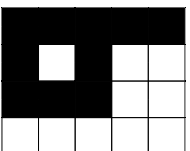
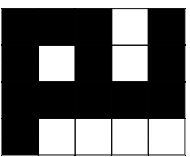


Simple Classification

- Based on images, perform letter classification
- Get lots of images of each letter
- Train a classifier



1.59.302

3.2

Stephen Marsland

Probability

- We are dealing with probabilities
- We make the histogram from our examples
- Joint probability $P(C_i, X_j)$
- Conditional probability $P(X_j | C_i)$

1.59.302

3.4

Stephen Marsland

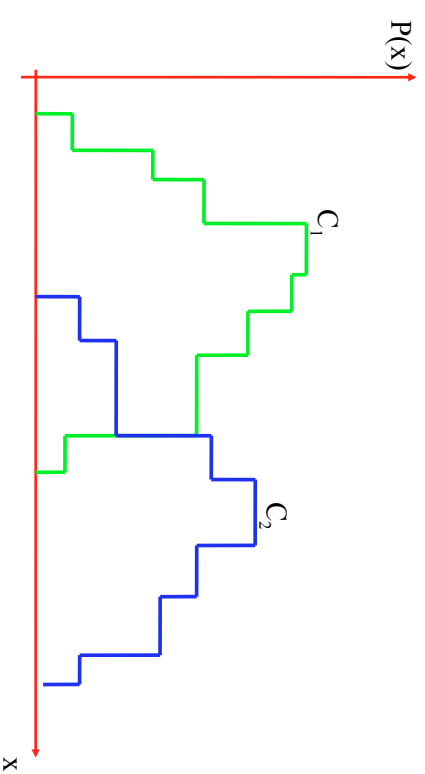
Bayesian Classification

1.59.302

3.1

Stephen Marsland

Feature Histograms



1.59.302

3.3

Stephen Marsland

Bayes' Rule

$$P(C_i|X_j) = \frac{P(X_j|C_i)P(C_i)}{P(X_j)}$$

- Most important equation in machine learning
- Combine things that are easy to find to get useful answers
- Denominator normalises it so probabilities sum to 1

1.59.302

3.6

Stephen Marsland

Classification Process

- Inference
 - ❖ Compute posterior probabilities from data
 - ❖ Make decisions
- ❖ Use posterior probabilities to classify new data
- We do this here by maximising the posterior probability

1.59.302

3.8

Stephen Marsland

Prior Knowledge

- Suppose we know that Class 1 is more likely than Class 2
 - ❖ Distribution of letters in English text
- We should be able to include this information into the classifier: $P(C_j)$

$$\begin{aligned}P(C_i, X_j) &= P(C_i|X_j)P(X_j) \\ P(C_i, X_j) &= P(X_j|C_i)P(C_i)\end{aligned}$$

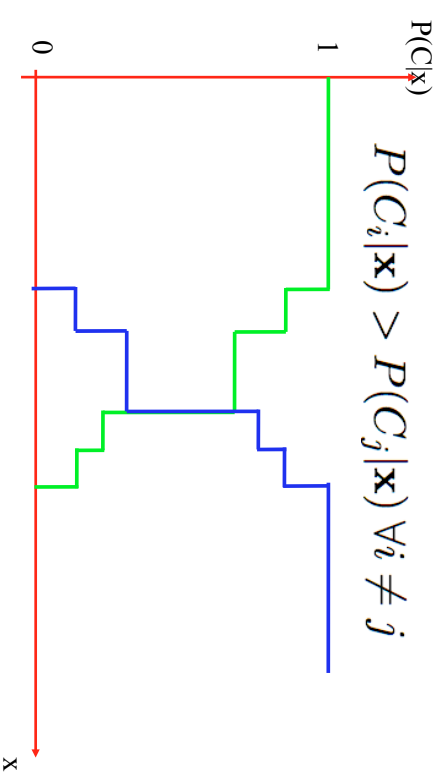
1.59.302

3.5

Stephen Marsland

Posterior Probability

$$P(C_i|\mathbf{x}) > P(C_j|\mathbf{x}) \forall i \neq j$$



1.59.302

3.7

Stephen Marsland

Näive Bayes Classifier

- What if we assume that the features are independent?
- Then can just compute:
$$P(X_j^1 = a_1 | C_i) \times P(X_j^2 = a_2 | C_i) \times \dots \times P(X_j^n = a_n | C_i)$$
$$= \prod_k P(X_j^k = a_k | C_i)$$
- Gross simplification
- Surprisingly effective

1.59.302

3.10

Stephen Marsland

Bayes Classifier

- $$P(C_i | \mathbf{x}) > P(C_j | \mathbf{x}) \forall i \neq j$$
- Need to compute $P(\mathbf{x} | C_j)$
 - Often have high dimensional feature vectors
$$P(X_j^1 = a_1, X_j^2 = a_2, \dots, X_j^n = a_n | C_i)$$
 - Curse of dimensionality applies - need lots and lots of data

1.59.302

3.9

Stephen Marsland

Decision Trees

1.59.302

3.12

Stephen Marsland

What to Maximise

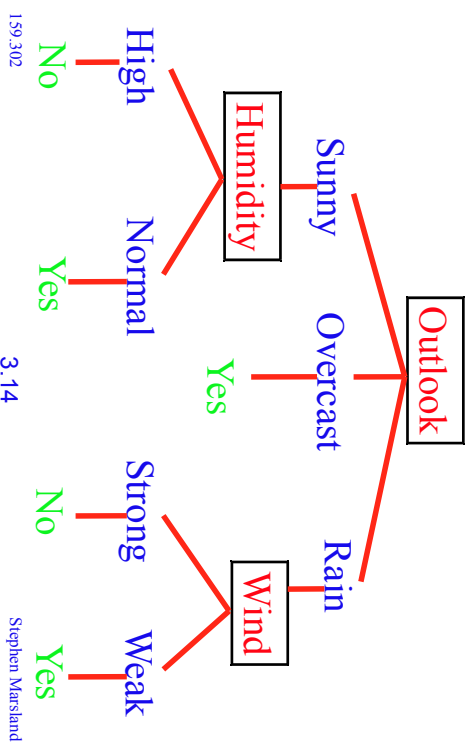
- We maximised posterior probability
- There are other choices
 - ❖ Maximise likelihood
 - ❖ Minimise risk
- Medical data - better to think somebody has a disease than not if unsure
- Loss matrix

1.59.302

3.11

Stephen Marsland

Example



1.59.302

3.14

Stephen Marsland

Entropy

- Tells us how much extra information we get from knowing p_i
- Measures the amount in impurity in the set of features
- Makes sense to pick the features that provides the most information

1.59.302

3.16

Stephen Marsland

Decision Trees

- Split classification down into a series of choices about features in turn
- Lay them out in a tree
- Progress down the tree to the leaves

1.59.302

3.13

Stephen Marsland

Rules and Decision Trees

- Can turn the tree into a set of rules:
 - ❖ (outlook = sunny & humidity = normal) | (outlook = overcast) | (outlook = rain & wind = weak)
- How do we generate the trees?
 - ❖ Need to choose features
 - ❖ Need to choose order of features

1.59.302

3.15

Stephen Marsland

Information Gain

$$\text{Gain}(S, F) = \text{Entropy}(S) - \sum_{f \in \text{values}(F)} \frac{|S_f|}{|S|} \text{Entropy}(S_f)$$

- Choose the feature that provides the highest information gain over all examples
- That is all there is to ID3:
 - ❖ At each stage, pick the feature with the highest information gain

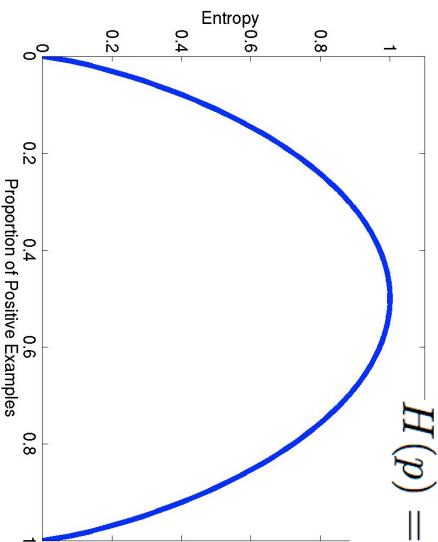
159.302

3.18

Stephen Marsland

Entropy

$$H(p) = \sum_i p_i \log_2 p_i$$

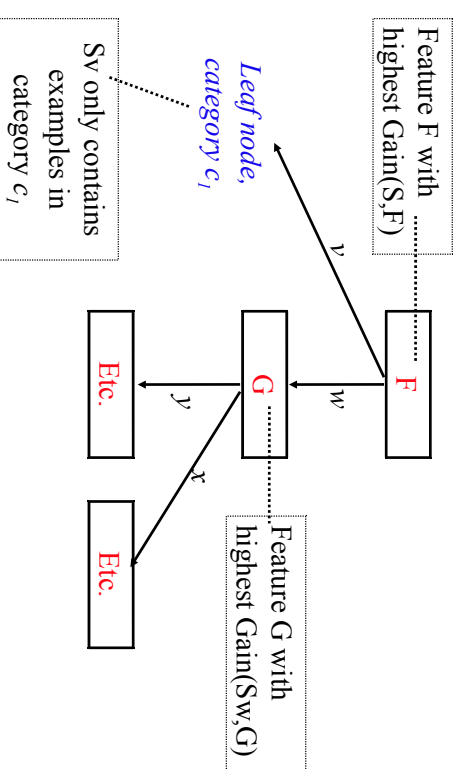


159.302

3.17

Stephen Marsland

ID3



159.302

3.20

Stephen Marsland

ID3 (Quinlan)

- Search over all possible trees
 - ❖ Greedy search - no backtracking
 - ❖ Susceptible to local minima
 - ❖ Uses all features - no pruning
- Can deal with noise
 - ❖ Labels are most common value of examples

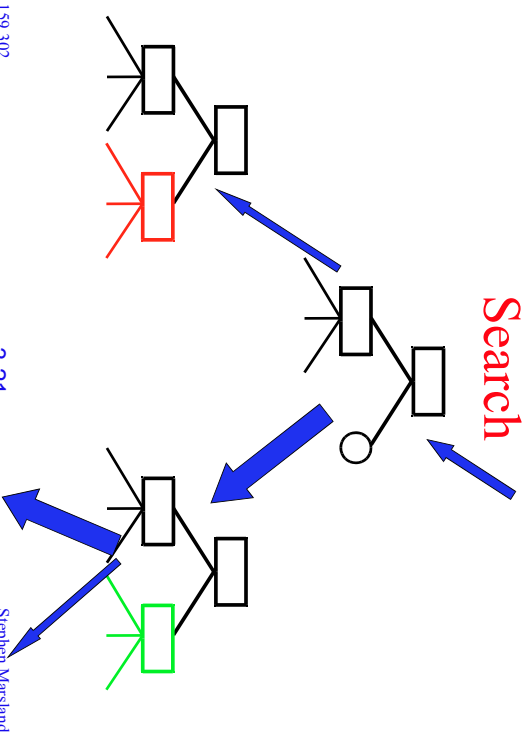
159.302

3.19

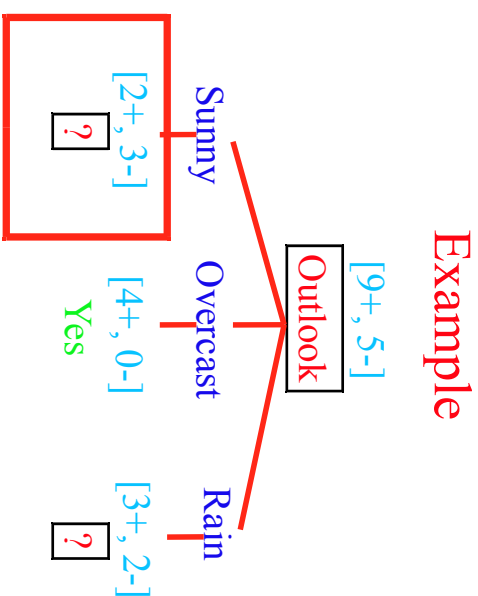
Stephen Marsland

Day	Outlook	Temp	Humid	Wind	Play?
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

1.59.302 3.22 Stephen Marsland



1.59.302 3.21 Stephen Marsland



1.59.302 3.24 Stephen Marsland

Example

- Values(Wind) = Weak, Strong
- $S = [9+, 5-]$
- $S(\text{Weak}) < [6+, 2-]$
- $S(\text{Strong}) < [3+, 3-]$
 - $\text{Gain}(S, \text{Wind}) = \text{Entropy}(S) - (8/14) \text{Entropy}(S(\text{Weak})) - (6/14) \text{Entropy}(S(\text{Strong}))$
 - $= 0.94 - (8/14)0.811 - (6/14)1.00$

1.59.302 3.23 Stephen Marsland

Missing Data

- Suppose that one feature has no value
- Can miss out that node and carry on down the tree, following all paths out of that node
- Can therefore still get a classification
- Virtually impossible with neural networks

1.59.302

3.26

Stephen Marsland

Post-Pruning

- Run over tree
- Prune each node by replacing subtree below with a leaf
- Evaluate error and keep if error same or better

1.59.302

3.28

Stephen Marsland

Inductive Bias

- How does the algorithm generalise from the training examples?
 - ❖ Choose features with highest information gain
 - ❖ Minimise amount of information is left
 - ❖ Bias towards shorter trees
 - ❖ Occam's Razor (KISS)
 - ❖ Put most useful features near root

1.59.302

3.25

Stephen Marsland

C4.5

- Improved version of ID3, also by Quinlan
- Use a validation set to avoid overfitting
 - ❖ Could just stop choosing features (early stopping)
- Better results from post-pruning
 - ❖ Make whole tree
 - ❖ Chop off some parts of tree afterwards

1.59.302

3.27

Stephen Marsland

Rule Post-Pruning

- IF ((outlook = sunny) & (humidity = high))
- THEN playTennis = no
- Remove preconditions:
 - ❖ Consider IF (outlook = sunny)
 - ❖ And IF (humidity = high)
 - ❖ Test accuracy
 - ❖ If one of them is better, try removing both

Rule Post-Pruning

- Turn tree into set of if-then rules
- Remove preconditions from each rule in turn, and check accuracy
- Sort rules according to accuracy
- Rules are easy to read

Day	Outlook	Temp	Humid	Wind	play?
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

Comparison Between Decision Trees and Naïve Bayes Classifier

Test Case

- Outlook = Sunny
- Temperature = Cool
- Humidity = High
- Wind = Strong

1.59.302

3.34

Stephen Marsland

Naïve Bayes

- Yes: 0.0053
- No: 0.0206
- So solution is no
- Conditional probability is:

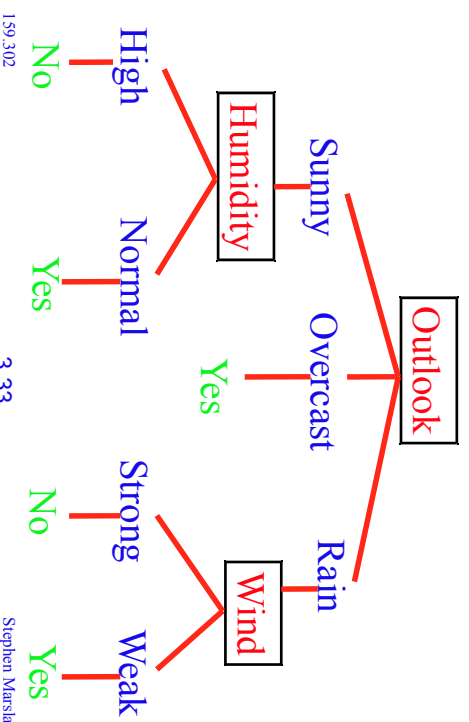
$$\frac{0.0206}{0.0206 + 0.0053} = 0.795$$

1.59.302

3.36

Stephen Marsland

ID3 Decision Tree



1.59.302

3.33

Stephen Marsland

Naïve Bayes

- $P(\text{yes}) * P(\text{Outlook}=\text{Sunny} | \text{yes}) * P(\text{Temperature}=\text{Cool} | \text{yes}) * P(\text{Humidity}=\text{High} | \text{yes}) * P(\text{Wind}=\text{Strong} | \text{yes})$
- Similar for no
- Count all the probabilities from the table

1.59.302

3.35

Stephen Marsland

