\mathcal{D}$  we have observed—this is known as the **posteriori** probability
- This depends on the **likelihood** of the data given the world  $\mathbb{P}(\mathcal{D}|W)$
- Multiplied by the **prior probability** of the world,  $\mathbb{P}(W)$

# Bayes' Rule

- A trivial identity in probability known as Bayes' rules tells you how to solve inverse problems

$$\mathbb{P}(W|\mathcal{D}) = \frac{\mathbb{P}(\mathcal{D}|W) \mathbb{P}(W)}{\mathbb{P}(D)}$$

- What we want is to know the probability of the world,  $W$ , given the data,  $\mathcal{D}$  we have observed—this is known as the **posteriori** probability
- This depends on the **likelihood** of the data given the world  $\mathbb{P}(\mathcal{D}|W)$
- Multiplied by the **prior probability** of the world,  $\mathbb{P}(W)$

# Normalisation

- The denominator is a normalisation constant

$$\mathbb{P}(D) = \sum_W \mathbb{P}(\mathcal{D}|W) \mathbb{P}(W)$$

- It is useful for comparing between different models of the world,  $W$ , and is sometimes called the **posterior likelihood** or **evidence**
- The model with the largest evidence is the most likely model to be correct



# Normalisation

- The denominator is a normalisation constant

$$\mathbb{P}(D) = \sum_W \mathbb{P}(\mathcal{D}|W) \mathbb{P}(W)$$

- It is useful for comparing between different models of the world,  $W$ , and is sometimes called the **posterior likelihood** or **evidence**
- The model with the largest evidence is the most likely model to be correct

# Normalisation

- The denominator is a normalisation constant

$$\mathbb{P}(D) = \sum_W \mathbb{P}(\mathcal{D}|W) \mathbb{P}(W)$$

- It is useful for comparing between different models of the world,  $W$ , and is sometimes called the **posterior likelihood** or **evidence**
- The model with the largest evidence is the most likely model to be correct

# Normalisation

- The denominator is a normalisation constant

$$\mathbb{P}(D) = \sum_W \mathbb{P}(\mathcal{D}|W) \mathbb{P}(W)$$

- It is useful for comparing between different models of the world,  $W$ , and is sometimes called the **posterior likelihood** or **evidence**
- The model with the largest evidence is the most likely model to be correct—used for model selection

# Worlds

- $W$  is some model of the world—technically they are parameters of the likelihood function
- If we trying to infer the length of rod then it would be a length
- If we are trying to infer the DNA sequence given a set of short reads then  $W$  would be a DNA sequence
- If we are trying to infer the crystal structure of molecules then  $W$  would be a structure

# Worlds

- $W$  is some model of the world—technically they are parameters of the likelihood function
- If we trying to infer the length of rod then it would be a length
- If we are trying to infer the DNA sequence given a set of short reads then  $W$  would be a DNA sequence
- If we are trying to infer the crystal structure of molecules then  $W$  would be a structure

# Worlds

- $W$  is some model of the world—technically they are parameters of the likelihood function
- If we trying to infer the length of rod then it would be a length
- If we are trying to infer the DNA sequence given a set of short reads then  $W$  would be a DNA sequence
- If we are trying to infer the crystal structure of molecules then  $W$  would be a structure

# Worlds

- $W$  is some model of the world—technically they are parameters of the likelihood function
- If we trying to infer the length of rod then it would be a length
- If we are trying to infer the DNA sequence given a set of short reads then  $W$  would be a DNA sequence
- If we are trying to infer the crystal structure of molecules then  $W$  would be a structure

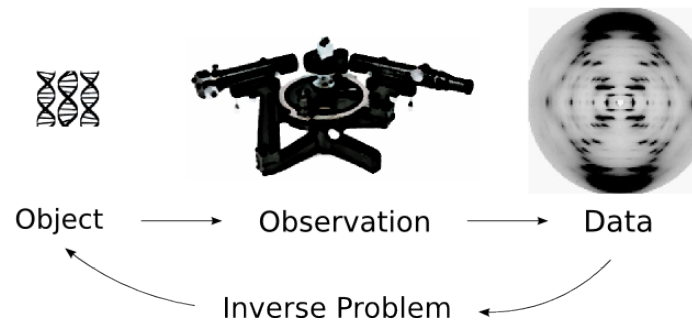
# Worlds

- $W$  is some model of the world—technically they are parameters of the likelihood function
- If we trying to infer the length of rod then it would be a length
- If we are trying to infer the DNA sequence given a set of short reads then  $W$  would be a DNA sequence
- If we are trying to infer the crystal structure of molecules then  $W$  would be a structure



# Advantages of Bayesian Inference

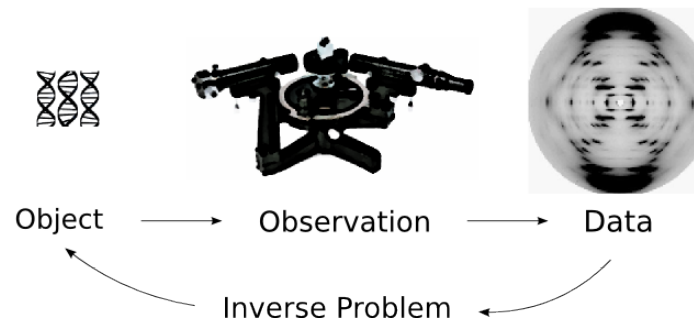
- Its optimal (follows from a simple probabilistic identity)
- It is relatively easy to model the likelihood (this is a forward process not an inverse process)



- No problem with missing data—we just calculate the likelihood of the data we see
- Won't over-fit
- Provides a full probabilistic description of the inference

# Advantages of Bayesian Inference

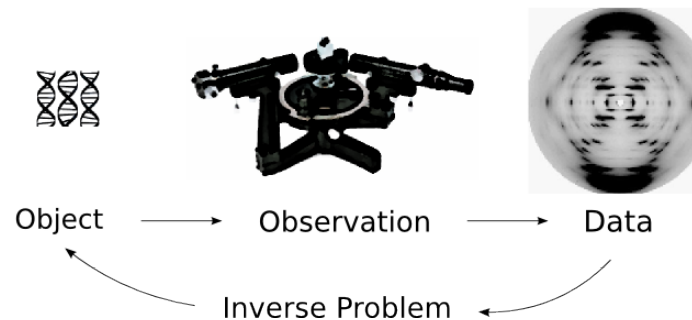
- Its optimal (follows from a simple probabilistic identity)
- It is relatively easy to model the likelihood (this is a forward process not an inverse process)



- No problem with missing data—we just calculate the likelihood of the data we see
- Won't over-fit
- Provides a full probabilistic description of the inference

# Advantages of Bayesian Inference

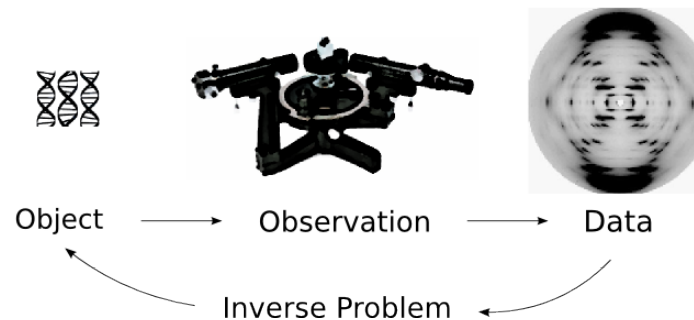
- Its optimal (follows from a simple probabilistic identity)
- It is relatively easy to model the likelihood (this is a forward process not an inverse process)



- No problem with missing data—we just calculate the likelihood of the data we see
- Won't over-fit
- Provides a full probabilistic description of the inference

# Advantages of Bayesian Inference

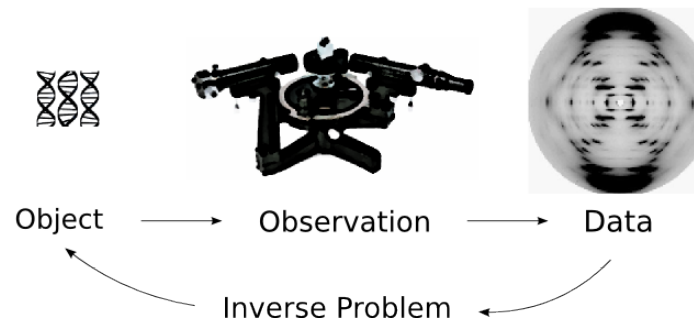
- Its optimal (follows from a simple probabilistic identity)
- It is relatively easy to model the likelihood (this is a forward process not an inverse process)



- No problem with missing data—we just calculate the likelihood of the data we see
- Won't over-fit
- Provides a full probabilistic description of the inference

# Advantages of Bayesian Inference

- Its optimal (follows from a simple probabilistic identity)
- It is relatively easy to model the likelihood (this is a forward process not an inverse process)



- No problem with missing data—we just calculate the likelihood of the data we see
- Won't over-fit
- Provides a full probabilistic description of the inference

# Disadvantages of Bayesian Inference

- It is a bespoke method, you need to model the likelihood and put a prior probability on the world
- Need to be able to sum over all possible worlds
- This is technically challenging and/or computationally slow
- Markov Chain Monte Carlo (MCMC) is a very established framework for summing over worlds
  - ★ There exist mature libraries
  - ★ Large number of tricks to speed things up
- Various approximations also exist for averaging over worlds

# Disadvantages of Bayesian Inference

- It is a bespoke method, you need to model the likelihood and put a prior probability on the world—not out of the box
- Need to be able to sum over all possible worlds
- This is technically challenging and/or computationally slow
- Markov Chain Monte Carlo (MCMC) is a very established framework for summing over worlds
  - ★ There exist mature libraries
  - ★ Large number of tricks to speed things up
- Various approximations also exist for averaging over worlds

# Disadvantages of Bayesian Inference

- It is a bespoke method, you need to model the likelihood and put a prior probability on the world—not out of the box
- Need to be able to sum over all possible worlds
- This is technically challenging and/or computationally slow
- Markov Chain Monte Carlo (MCMC) is a very established framework for summing over worlds
  - ★ There exist mature libraries
  - ★ Large number of tricks to speed things up
- Various approximations also exist for averaging over worlds



# Disadvantages of Bayesian Inference

- It is a bespoke method, you need to model the likelihood and put a prior probability on the world—not out of the box
- Need to be able to sum over all possible worlds
- This is technically challenging and/or computationally slow
- Markov Chain Monte Carlo (MCMC) is a very established framework for summing over worlds
  - ★ There exist mature libraries
  - ★ Large number of tricks to speed things up
- Various approximations also exist for averaging over worlds

# Disadvantages of Bayesian Inference

- It is a bespoke method, you need to model the likelihood and put a prior probability on the world—not out of the box
- Need to be able to sum over all possible worlds
- This is technically challenging and/or computationally slow
- Markov Chain Monte Carlo (MCMC) is a very established framework for summing over worlds
  - ★ There exist mature libraries
  - ★ Large number of tricks to speed things up
- Various approximations also exist for averaging over worlds

# Disadvantages of Bayesian Inference

- It is a bespoke method, you need to model the likelihood and put a prior probability on the world—not out of the box
- Need to be able to sum over all possible worlds
- This is technically challenging and/or computationally slow
- Markov Chain Monte Carlo (MCMC) is a very established framework for summing over worlds
  - ★ There exist mature libraries
  - ★ Large number of tricks to speed things up
- Various approximations also exist for averaging over worlds

# Simple Bayes

- There are a few cases where Bayesian Inference is relatively simple
- When we are inferring quantities with a very simple likelihood, sometimes Bayesian approach becomes trivial
- There are a few cases where we can use the same model for different situations and the model is easy to compute
  - ★ Gaussian processes for regression
  - ★ Naive Bayes (e.g. for text)

# Simple Bayes

- There are a few cases where Bayesian Inference is relatively simple
- When we are inferring quantities with a very simple likelihood, sometimes Bayesian approach becomes trivial
- There are a few cases where we can use the same model for different situations and the model is easy to compute
  - ★ Gaussian processes for regression
  - ★ Naive Bayes (e.g. for text)

# Simple Bayes

- There are a few cases where Bayesian Inference is relatively simple
- When we are inferring quantities with a very simple likelihood, sometimes Bayesian approach becomes trivial
- There are a few cases where we can use the same model for different situations and the model is easy to compute
  - ★ Gaussian processes for regression
  - ★ Naive Bayes (e.g. for text)

# Simple Bayes

- There are a few cases where Bayesian Inference is relatively simple
- When we are inferring quantities with a very simple likelihood, sometimes Bayesian approach becomes trivial
- There are a few cases where we can use the same model for different situations and the model is easy to compute
  - ★ Gaussian processes for regression
  - ★ Naive Bayes (e.g. for text)

# Simple Bayes

- There are a few cases where Bayesian Inference is relatively simple
- When we are inferring quantities with a very simple likelihood, sometimes Bayesian approach becomes trivial
- There are a few cases where we can use the same model for different situations and the model is easy to compute
  - ★ Gaussian processes for regression
  - ★ Naive Bayes (e.g. for text)



# Gaussian Processes

- Gaussian Processes is an off the shelf method for regression (multi-dimensional curve fitting)
- It makes a Gaussian (normal distribution) assumption
  - ★ The probability of a point taking a value is normally distributed with a mean and variance that depends on all the observed points
  - ★ The dependence is controlled by a covariance function which you are allowed to choose
- Because everything is Gaussian, all sums over worlds can be done
- Mathematics is all doable

# Gaussian Processes

- Gaussian Processes is an off the shelf method for regression (multi-dimensional curve fitting)
- It makes a Gaussian (normal distribution) assumption
  - ★ The probability of a point taking a value is normally distributed with a mean and variance that depends on all the observed points
  - ★ The dependence is controlled by a covariance function which you are allowed to choose
- Because everything is Gaussian, all sums over worlds can be done
- Mathematics is all doable

# Gaussian Processes

- Gaussian Processes is an off the shelf method for regression (multi-dimensional curve fitting)
- It makes a Gaussian (normal distribution) assumption
  - ★ The probability of a point taking a value is normally distributed with a mean and variance that depends on all the observed points
  - ★ The dependence is controlled by a covariance function which you are allowed to choose
- Because everything is Gaussian, all sums over worlds can be done
- Mathematics is all doable

# Gaussian Processes

- Gaussian Processes is an off the shelf method for regression (multi-dimensional curve fitting)
- It makes a Gaussian (normal distribution) assumption
  - ★ The probability of a point taking a value is normally distributed with a mean and variance that depends on all the observed points
  - ★ The dependence is controlled by a covariance function which you are allowed to choose
- Because everything is Gaussian, all sums over worlds can be done
- Mathematics is all doable

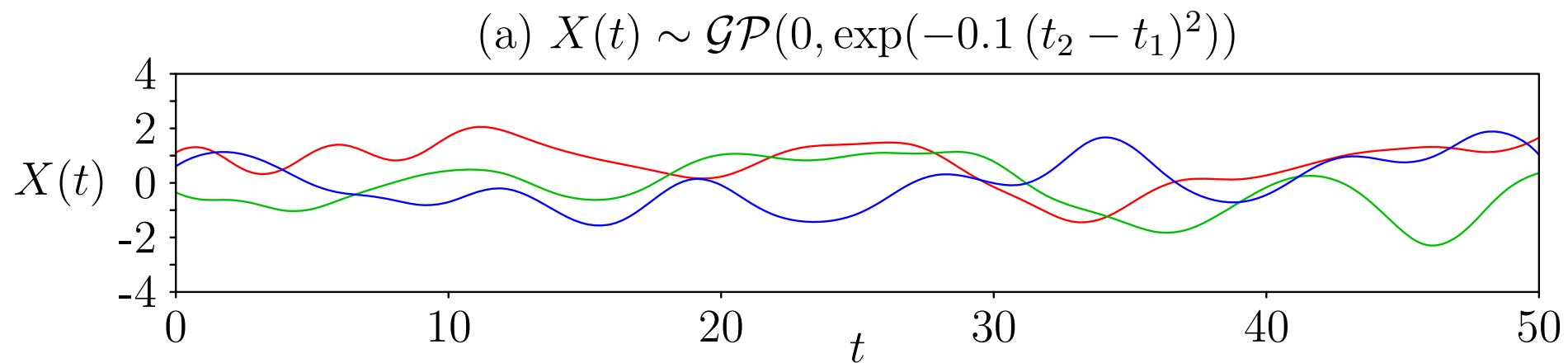
# Gaussian Processes

- Gaussian Processes is an off the shelf method for regression (multi-dimensional curve fitting)
- It makes a Gaussian (normal distribution) assumption
  - ★ The probability of a point taking a value is normally distributed with a mean and variance that depends on all the observed points
  - ★ The dependence is controlled by a covariance function which you are allowed to choose
- Because everything is Gaussian, all sums over worlds can be done
- Mathematics is all doable

# Gaussian Processes

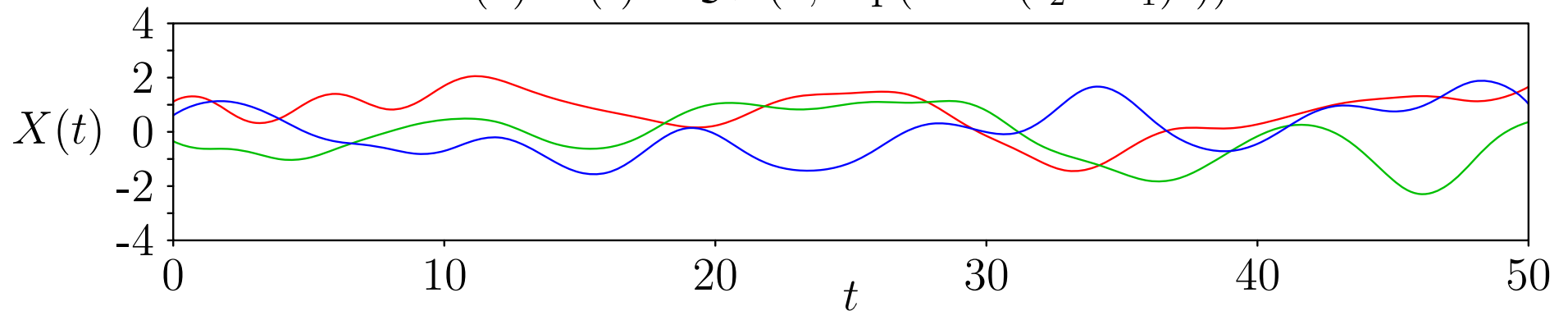
- Gaussian Processes is an off the shelf method for regression (multi-dimensional curve fitting)
- It makes a Gaussian (normal distribution) assumption
  - ★ The probability of a point taking a value is normally distributed with a mean and variance that depends on all the observed points
  - ★ The dependence is controlled by a covariance function which you are allowed to choose
- Because everything is Gaussian, all sums over worlds can be done
- Mathematics is all doable, but a real pain!

# Gaussian Process Worlds

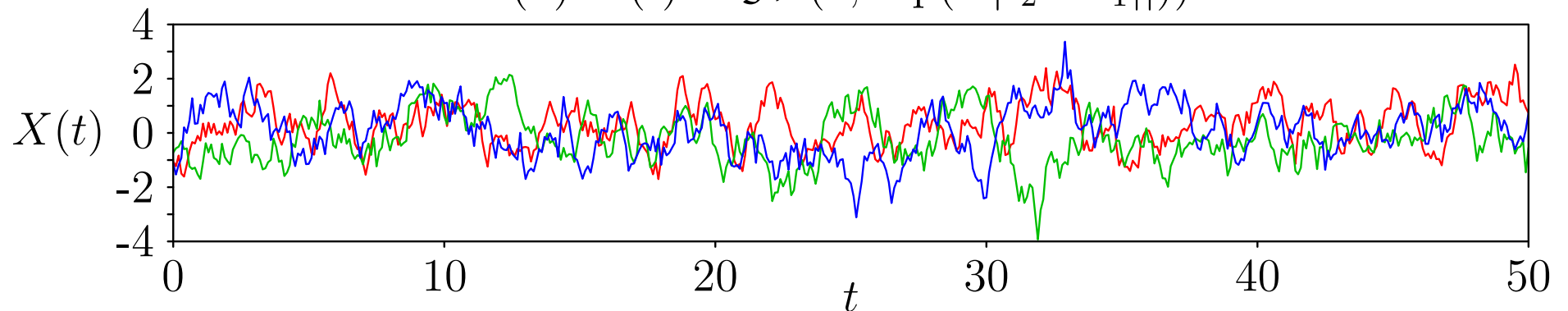


# Gaussian Process Worlds

(a)  $X(t) \sim \mathcal{GP}(0, \exp(-0.1(t_2 - t_1)^2))$

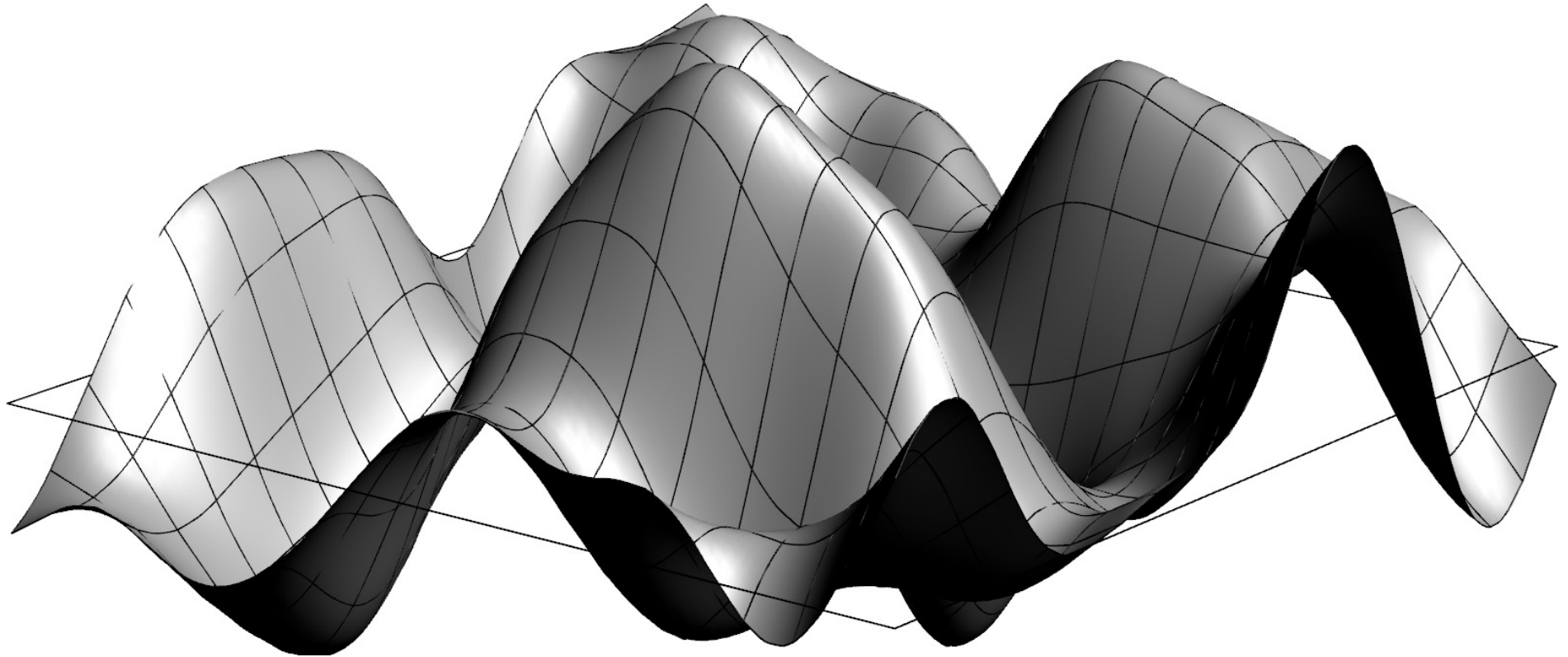


(b)  $X(t) \sim \mathcal{GP}(0, \exp(-|t_2 - t_1|))$

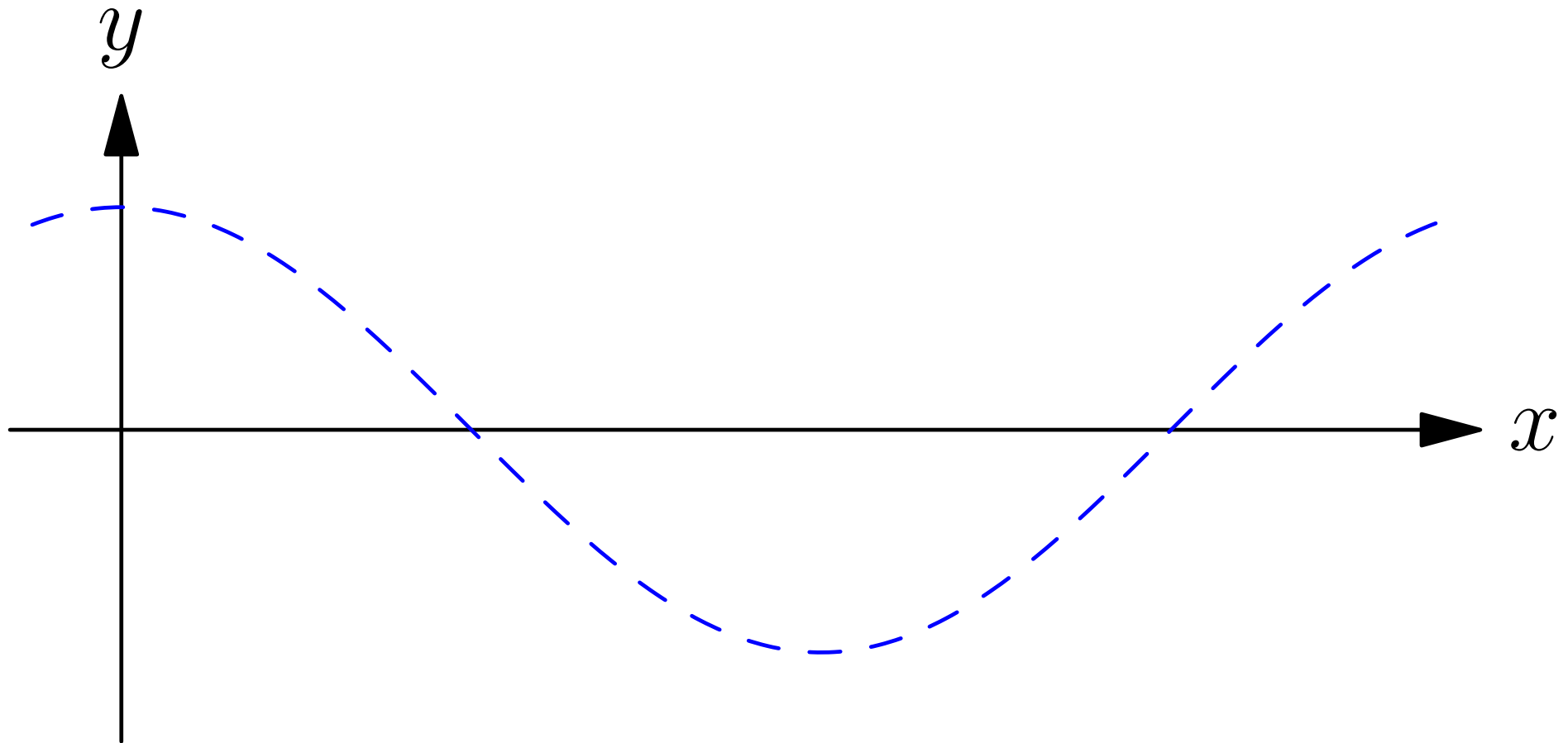




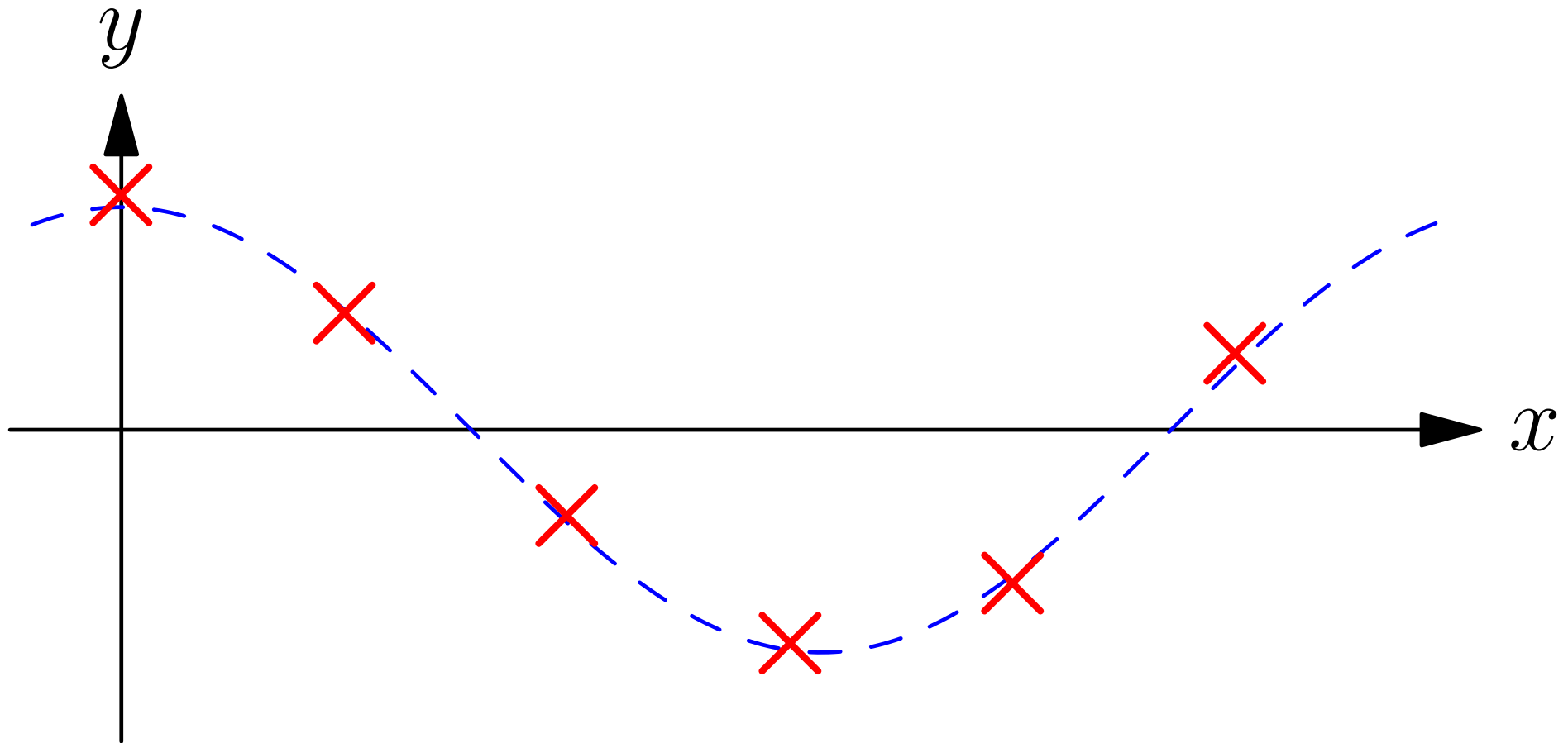
# 2-D Gaussian Processes



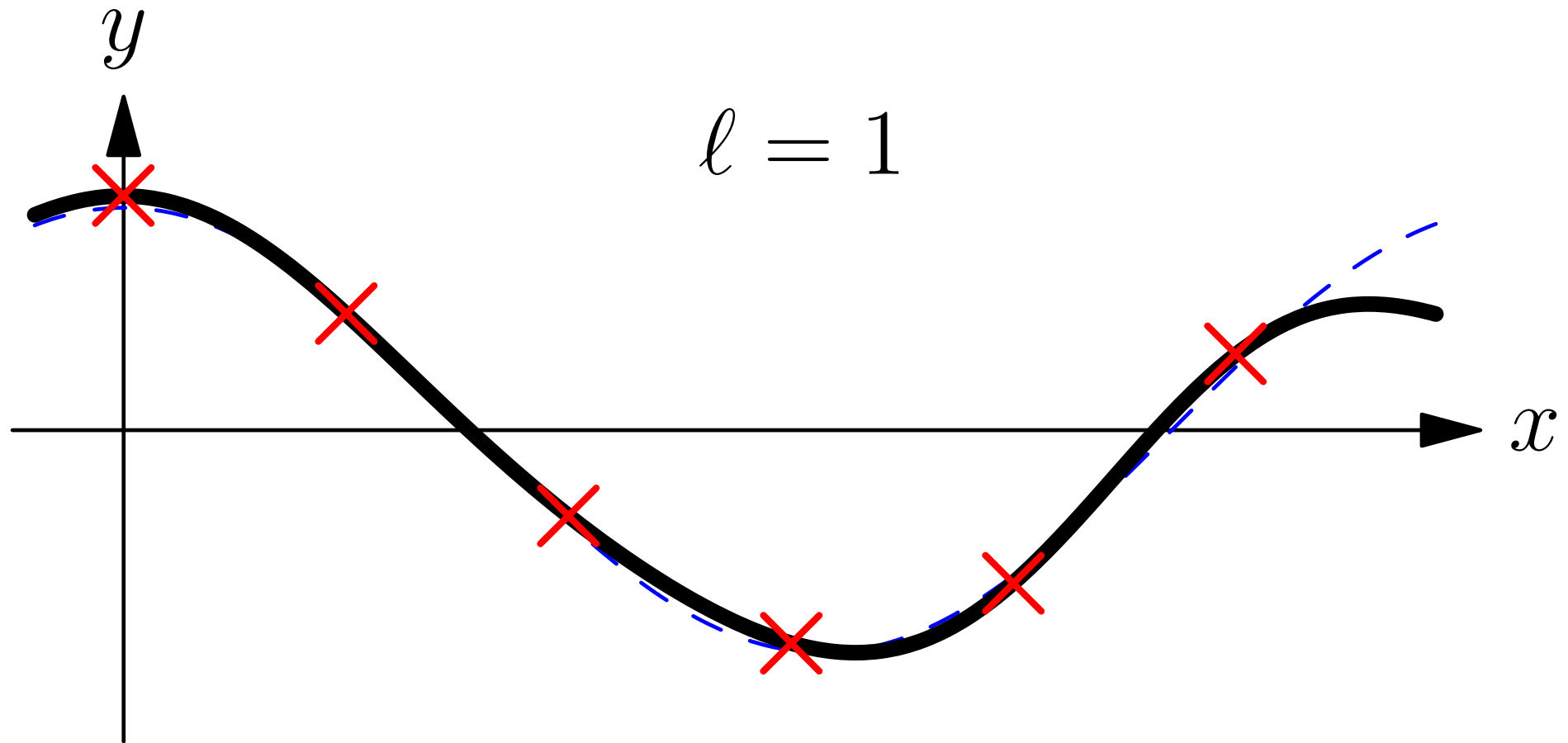
$$K(x, x') = \exp(-(x - x')^2 / (2 \ell^2))$$



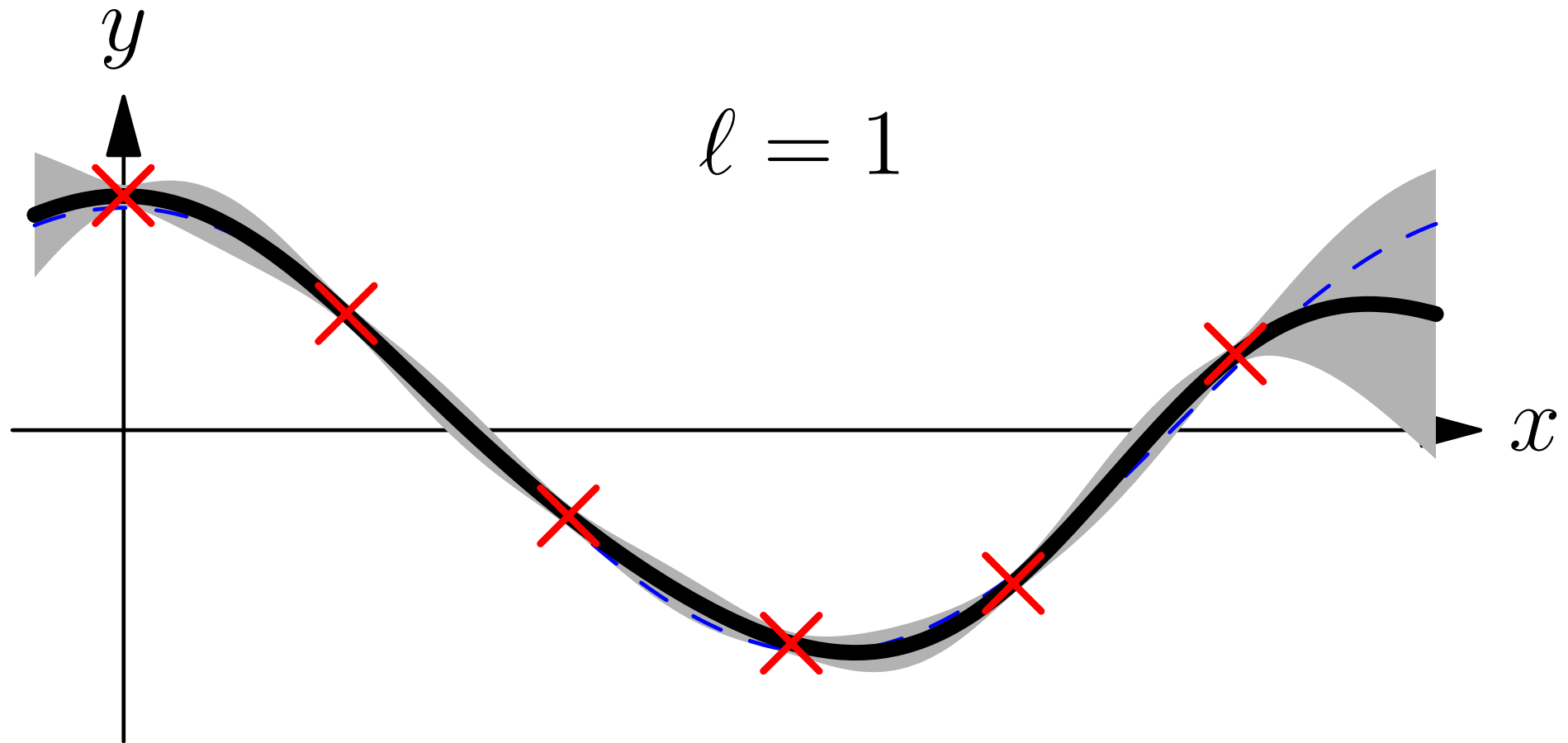
$$K(x, x') = \exp(-(x - x')^2 / (2 \ell^2))$$



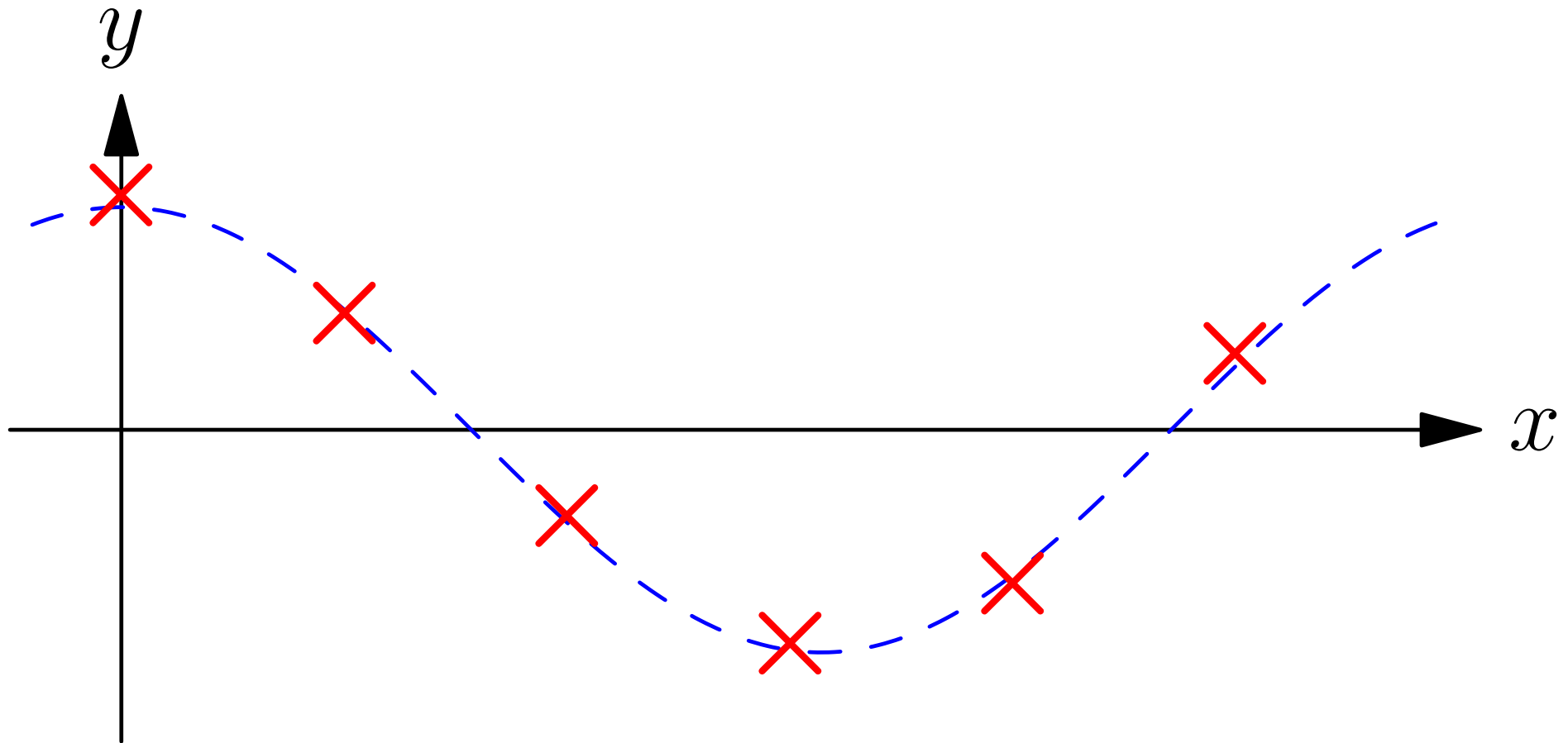
$$K(x, x') = \exp(-(x - x')^2 / (2 \ell^2))$$



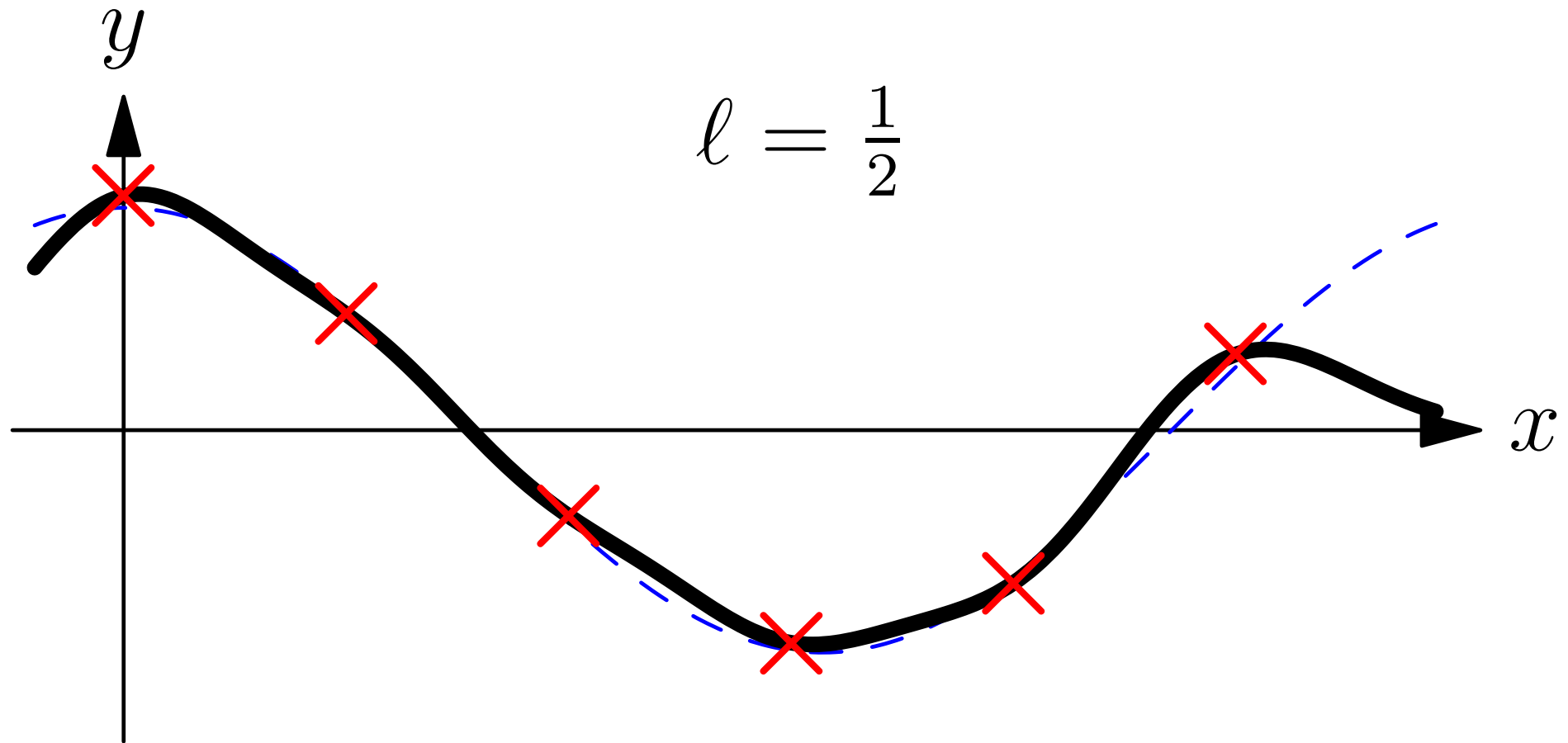
$$K(x, x') = \exp(-(x - x')^2 / (2 \ell^2))$$



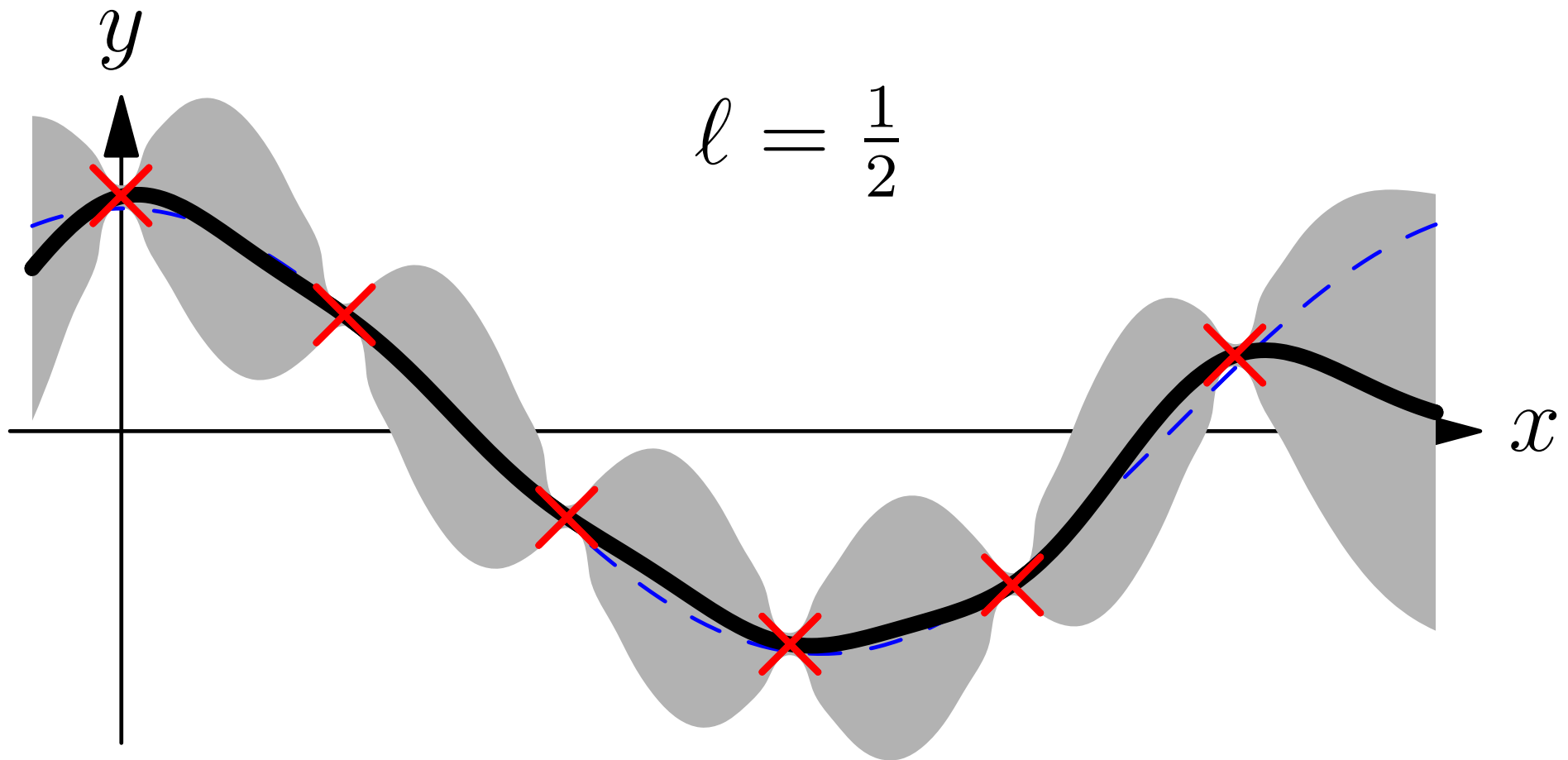
$$K(x, x') = \exp(-(x - x')^2 / (2 \ell^2))$$



$$K(x, x') = \exp(-(x - x')^2 / (2 \ell^2))$$



$$K(x, x') = \exp(-(x - x')^2 / (2 \ell^2))$$





# Model Selection

- In GP we have to choose hyper-parameters
  - ★ Covariance function (length parameter  $\ell$ )
  - ★ Noise in measurements,  $\sigma$
- To do this in a Bayesian framework we can compute the evidence,  $\mathbb{P}(\mathcal{D}|\ell, \sigma)$
- For GP we can compute this in closed form

# Model Selection

- In GP we have to choose hyper-parameters
  - ★ Covariance function (length parameter  $\ell$ )
  - ★ Noise in measurements,  $\sigma$
- To do this in a Bayesian framework we can compute the evidence,  $\mathbb{P}(\mathcal{D}|\ell, \sigma)$
- For GP we can compute this in closed form

# Model Selection

- In GP we have to choose hyper-parameters
  - ★ Covariance function (length parameter  $\ell$ )
  - ★ Noise in measurements,  $\sigma$
- To do this in a Bayesian framework we can compute the evidence,  $\mathbb{P}(\mathcal{D}|\ell, \sigma)$
- For GP we can compute this in closed form

# Using Gaussian Processes

- We are usually not doing 1-d curve fitting but are doing high-dimensional regression (trying to infer a number given a set of (numerical) features)
- You need to give it the right covariance functional, although (hyper-)parameters can be chosen using automatic model selection
- Gaussian processes are often the best method for regression
- Optimal if and only if Gaussian assumption is correct

# Using Gaussian Processes

- We are usually not doing 1-d curve fitting but are doing high-dimensional regression (trying to infer a number given a set of (numerical) features)
- You need to give it the right covariance functional, although (hyper-)parameters can be chosen using automatic model selection
- Gaussian processes are often the best method for regression
- Optimal if and only if Gaussian assumption is correct

# Using Gaussian Processes

- We are usually not doing 1-d curve fitting but are doing high-dimensional regression (trying to infer a number given a set of (numerical) features)
- You need to give it the right covariance functional, although (hyper-)parameters can be chosen using automatic model selection
- Gaussian processes are often the best method for regression
- Optimal if and only if Gaussian assumption is correct

# Using Gaussian Processes

- We are usually not doing 1-d curve fitting but are doing high-dimensional regression (trying to infer a number given a set of (numerical) features)
- You need to give it the right covariance functional, although (hyper-)parameters can be chosen using automatic model selection
- Gaussian processes are often the best method for regression
- Optimal if and only if Gaussian assumption is correct

# Naive Bayes

- We finish with *Naive Bayes* which is unreasonably successful at many text problems
- It is the basis of one of the first big successes of machine learning, namely spam filtering
- It makes the crude approximation that the probability of a document depends only on the words, with all words being independent

$$\mathbb{P}(\text{Doc}|\text{Something}) = \prod_{\text{word} \in \text{Doc}} \mathbb{P}(\text{word}|\text{Something})$$



# Naive Bayes

- We finish with *Naive Bayes* which is unreasonably successful at many text problems
- It is the basis of one of the first big successes of machine learning, namely spam filtering
- It makes the crude approximation that the probability of a document depends only on the words, with all words being independent

$$\mathbb{P}(\text{Doc}|\text{Something}) = \prod_{\text{word} \in \text{Doc}} \mathbb{P}(\text{word}|\text{Something})$$

# Naive Bayes

- We finish with *Naive Bayes* which is unreasonably successful at many text problems
- It is the basis of one of the first big successes of machine learning, namely spam filtering
- It makes the crude approximation that the probability of a document depends only on the words, with all words being independent

$$\mathbb{P}(\text{Doc}|\text{Something}) = \prod_{\text{word} \in \text{Doc}} \mathbb{P}(\text{word}|\text{Something})$$

# Naive Bayes

- We finish with *Naive Bayes* which is unreasonably successful at many text problems
- It is the basis of one of the first big successes of machine learning, namely spam filtering
- It makes the crude approximation that the probability of a document depends only on the words, with all words being independent

$$\mathbb{P}(\text{Doc}|\text{Something}) = \prod_{\text{word} \in \text{Doc}} \mathbb{P}(\text{word}|\text{Something})$$

# Estimating $\mathbb{P}(\text{word}|\text{Something})$

- We can estimate  $\mathbb{P}(\text{word}|\text{Something})$  from the empirical frequency of occurrence of the word
- Also need to estimate  $\mathbb{P}(\text{word}|\neg\text{Something})$
- Often add a pseudo-count (pretend we have seen all words one more time than we have)
- This acts as a regulariser (makes the learning machine less susceptible to the training data)

# Estimating $\mathbb{P}(\text{word}|\text{Something})$

- We can estimate  $\mathbb{P}(\text{word}|\text{Something})$  from the empirical frequency of occurrence of the word
- Also need to estimate  $\mathbb{P}(\text{word}|\neg\text{Something})$
- Often add a pseudo-count (pretend we have seen all words one more time than we have)
- This acts as a regulariser (makes the learning machine less susceptible to the training data)

# Estimating $\mathbb{P}(\text{word}|\text{Something})$

- We can estimate  $\mathbb{P}(\text{word}|\text{Something})$  from the empirical frequency of occurrence of the word
- Also need to estimate  $\mathbb{P}(\text{word}|\neg\text{Something})$
- Often add a pseudo-count (pretend we have seen all words one more time than we have)
- This acts as a regulariser (makes the learning machine less susceptible to the training data)

# Estimating $\mathbb{P}(\text{word}|\text{Something})$

- We can estimate  $\mathbb{P}(\text{word}|\text{Something})$  from the empirical frequency of occurrence of the word
- Also need to estimate  $\mathbb{P}(\text{word}|\neg\text{Something})$
- Often add a pseudo-count (pretend we have seen all words one more time than we have)
- This acts as a regulariser (makes the learning machine less susceptible to the training data)

# Unreasonable Success

- Naive Bayes is relatively simple to implement
- It uses the crude approximation that the order of words is irrelevant!
- Also assume that words are independent {not, realistic}
- However, often gives remarkably good performance



# Unreasonable Success

- Naive Bayes is relatively simple to implement
- It uses the crude approximation that the order of words is irrelevant!
- Also assume that words are independent {not, realistic}
- However, often gives remarkably good performance

# Unreasonable Success

- Naive Bayes is relatively simple to implement
- It uses the crude approximation that the order of words is irrelevant!—same assumption as bag-of-words
- Also assume that words are independent {not, realistic}
- However, often gives remarkably good performance

# Unreasonable Success

- Naive Bayes is relatively simple to implement
- It uses the crude approximation that the order of words is irrelevant!—same assumption as bag-of-words
- Also assume that words are independent {not, realistic}
- However, often gives remarkably good performance

# Unreasonable Success

- Naive Bayes is relatively simple to implement
- It uses the crude approximation that the order of words is irrelevant!—same assumption as bag-of-words
- Also assume that words of independent {not, realistic}
- However, often gives remarkably good performance

# Recap

- SVMs are often best when you have smallish data sets and can encode all the features numerically
- Random Forest and Gradient Boosted Trees are used for large, messy (tabular) data sets
- Bayesian inference is used when you can build a strong model of the likelihood of the data, but it takes work
- Gaussian processes often gives excellent results for regression of numeric inputs
- Deep learning (which Jon covers later) is often the model of choice for images and signals (including text)

# Recap

- SVMs are often best when you have smallish data sets and can encode all the features numerically
- Random Forest and Gradient Boosted Trees are used for large, messy (tabular) data sets
- Bayesian inference is used when you can build a strong model of the likelihood of the data, but it takes work
- Gaussian processes often gives excellent results for regression of numeric inputs
- Deep learning (which Jon covers later) is often the model of choice for images and signals (including text)

# Recap

- SVMs are often best when you have smallish data sets and can encode all the features numerically
- Random Forest and Gradient Boosted Trees are used for large, messy (tabular) data sets
- Bayesian inference is used when you can build a strong model of the likelihood of the data, but it takes work
- Gaussian processes often gives excellent results for regression of numeric inputs
- Deep learning (which Jon covers later) is often the model of choice for images and signals (including text)

# Recap

- SVMs are often best when you have smallish data sets and can encode all the features numerically
- Random Forest and Gradient Boosted Trees are used for large, messy (tabular) data sets
- Bayesian inference is used when you can build a strong model of the likelihood of the data, but it takes work
- Gaussian processes often gives excellent results for regression of numeric inputs
- Deep learning (which Jon covers later) is often the model of choice for images and signals (including text)



# Recap

- SVMs are often best when you have smallish data sets and can encode all the features numerically
- Random Forest and Gradient Boosted Trees are used for large, messy (tabular) data sets
- Bayesian inference is used when you can build a strong model of the likelihood of the data, but it takes work
- Gaussian processes often gives excellent results for regression of numeric inputs
- Deep learning (which Jon covers later) is often the model of choice for images and signals (including text)

# Conclusion

- We started by discussing generalisation performance and the bias-variance dilemma
- We've looked at many of the machine learning methods that are currently state-of-the-art
- They nearly all try to address the bias-variance dilemma one way or another
- When to use them depends on what data you have

# Conclusion

- We started by discussing generalisation performance and the bias-variance dilemma
- We've looked at many of the machine learning methods that are currently state-of-the-art
- They nearly all try to address the bias-variance dilemma one way or another
- When to use them depends on what data you have

# Conclusion

- We started by discussing generalisation performance and the bias-variance dilemma
- We've looked at many of the machine learning methods that are currently state-of-the-art
- They nearly all try to address the bias-variance dilemma one way or another
- When to use them depends on what data you have

# Conclusion

- We started by discussing generalisation performance and the bias-variance dilemma
- We've looked at many of the machine learning methods that are currently state-of-the-art
- They nearly all try to address the bias-variance dilemma one way or another
- When to use them depends on what data you have