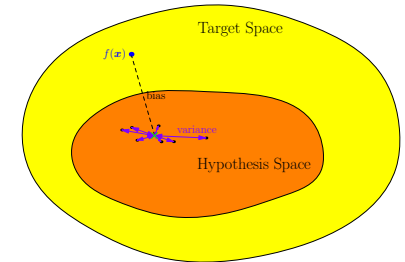


*When ML Works, SVMs, Decision Trees, Ensemble Methods, Bayesian Inference*

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



## What Makes a Good Learning Machine?

- We are going to cover some advanced machine learning techniques
- To understand why these work 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)—can minimise 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

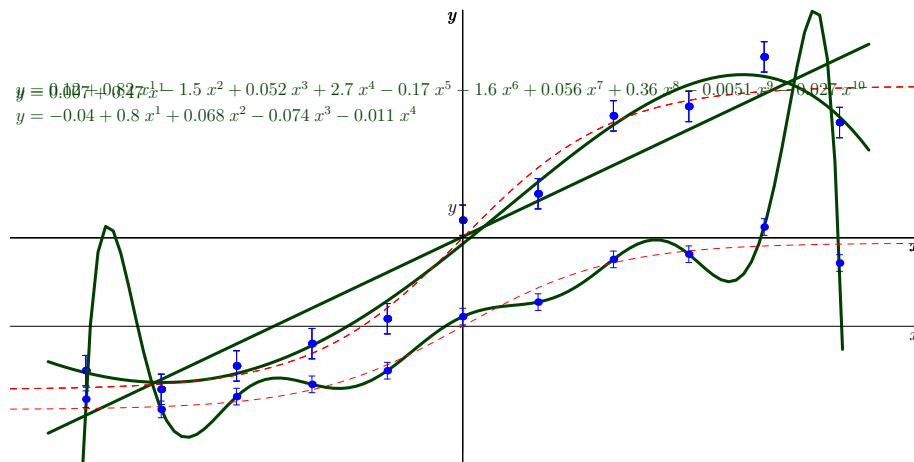
- Suppose we want to learn some function  $f(x)$
- We construct a learning machine that makes a prediction  $\hat{f}(x|w)$ , where  $w$  are weights we want to learn
- We typically choose the weights to minimise a *training error*

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

where  $\mathcal{D}$  is a finite data set of size  $N$ , sampled from the set of all inputs,  $\mathcal{X}$ , according to a probability distribution  $p(x)$  describing where our data is

## Too Simple or Too Complex?

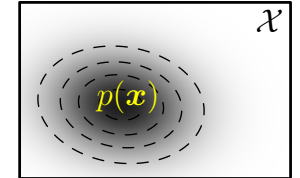
- Fit  $\hat{f}(x, w)$  to data



## Generalisation Error

- We want to minimise the *generalisation error* which in this case we can measure as

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

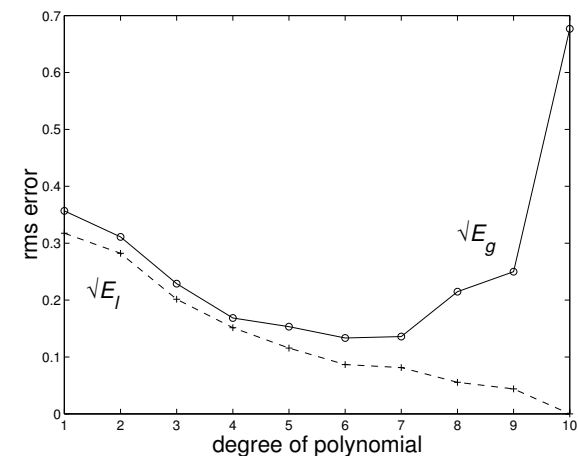


(we can estimate this if we have some examples with known labels  $y_i = f(x_i)$  which we have not trained on)

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

## Measuring Generalisation Error for Regression

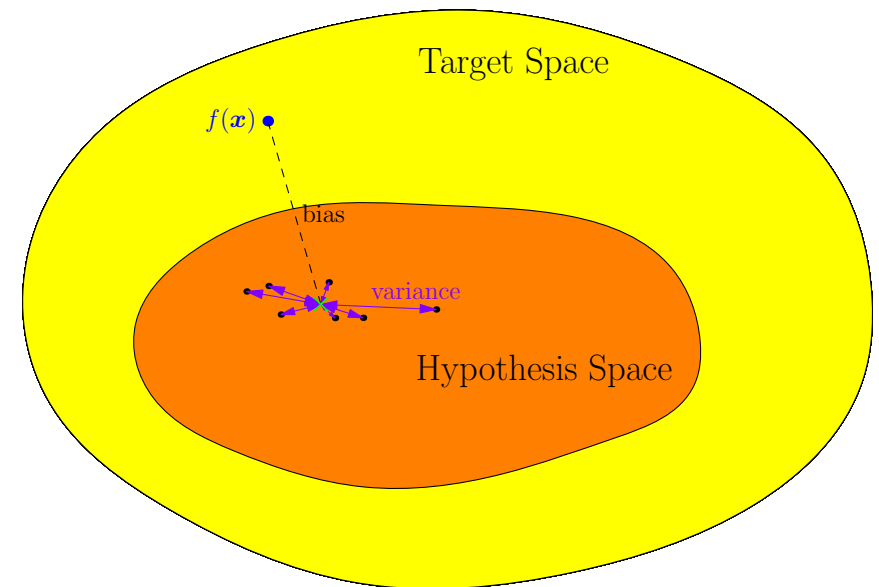
- Consider the regression example. The root mean squared error is



## 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

## Approximation and Estimation Errors



## Mean Machine

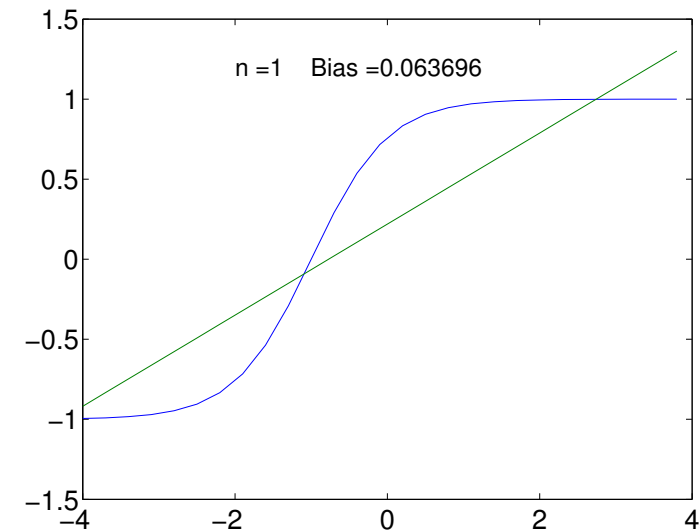
- To help understand generalisation we can consider the mean prediction with respect to machines trained with all data sets of size  $m$

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

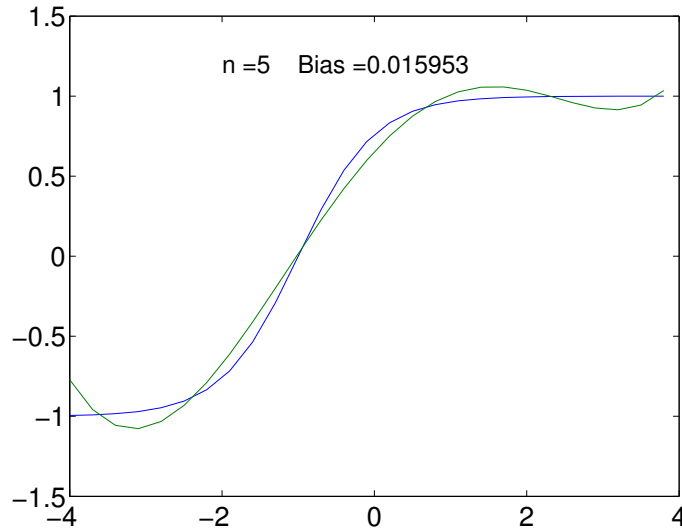
- We can define the **bias** to be generalisation performance of the mean machine

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

## Regression Example $n = 1$



## Regression Example $n = 5$



## Bias and Variance

- We can write the expected generalisation as

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

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

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

## Bias and Variance

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

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

## 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 look at machines that achieve this balance

## Binary Classification Task for You



Class 1



Class 2

## 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. . .
- If we used an infinitely flexible machine we can fit our training data perfectly, but would have no generalisation ability

## Which Category?

- Which category does the following image belong to?



- 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
- (Labelled) data is often expensive to collect so we sometimes have no choice
- Need to control the complexity of our learning machine

## 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

- 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

## 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

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

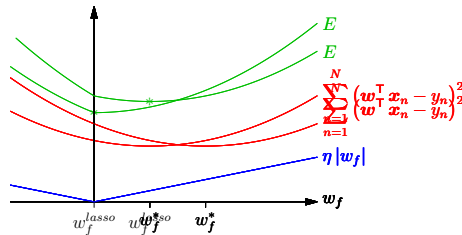
varies rapidly as we change  $x_i$

## Lasso

- We can use other regularisers

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

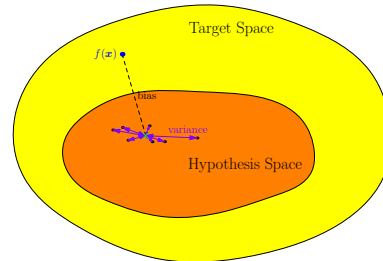
- Spurious features (e.g. shoe size) will give us a small improvement in training error



- Does automatic feature selection

## Outline

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



## Implicit Regularisation

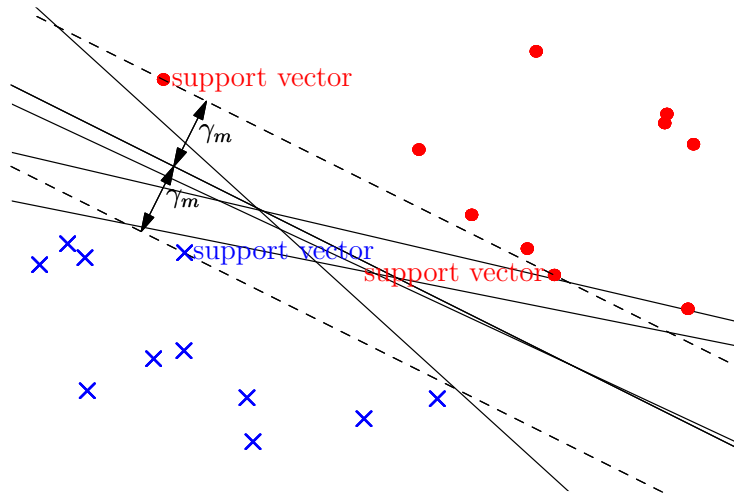
- In the last two examples we added an explicit regularisation term that made the function we learnt simpler
- Some learning machines do this less explicitly
- Some deep learning architectures do subtle averaging
- Sometimes the architecture biases the machine to find a simple solution
- We will see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

## 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



- Finds maximum-margin separating plane

## 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!

## Extended Feature Space

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

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

$$m \gg p$$

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

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

## Kernel Functions

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

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

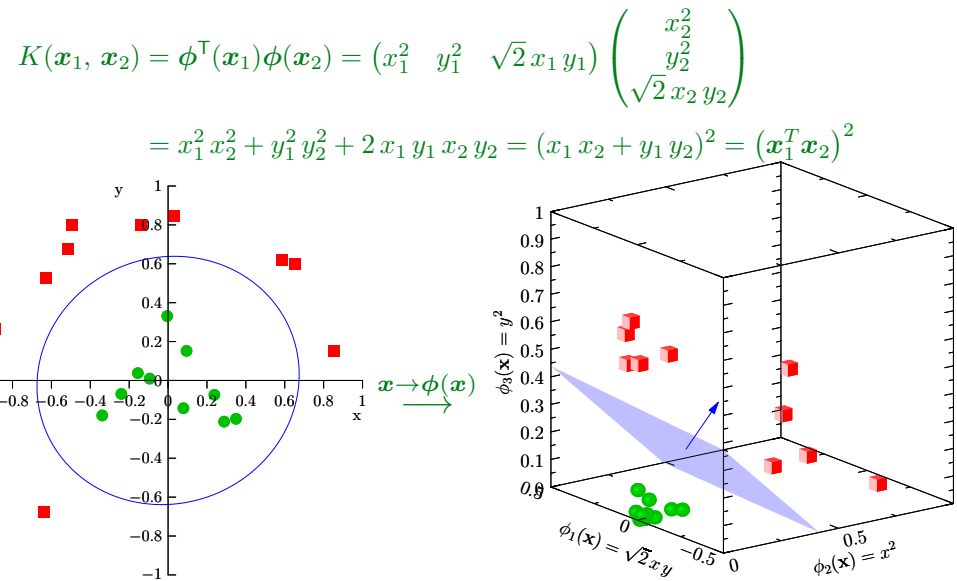
Gaussian (RBF) kernel:  $K(\mathbf{x}_1, \mathbf{x}_2) = e^{-\gamma \|\mathbf{x}_1 - \mathbf{x}_2\|^2}$

- 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



## 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)

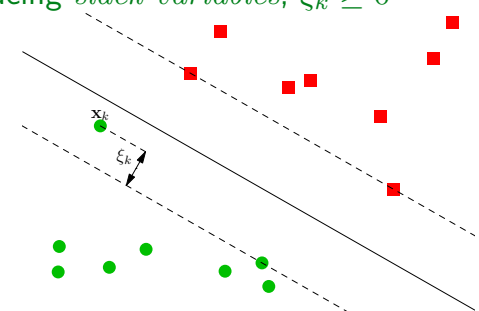
## Computing the Maximum-Margin Hyper-plane

- 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  $\phi^T(\mathbf{x}_i)\phi(\mathbf{x}_j)$
- 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  $pN^3$  where  $N$  is the number of training examples and  $p$  is the number of features

## Soft Margins

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

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



- Minimise  $\frac{\|\mathbf{w}\|_2^2}{2} + C \sum_{k=1}^n \xi_k$  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  $2^{15}$ )

## 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

## Choosing the Right Kernel Function

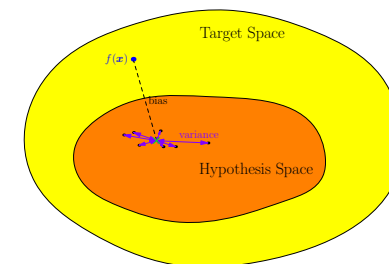
- There are kernels design for particular data types (e.g. string kernels for text or biological sequences)
- For numerical data people tend to look at using no kernel (linear SVM), a radial basis function (Gaussian) kernel or polynomial kernels
- Kernel's often come with parameters, e.g. the popular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

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

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

## Outline

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



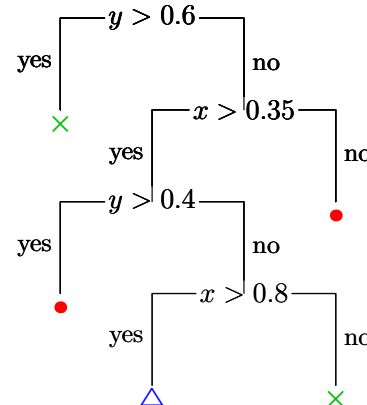
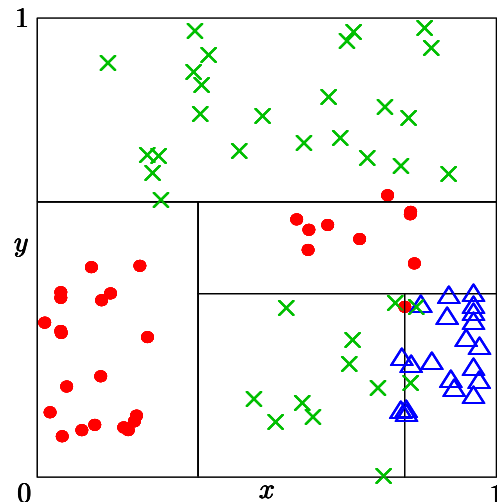
# Removing Variance By Averaging

- 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)

# Ensembling of Decision Trees

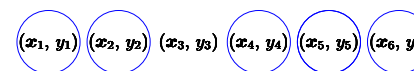
- One set of algorithms where ensembling are common place are decision trees
- These are particularly good for handling messy data
  - ★ categorical data
  - ★ mixture of data types
  - ★ missing data
  - ★ large data sets
  - ★ multiclass
- In many competitions ensembled trees, particularly *random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

## Building Decision Trees



## 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**



$(x_5, y_5)$   $(x_4, y_4)$   $(x_5, y_5)$   $(x_6, y_6)$   $(x_2, y_2)$   $(x_1, y_1)$   
 $(x_1, y_1)$   $(x_5, y_5)$   $(x_3, y_3)$   $(x_2, y_2)$   $(x_5, y_5)$   $(x_2, y_2)$   
 $(x_6, y_6)$   $(x_1, y_1)$   $(x_3, y_3)$   $(x_6, y_6)$   $(x_3, y_3)$   $(x_4, y_4)$   
 $(x_4, y_4)$   $(x_3, y_3)$   $(x_1, y_1)$   $(x_6, y_6)$   $(x_4, y_4)$   $(x_6, y_6)$   
 $(x_3, y_3)$   $(x_2, y_2)$   $(x_3, y_3)$   $(x_4, y_4)$   $(x_6, y_6)$   $(x_2, y_2)$

- 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

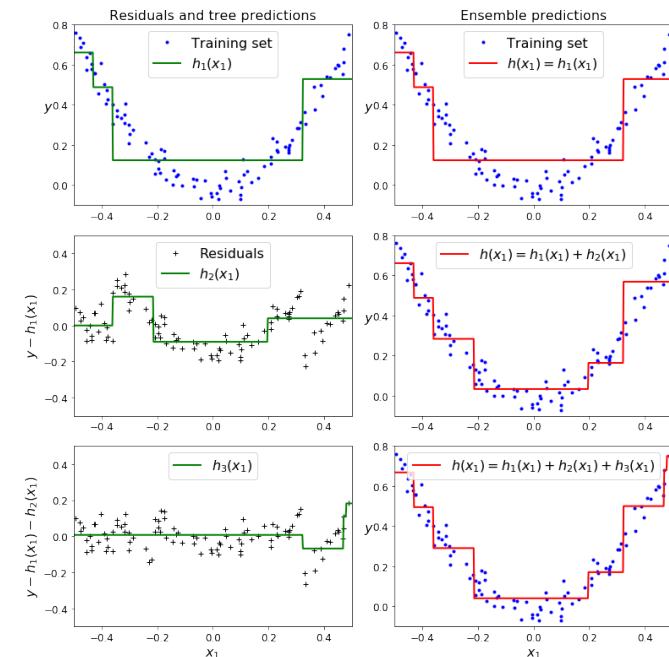
- In boosting we make a **strong learner** by using a weighted sum of **weak learners**

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

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

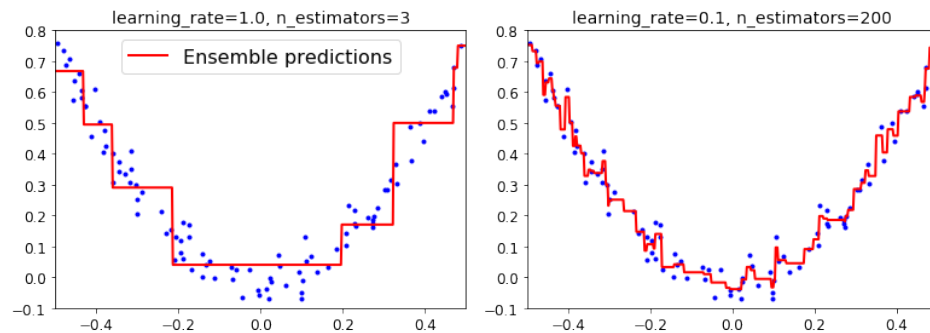
## Shallow Trees

- 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—a classic algorithm for binary problems
  - ★ gradient boosting—now the dominant method, train a classifier on the residual errors
- XGBoost is a very efficient library for gradient boosting on large data sets and often wins competition



## Keep On Going

- We can keep on going



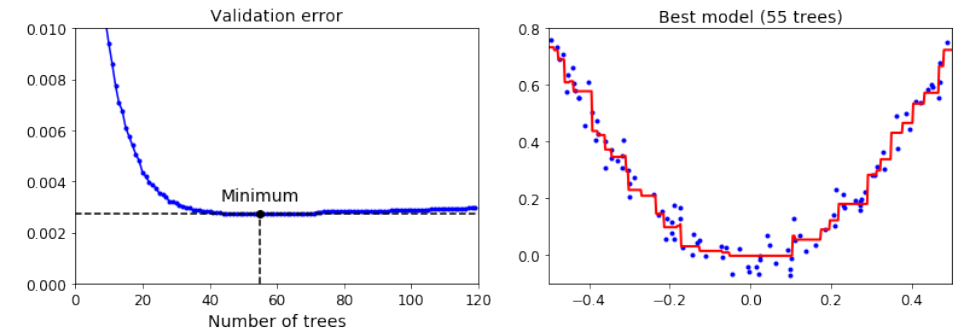
- But we will over-fit eventually

## 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

## Early Stopping

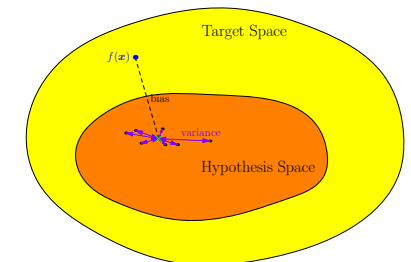
- Like many algorithms we often get better results by early stopping



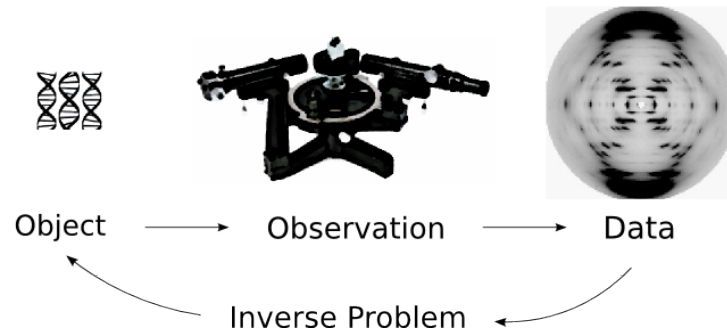
- Use cross-validation against a validation set to decide when to stop

## 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

## Normalisation

- The denominator is a normalisation constant

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

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

## Bayes' Rule

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

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

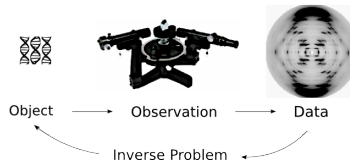
- What we want is to know the probability of the world,  $W$ , given the data,  $D$  we have observed—this is known as the **posteriori** probability
- This depends on the **likelihood** of the data given the world  $\mathbb{P}(D|W)$
- Multiplied by the **prior probability** of the world,  $\mathbb{P}(W)$

## 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

## 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)

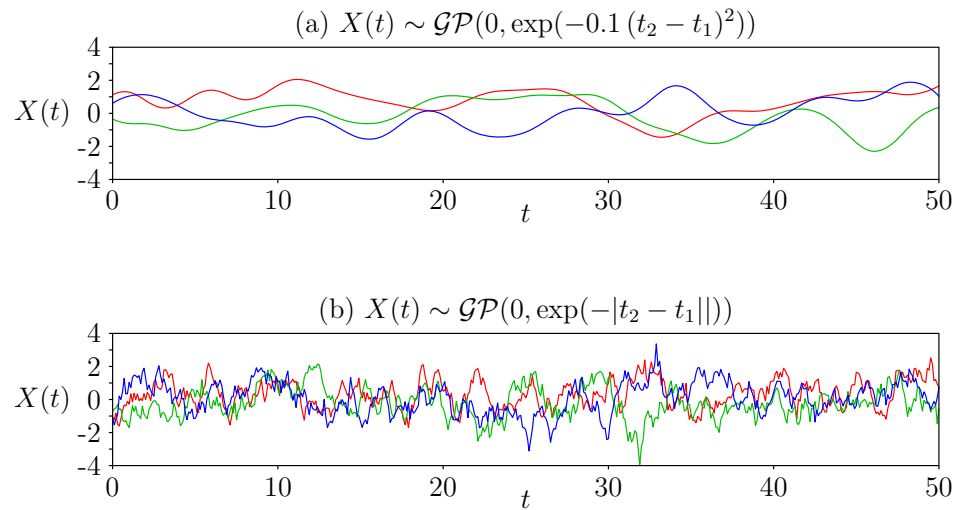
## Disadvantages of Bayesian Inference

- 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

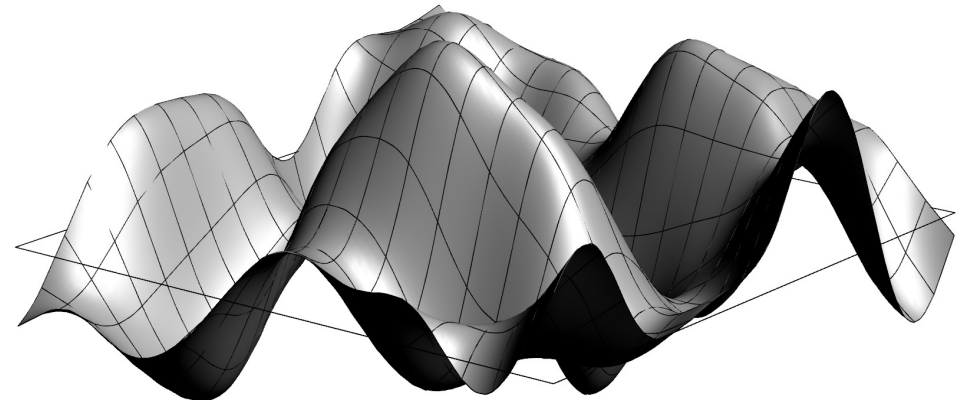
## Gaussian Processes

- 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!

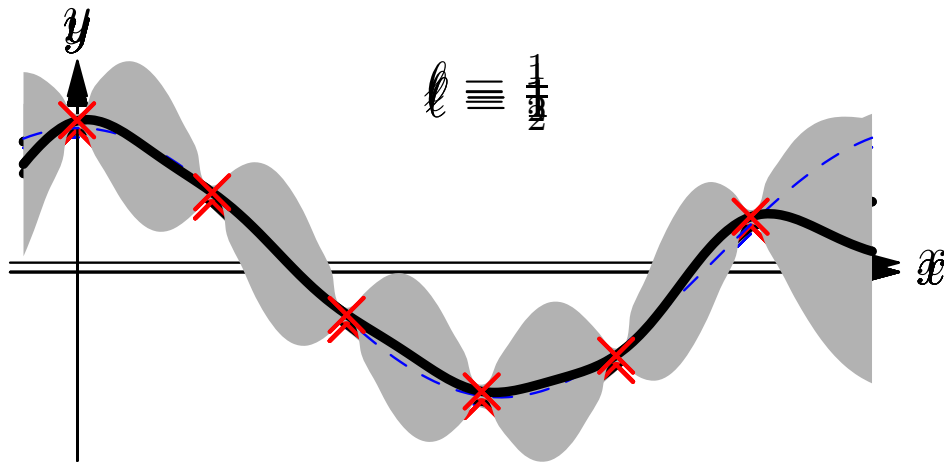
## Gaussian Process Worlds



## 2-D Gaussian Processes



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



## 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



- 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

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

## 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

- 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)

## 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