

Foundations: Machine Learning

* Introduction to Machine Learning

- Outline:
- Introduction to Machine Learning
 - Applications of Machine Learning
 - Useful pre-requisites

* Introduction

* Introduction Part II

* Prerequisites

- Programming experience (esp. Python)
- Statistics
-

* ML is the ROX

* Definition of Machine Learning

Building computational artifacts that learn over time based on experience. Not just building but the mathematics, science, engineering and computing behind the artifacts

* Supervised Learning

Using information from labelled datasets to label new datasets.

Function approximation or function induction

* Induction and Deduction

You assume you have a well behaved function which is consistent with your data and you use this to generalise.

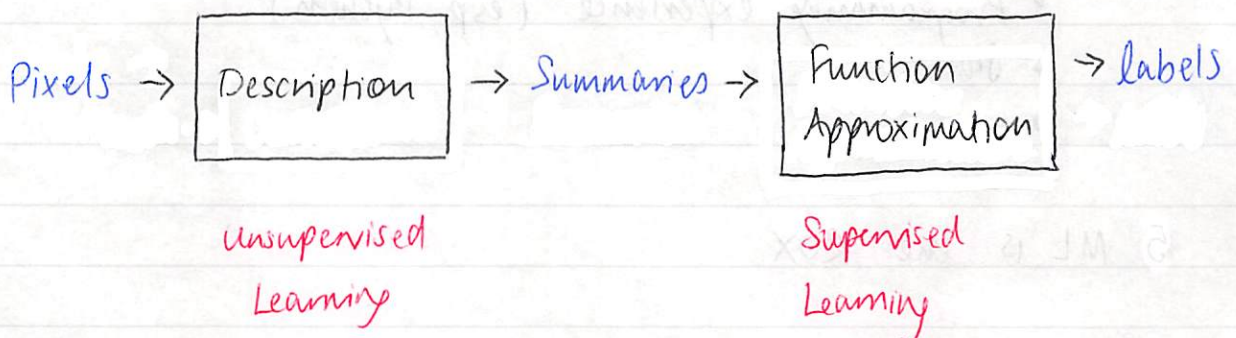
Problem: Inductive bias

Induction:- Example \rightarrow General Rule

Deduction:- General Rule \rightarrow Example

* Unsupervised Learning

Derive structure from the data.



* Reinforcement Learning

Learning from delayed reward

* Comparison of These Parts of Machine Learning

All these problems can be formulated as some kind of optimisation problem

supervised learning	\rightarrow	labels data well
reinforcement learning	\rightarrow	behaviour scores well
unsupervised learning	\rightarrow	cluster scores well

DATA
IS
CENTRAL!

* Machine Learning Basics

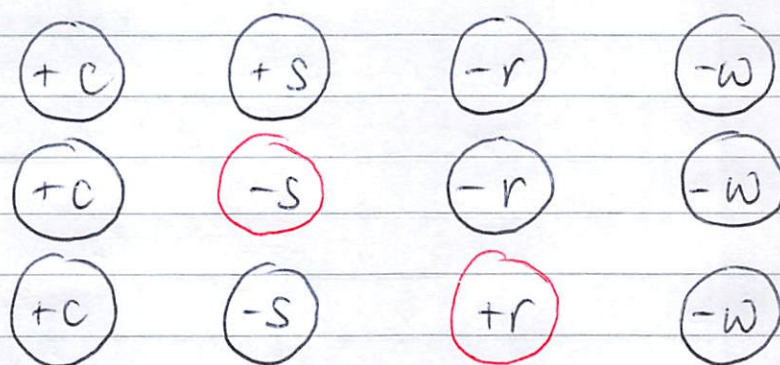
- Outline :
- How modern companies use Machine Learning
 - Stanley at the DARPA Grand Challenge
 - Machine Learning taxonomy
 - Overfitting: Occam's Razor
 - Example problem: Spam Detection
 - Linear methods for supervised learning

* What is Machine Learning

* Gibbs Sampling

... using Markov Chain Monte Carlo (MCMC)

Here we sample one variable at a time



this is consistent

Unit 5: Machine Learning I (Supervised Learning)

- Bayes networks = reason with known models
- Machine learning = learn models from data

Taxonomy of machine learning:

What?	parameters, structure, hidden concepts
What from?	supervised, unsupervised, reinforcement
What for?	prediction, diagnosis, summarisation
How?	passive, active, offline, online <i>the learning agent is... while/not while the data is being generated</i>
Outputs?	classification, regression
Details?	generative, discriminative <i>generates a model distinguish between samples</i>

Supervised Learning

$$\underbrace{\{x_1, x_2, \dots, x_n\}}_{\text{features of the data}} \rightarrow \underbrace{y}_{\text{feature we want to predict}}$$

$$\underbrace{\begin{bmatrix} X_{11} & X_{12} & \dots & X_{1n} \\ X_{21} & X_{22} & \dots & X_{2n} \\ \vdots & \vdots & & \vdots \\ X_{m1} & X_{m2} & \dots & X_{mn} \end{bmatrix}}_{X_m} \rightarrow \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_m \end{bmatrix} \quad \text{data}$$

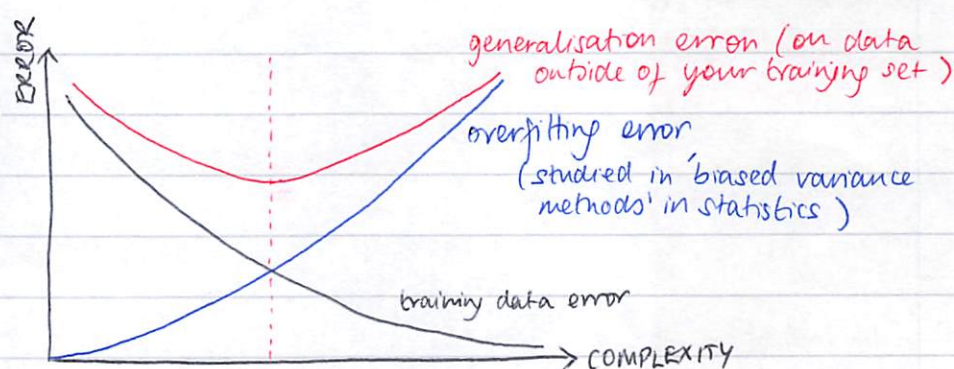
Want to find the function $f(X_m) = Y_m$ so it predicts well Y for X not in our training set

Occam's / Okham's Razor

\Rightarrow Everything else being equal, choose the less complex hypothesis

There tends to be a trade off between fit (to the training data) and complexity (of the hypothesis)

(good) FIT \longleftrightarrow COMPLEXITY (Low)



$$\text{generalisation error} = \text{overfitting error} + \text{training data error}$$

Spam detection - A classification problem

Email
 $\begin{cases} \longrightarrow \text{SPAM (don't show it to the reader)} \\ \longrightarrow \text{HAM (worth passing on to the reader)} \end{cases}$

Y can take one of two values - SPAM/HAM hence a classification problem.

Maximum likelihood

suppose

S S S H H H H H (8 messages)

let

$$p(S) = \pi$$

we want to find π which maximises the likelihood of the training set assuming that each email is drawn independently from an identical distribution

$$p(y_i) = \begin{cases} \pi & \text{if } y_i = S = 1 \\ 1 - \pi & \text{if } y_i = H = 0 \end{cases}$$

$$p(y_i) = \pi^{y_i} (1 - \pi)^{(1 - y_i)}$$

$$\begin{aligned} p(\text{data}) &= \prod_{i=1}^m p(y_i) \\ &= \pi^{\text{count}(y_i=1)} (1 - \pi)^{\text{count}(y_i=0)} \end{aligned}$$

$$\begin{aligned} \ln(p(\text{data})) &= \text{count}(y_i=1) \ln \pi \\ &\quad + \text{count}(y_i=0) \ln(1 - \pi) \end{aligned}$$

$$\frac{d}{d\pi} (\ln(p(\text{data}))) = \frac{\text{count}(y_i=1)}{\pi} - \frac{\text{count}(y_i=0)}{1 - \pi}$$

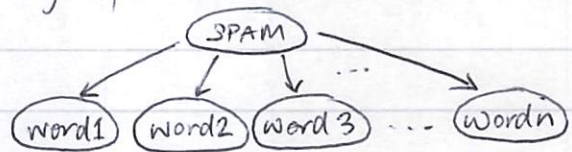
$$\frac{d}{d\pi} (\ln(p(\text{data}))) = 0 \Rightarrow \frac{\text{count}(y_i=1)}{\pi} = \frac{\text{count}(y_i=0)}{1 - \pi}$$

$$\Rightarrow \text{count}(y_i=0) \pi = \text{count}(y_i=1) (1 - \pi)$$

$$\Rightarrow [\text{count}(y_i=0) + \text{count}(y_i=1)] \pi = \text{count}(y_i=1)$$

$$\Rightarrow \pi = \frac{\text{count}(y_i=1)}{\text{count}(y_i=1) + \text{count}(y_i=0)}$$

* Bayes networks examples: Bag of words model
 Naive Bayes Network:
 Detecting SPAM emails:



SPAM	HAM
OFFER IS SECRET	PLAY SPORTS TODAY
CLICK SECRET LINK	WENT PLAY SPORTS
SECRET SPORTS LINK	SECRET SPORTS EVENT
	SPORTS IS TODAY
	SPORTS COSTS MONEY

	Dictionary	SPAM	HAM
1	OFFER		
2	IS		
3	SECRET		
4	CLICK		
5	LINK		
6	SPORT		
7	PLAY		
8	TODAY		
9	WENT		
10	EVENT		
11	COSTS		
12	MONEY		
	Total	9	15

Using Maximum likelihood...

$$P(\text{SPAM}) = \frac{3}{8}$$

$$P(\text{"SECRET"} | \text{SPAM}) = \frac{1}{3}$$

$$P(\text{"SECRET"} | \text{HAM}) = \frac{1}{15}$$

$$P(\text{"IS"} | \text{SPAM}) = \frac{1}{9}$$

$$P(\text{"IS"} | \text{HAM}) = \frac{1}{15}$$

$$P(\text{"TODAY"} | \text{SPAM}) = 0$$

$$P(\text{"TODAY"} | \text{HAM}) = \frac{2}{15}$$

Here it's the probability of the word in a message being "... given ...

Q1 Message $M = \text{"SPORTS"}$

$$P(\text{SPAM} | M) = \frac{\frac{3}{8} \cdot \frac{1}{9}}{\frac{3}{8} \cdot \frac{1}{9} + \frac{5}{8} \cdot \frac{1}{3}} = \frac{3}{3 + 15} = \frac{1}{6}$$

Q2 $M = \text{"SECRET IS SECRET"}$

$$\begin{aligned} P(\text{SPAM} | M) &= \frac{\frac{3}{8} \cdot \frac{1}{3} \cdot \frac{1}{9} \cdot \frac{1}{3}}{\frac{3}{8} \cdot \frac{1}{3} \cdot \frac{1}{9} \cdot \frac{1}{3} + \frac{5}{8} \cdot \frac{1}{15} \cdot \frac{1}{15} \cdot \frac{1}{15}} = \frac{\cancel{3} \times 5 \times \cancel{1} \times 5}{\cancel{3} \times 5 \times \cancel{1} \times 5 + \cancel{5} \times 1 \times \cancel{1} \times 1} \\ &= \frac{25}{26} \end{aligned}$$

Q3 $M = \text{"TODAY IS SECRET"}$

$$P(\text{SPAM} | M) = \frac{\frac{3}{8} \cdot 0 \cdot \frac{1}{9} \cdot \frac{1}{3}}{\frac{3}{8} \cdot 0 \cdot \frac{1}{9} \cdot \frac{1}{3} + \frac{5}{8} \cdot \frac{2}{15} \cdot \frac{1}{15} \cdot \frac{1}{15}} = 0$$

The third example is a poor estimate clearly. Here we are overfitting.

* Laplace Smoothing

Technique to deal with overfitting.

Maximum Likelihood : $p(x) = \frac{\text{count}(x)}{N}$

Laplace Smoothing : $p(x) = \frac{\text{count}(x) + k}{N + k|\mathcal{X}|}$

$\text{count}(x)$ = # of occurrences of this value of x
 $|\mathcal{X}|$ is the number of values which x can have
 k is the smoothing parameter

This is like adding k dummy emails to each class (SPAM and HAM) each containing all the words in the dictionary.

Using Laplace Smoothing for the Bayes networks example...

$$k=1 \Rightarrow P(\text{SPAM}) = \frac{3+1}{8+2} = \frac{4}{10} = \frac{2}{5}$$

$$\Rightarrow P(\text{HAM}) = \frac{3}{5}$$

$$P(\text{"TODAY"} | \text{SPAM}) = \frac{0+1}{9+12} = \frac{1}{21}$$

$$P(\text{"TODAY"} | \text{HAM}) = \frac{2+1}{15+12} = \frac{3}{27} = \frac{1}{9}$$

$$P(\text{"IS"} | \text{SPAM}) = \frac{1+1}{9+12} = \frac{2}{21}$$

$$P(\text{"IS"} | \text{HAM}) = \frac{1+1}{15+12} = \frac{2}{27}$$

$$P(\text{"SECRET"} | \text{SPAM}) = \frac{3+1}{9+12} = \frac{4}{21}$$

$$P(\text{"SECRET"} | \text{HAM}) = \frac{1+1}{15+12} = \frac{2}{27}$$

$$P(\text{SPAM} | \text{"TODAY IS SECRET"})$$

$$= \frac{\cancel{\frac{2}{5}} \cdot \frac{1}{21} \cdot \cancel{\frac{2}{21}} \cdot \frac{4}{21}}{\cancel{\frac{2}{5}} \cdot \frac{1}{21} \cdot \cancel{\frac{2}{21}} \cdot \frac{4}{21} + \frac{3}{5} \cdot \frac{1}{9} \cdot \cancel{\frac{2}{27}} \cdot \cancel{\frac{2}{27}}}$$

$$= \frac{4}{4 + \frac{1}{3} \cdot \frac{1}{27} \cdot \frac{1}{27} \cdot 21 \cdot 21 \cdot 21}$$

$$= \frac{4}{4 + \frac{1}{9} \cdot \frac{1}{9} \cdot 7 \cdot 7 \cdot 7} = \frac{4}{4 + 343/81}$$

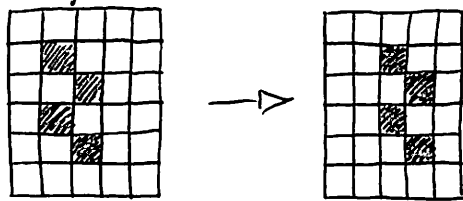
Advanced Spam Filters

- known spamming IP?
- have you emailed the person before?
- have 1K other people received the message?
- is the email header and IP consistent?
- is the email all caps
- do inline URLs point to where they say?
- are you addressed by name?

All these can be used in Naive Bayes

* Digital Recognition

Recognising handwritten digits - could try to use Naive-Bayes on 16×16 pixels but this doesn't deal



well with shifting. To deal with this one could use input smoothing in a different way. One could mix pixel counts with those of the neighbouring pixels so if they are shifted we get similar statistics. Here we convolve the input with a Gaussian variable. This may give better results than just using the raw pixel values themselves.

Actually Naive-Bayes is not a great choice for this problem since it turns out that conditional independence of the pixels in this case is too strong an assumption.

* Overfitting Prevention

We talked about

- Occam's Razor (trade off between fit / prediction)
 - Laplace smoothing
 - Input smoothing
- } How do we choose the smoothing parameter?

Cross-validation:

TRAINING DATA		
TRAIN	CROSS-VALIDATE	TEST
80%	10%	10%

↓
parameters

↓
k

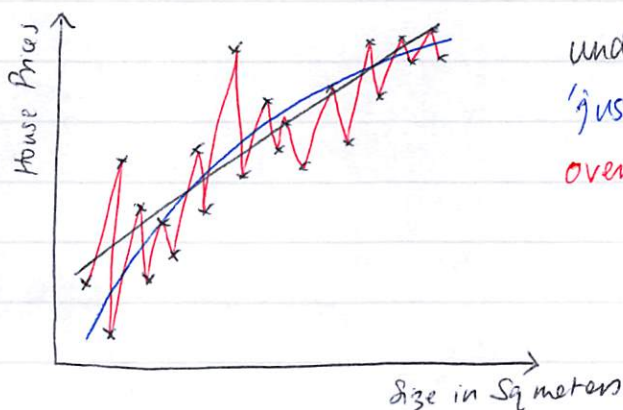
↓
verify performance and report

find k which gives optimal performance on the predictions for the cross-validation data maybe iteratively

← Typical split of your training data.

* Classification vs Regression

y is now continuous rather than discrete
⇒ regression problem



DATA

$$\begin{bmatrix} X_{11} & X_{12} & \dots & X_{1n} & \rightarrow & Y_1 \\ X_{21} & X_{22} & \dots & X_{2n} & \rightarrow & Y_2 \\ \vdots & \vdots & & \vdots & & \vdots \\ X_{m1} & X_{m2} & \dots & X_{mn} & \rightarrow & Y_m \end{bmatrix}$$

$$f(x) = w_1 x + w_0$$

$$f(x) = wX + w_0$$

We want to find f | $f(x) = y$. We do this by minimising the loss function

$$\text{Loss} = \sum_j (y_j - w_1 x_j - w_0)^2 \quad \text{sum square error}$$

$$\text{solution } w^* = \underset{w}{\text{argmin}} \{ \text{Loss} \}$$

* Minimising quadratic loss

$$\min_w \sum_{i=1}^m (y_i - w_1 x_i - w_0)^2 = L$$

$$\frac{\partial L}{\partial w_0} = -2 \sum_{i=1}^m (y_i - w_1 x_i - w_0)$$

$$\frac{\partial L}{\partial w_0} = 0 \Rightarrow m w_0 = \sum_{i=1}^m (y_i - w_1 x_i)$$

$$\Rightarrow w_0 = \frac{1}{m} \left(\sum_{i=1}^m y_i - w_1 \sum_{i=1}^m x_i \right)$$

$$\frac{\partial L}{\partial w_1} = -2 \sum_{i=1}^m (y_i - w_1 x_i - w_0) x_i$$

$$\frac{\partial L}{\partial w_1} = 0 \Rightarrow w_1 \sum_{i=1}^m x_i^2 = \sum_{i=1}^m (y_i - w_0) x_i$$

$$\Rightarrow w_1 \sum_{i=1}^m x_i^2 + w_0 \sum_{i=1}^m x_i = \sum_{i=1}^m x_i y_i$$

$$\Rightarrow w_1 \sum_{i=1}^m x_i^2 + \frac{1}{m} \sum_{i=1}^m (y_i - w_1 x_i) \sum_{i=1}^m x_i = \sum_{i=1}^m x_i y_i$$

$$\Rightarrow w_1 \sum_{i=1}^m x_i^2 + \frac{1}{m} \sum_{i=1}^m y_i \sum_{i=1}^m x_i - \frac{w_1}{m} \left(\sum_{i=1}^m x_i \right)^2 = \sum_{i=1}^m x_i y_i$$

$$\Rightarrow w_1 = \frac{\sum_{i=1}^m x_i y_i - \frac{1}{m} \sum_{i=1}^m x_i \sum_{i=1}^m y_i}{\sum_{i=1}^m x_i^2 - \frac{1}{m} \left(\sum_{i=1}^m x_i \right)^2}$$

Problems with linear regression

- non-linear data
- outliers
- classification problems

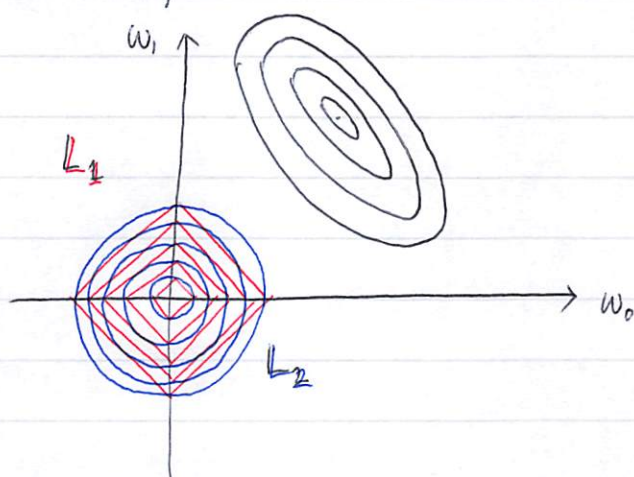
For classification problems we can use logistic regression.

$$\frac{1}{1 + e^{-f(x)}}$$

Regularisation is used for complexity control

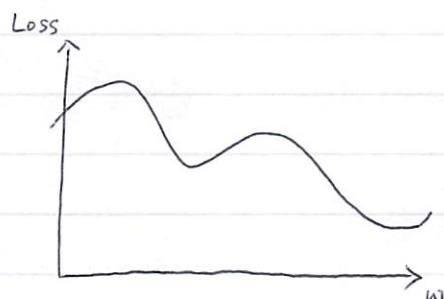
$$\text{Loss} = \text{Loss}(\text{data}) + \text{Loss}(\text{parameters})$$

$$= \sum_j (y_j - w_1 x_j - w_0)^2 + \sum_i |w_i|^p$$



$p=1$: L_1 regularisation
 $p=2$: L_2 regularisation

* Minimisation of more complicated loss functions



Gradient descent

$$w^0$$

$$w^{i+1} \leftarrow w^i - \alpha \frac{\partial L}{\partial w_i}$$

* Gradient descent implementation

$$L = \sum_j (y_j - w_1 x_j - w_0)^2 \rightarrow \min$$

$$\frac{\partial L}{\partial w_1} = -2 \sum_j (y_j - w_1 x_j - w_0) x_j$$

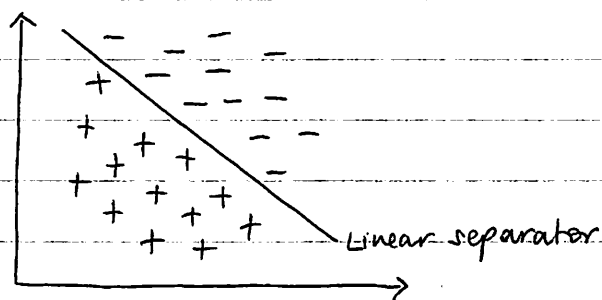
$$\frac{\partial L}{\partial w_0} = -2 \sum_j (y_j - w_1 x_j - w_0)$$

$$w_0 = w_0^0 \text{ and } w_1 = w_1^0$$

$$w_1^i = w_1^{i-1} - \alpha \frac{\partial L}{\partial w_1} (w_1^{i-1})$$

$$w_0^i = w_0^{i-1} - \alpha \frac{\partial L}{\partial w_0} (w_0^{i-1})$$

* Perceptron algorithm



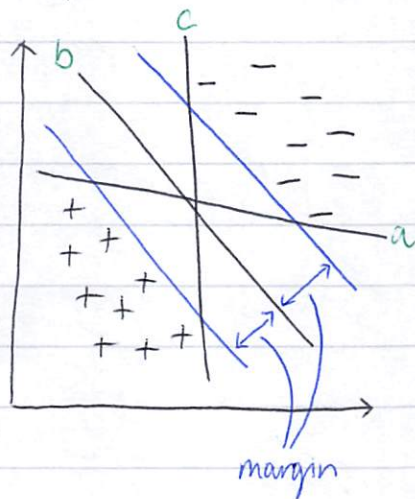
$$f(x) = \begin{cases} 0 & \text{if } w_1 x + w_0 < 0 \\ 1 & \text{if } w_1 x + w_0 \geq 0 \end{cases}$$

linear function

Start with a random guess for w_0 and w_1

$$w_i^k \leftarrow w_i^{k-1} + \alpha (y_i - f(x_i))$$

* Linear separators



separator b is preferable to a and c because of the large margin.

Perceptron only finds a linear separator - not the "best" one.

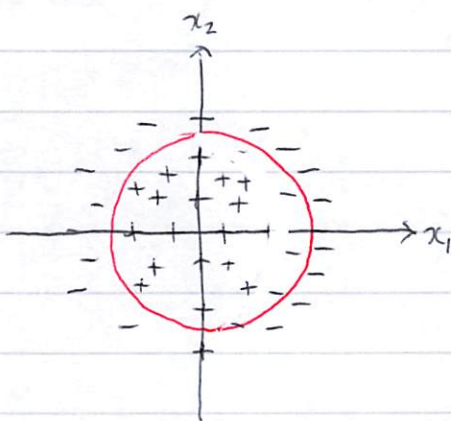
Maximum margin algorithms:

- support vector machines
- boosting

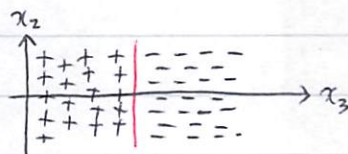
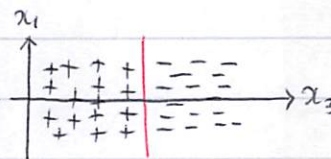
Support vector machines:

These use a "Kernel Trick" to find features which turn complex non-linear decision boundaries into linear ones

Illustration:



$$x_3 = \sqrt{x_1^2 + x_2^2}$$



Linear Methods

- Regression vs Classification
- Exact vs Iterative solutions
- Smoothing
- Non-linear problems