

Because $N \rightarrow$ is huge
"big data"

Stochastic Gradient Descent

Machine Learning – CSE446

Carlos Guestrin

University of Washington

April 19, 2013

©Carlos Guestrin 2005-2013

1

Logistic Regression

Logistic function
(or Sigmoid): $\frac{1}{1 + \exp(-z)}$

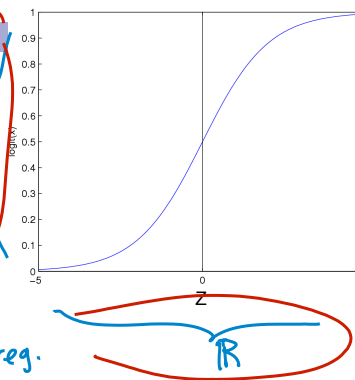
Learn $P(Y|\mathbf{X})$ directly

- Assume a particular functional form for link function
- Sigmoid applied to a linear function of the input features:

$$P(Y = 0|X, W) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

choice

z : linear just like in reg.



$w_0 + \sum_i w_i X_i$ ← not bounded, could be neg.

after logistic fcn, output is in $[0, 1]$

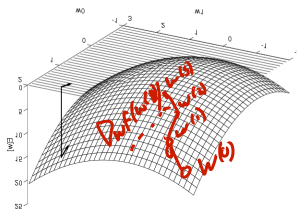
Features can be discrete or continuous!

©Carlos Guestrin 2005-2013

2

Optimizing concave function – Gradient ascent

- Conditional likelihood for Logistic Regression is concave. Find optimum with gradient ascent



Gradient: $\nabla_{\mathbf{w}} l(\mathbf{w}) = \left[\frac{\partial l(\mathbf{w})}{\partial w_0}, \dots, \frac{\partial l(\mathbf{w})}{\partial w_n} \right]'$

Step size, $\eta > 0$

Update rule: $\Delta \mathbf{w} = \eta \nabla_{\mathbf{w}} l(\mathbf{w})$

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \frac{\partial l(\mathbf{w}^{(t)})}{\partial w_i}$$

But how?
 η will be constant

- Gradient ascent is simplest of optimization approaches

- e.g., Conjugate gradient ascent can be much better

Often, especially in proof, η gets smaller with iterations
e.g. $\eta_t = \frac{\alpha}{t}$ $\alpha \leftarrow \text{constant}$

©Carlos Guestrin 2005-2013

3

Gradient Ascent for LR

Start from some $\mathbf{w}^{(0)}$ e.g. $\mathbf{0}$

revisit soon

Gradient ascent algorithm: iterate until change $< \epsilon$

$$w_0^{(t+1)} \leftarrow w_0^{(t)} + \eta \sum_j [y^j - \hat{P}(Y^j = 1 | \mathbf{x}^j, \mathbf{w}^{(t)})]$$

step size

For $i=1, \dots, k$,

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \sum_j x_i^j [y^j - \hat{P}(Y^j = 1 | \mathbf{x}^j, \mathbf{w}^{(t)})]$$

repeat

©Carlos Guestrin 2005-2013

4

The Cost, The Cost!!! Think about the cost...

$k, N \leftarrow$ terms of

- What's the cost of a gradient update step for LR???

expensive Approx to expected gradient

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \left\{ -\lambda w_i^{(t)} + \sum_{j=1}^N x_i^j [y^j - \hat{P}(Y^j = 1 | \mathbf{x}^j, \mathbf{w}^{(t)})] \right\}$$

for i
naively $O(Nk^2)$
 but cache \hat{P} (same for all i)
 $O(Nk) \in \dots$ *if N is huge very slow per little gradient step*

©Carlos Guestrin 2005-2013

5

Learning Problems as Expectations

- Minimizing loss in training data:

- Given dataset: x_1, x_2, \dots, x_N *$x_i \text{ iid } p(x)$*
 - Sampled iid from some distribution $p(x)$ on features:
- Loss function, e.g., hinge loss, logistic loss, ... *squared loss*
- We often minimize loss in training data:

loss over training data

$$\ell_D(\mathbf{w}) = \frac{1}{N} \sum_{j=1}^N \ell(\mathbf{w}, \mathbf{x}^j)$$

$\ln p(y^j | \mathbf{x}^j, \mathbf{w}) - \lambda \|\mathbf{w}\|^2$

- However, we should really minimize expected loss on all data:

Surrogate (empirical loss)

$$\ell(\mathbf{w}) = E_{\mathbf{x}} [\ell(\mathbf{w}, \mathbf{x})] = \int p(\mathbf{x}) \ell(\mathbf{w}, \mathbf{x}) d\mathbf{x}$$

expected loss

- So, we are approximating the integral by the average on the training data

©Carlos Guestrin 2005-2013

6

Gradient ascent in Terms of Expectations

- “True” objective function:

$$\ell(\mathbf{w}) = E_{\mathbf{x}} [\ell(\mathbf{w}, \mathbf{x})] = \int p(\mathbf{x}) \ell(\mathbf{w}, \mathbf{x}) d\mathbf{x}$$

- Taking the gradient:

$$\nabla_{\mathbf{w}} \ell(\mathbf{w}) = \nabla_{\mathbf{w}} [E_{\mathbf{x}} [\ell(\mathbf{w}, \mathbf{x})]] = E_{\mathbf{x}} [\nabla_{\mathbf{w}} \ell(\mathbf{w}, \mathbf{x})]$$

- “True” gradient ascent rule:

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + \eta E_{\mathbf{x}} [\nabla_{\mathbf{w}} \ell(\mathbf{w}, \mathbf{x})]$$

- How do we estimate expected gradient?

cannot compute this
approx. with sample data

©Carlos Guestrin 2005-2013

7

SGD: Stochastic Gradient Ascent (or Descent)

- “True” gradient: $\nabla \ell(\mathbf{w}) = E_{\mathbf{x}} [\nabla \ell(\mathbf{w}, \mathbf{x})]$

- Sample based approximation: take iid samples \mathbf{x}^j

$$\nabla_{\mathbf{w}} \ell(\mathbf{w}) \approx \frac{1}{N} \sum_{j=1}^N \nabla_{\mathbf{w}} \ell(\mathbf{w}, \mathbf{x}^j)$$

$\tilde{\nabla} \ell$

- What if we estimate gradient with just one sample???

- Unbiased estimate of gradient $E[\tilde{\nabla} \ell] = \nabla_{\mathbf{w}} \ell$
- Very noisy!
- Called stochastic gradient ascent (or descent)
 - Among many other names
- VERY useful in practice!!!

©Carlos Guestrin 2005-2013

8

Stochastic Gradient Ascent for Logistic Regression

- Logistic loss as a stochastic function:

$$E_{\mathbf{x}} [\ell(\mathbf{w}, \mathbf{x})] = E_{\mathbf{x}} [\ln P(y|\mathbf{x}, \mathbf{w}) - \lambda \|\mathbf{w}\|_2^2]$$

- Batch gradient ascent updates:

$$O(KN) \rightarrow w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \left\{ -\lambda w_i^{(t)} + \frac{1}{N} \sum_{j=1}^N x_i^{(j)} [y^{(j)} - P(Y=1|\mathbf{x}^{(j)}, \mathbf{w}^{(t)})] \right\}$$

- Stochastic gradient ascent updates:

- Online setting:

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta_t \left\{ -\lambda w_i^{(t)} + x_i^{(t)} [y^{(t)} - P(Y=1|\mathbf{x}^{(t)}, \mathbf{w}^{(t)})] \right\}$$

i: feature or coordinate
take one data point at a time: $x^{(t)}$
same data point $x_i^{(t)}$
gradient update for just $x^{(t)}$

©Carlos Guestrin 2005-2013

9

Stochastic Gradient Ascent: general case

- Given a stochastic function of parameters: $f(\mathbf{w}) = E_{\mathbf{x}} [f(\mathbf{w}, \mathbf{x})]$

- Want to find maximum

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} f(\mathbf{w}) \equiv \arg \max_{\mathbf{w}} E_{\mathbf{x}} [f(\mathbf{w}, \mathbf{x})]$$

- Start from $\mathbf{w}^{(0)}$ e.g. $\mathbf{w}^{(0)} = \mathbf{0}$

- Repeat until convergence:

- Get a sample data point \mathbf{x}^t

- Update parameters:

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + \eta \nabla f_{\mathbf{w}}(\mathbf{w}, \mathbf{x}^t)$$

- Works on the online learning setting!

- Complexity of each gradient step is constant in number of examples!

- In general, step size changes with iterations

η decrease with t , e.g. $\eta_t = \frac{\alpha}{t}$ for $\alpha > 0$ constant
in your homework use constant η

©Carlos Guestrin 2005-2013

10

What you should know...

- Classification: predict discrete classes rather than real values
- Logistic regression model: Linear model
 - Logistic function maps real values to $[0,1]$
- Optimize conditional likelihood
- Gradient computation
- Overfitting
- Regularization
- Regularized optimization
- Cost of gradient step is high, use stochastic gradient descent

©Carlos Guestrin 2005-2013

11

Decision Trees

Machine Learning – CSE446

Carlos Guestrin

University of Washington

April 19, 2013

©Carlos Guestrin 2005-2013

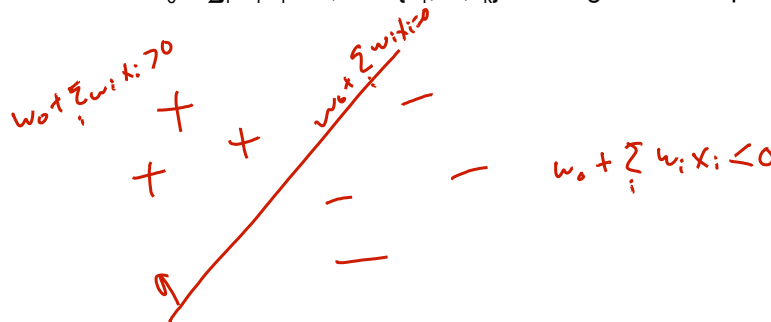
12

Linear separability

- A dataset is **linearly separable** iff there exists a **separating hyperplane**:

□ Exists \mathbf{w} , such that:

- $w_0 + \sum_i w_i x_i > 0$; if $\mathbf{x}=\{x_1, \dots, x_k\}$ is a positive example
- $w_0 + \sum_i w_i x_i < 0$; if $\mathbf{x}=\{x_1, \dots, x_k\}$ is a negative example



©Carlos Guestrin 2005-2013

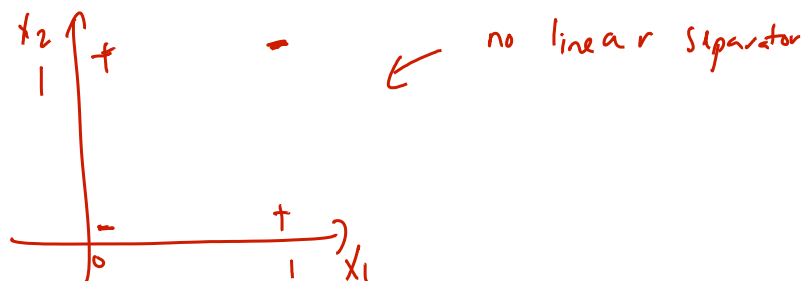
13

Not linearly separable data

- Some datasets are **not linearly separable**!

$$x_1 \oplus x_2$$

$$x_1 \wedge \neg x_2 \vee \neg x_1 \wedge x_2$$

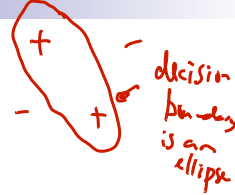


©Carlos Guestrin 2005-2013

14

Addressing non-linearly separable data – Option 1, non-linear features

- Choose non-linear features, e.g.,
 - Typical linear features: $w_0 + \sum_i w_i x_i$
 - Example of non-linear features:
 - Degree 2 polynomials, $w_0 + \sum_i w_i x_i + \sum_{ij} w_{ij} x_i x_j$
- Classifier $h_w(\mathbf{x})$ still linear in parameters \mathbf{w}
 - As easy to learn \leftarrow convex/concave problem
 - Data is linearly separable in higher dimensional spaces
 - More discussion later this quarter



©Carlos Guestrin 2005-2013

15

Addressing non-linearly separable data – Option 2, non-linear classifier

- Choose a classifier $h_w(\mathbf{x})$ that is non-linear in parameters \mathbf{w} , e.g.,
 - Decision trees, boosting, nearest neighbor, neural networks...
- More general than linear classifiers
- But, can often be harder to learn (non-convex/concave optimization required)
- But, but, often very useful
- (BTW. Later this quarter, we'll see that these options are not that different)

©Carlos Guestrin 2005-2013

16

A small dataset: Miles Per Gallon

Suppose we want to predict MPG

$X \rightarrow \text{MPG} \{ \text{Good, Bad} \}$

mpg	cylinders	displacement	horsepower	weight	acceleration	modelyear	maker
good	4	low	low	low	high	75to78	asia
bad	6	medium	medium	medium	medium	70to74	america
bad	4	medium	medium	medium	low	75to78	europa
bad	8	high	high	high	low	70to74	america
bad	6	medium	medium	medium	medium	70to74	america
bad	4	low	medium	low	medium	70to74	asia
bad	4	low	medium	low	low	70to74	asia
bad	8	high	high	high	low	75to78	america
:	:	:	:	:	:	:	:
:	:	:	:	:	:	:	:
:	:	:	:	:	:	:	:
bad	8	high	high	high	low	70to74	america
good	8	high	medium	high	high	79to83	america
bad	8	high	high	high	low	75to78	america
good	4	low	low	low	low	79to83	america
bad	6	medium	medium	medium	high	75to78	america
good	4	medium	low	low	low	79to83	america
good	4	low	low	medium	high	79to83	america
bad	8	high	high	high	low	70to74	america
good	4	low	medium	low	medium	75to78	europa
bad	5	medium	medium	medium	medium	75to78	europa

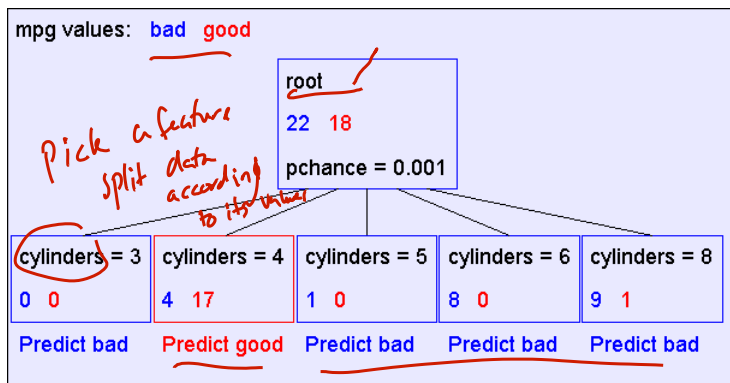
40 training examples

From the UCI repository (thanks to Ross Quinlan)

©Carlos Guestrin 2005-2013

17

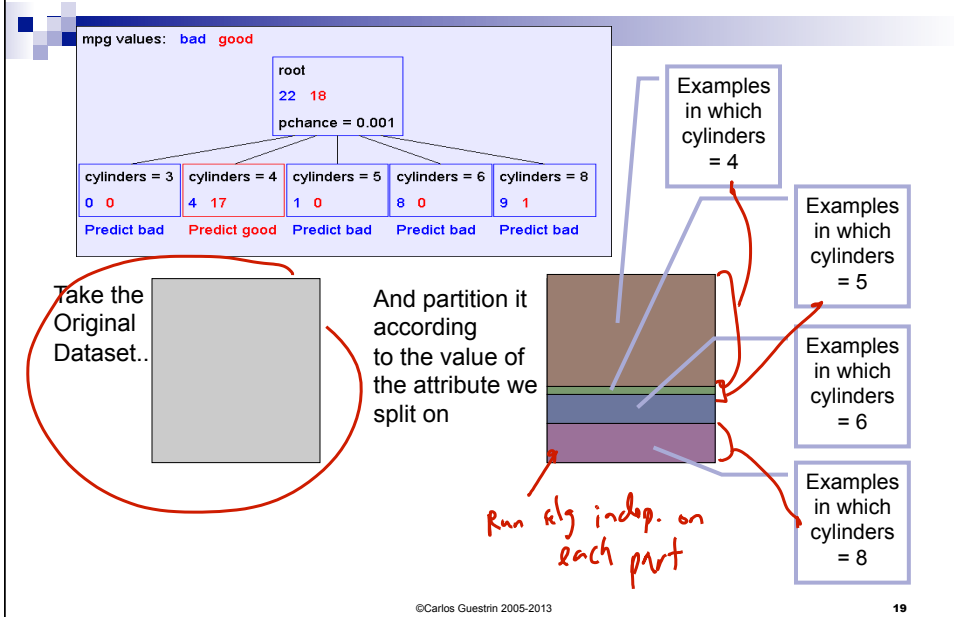
A Decision Stump



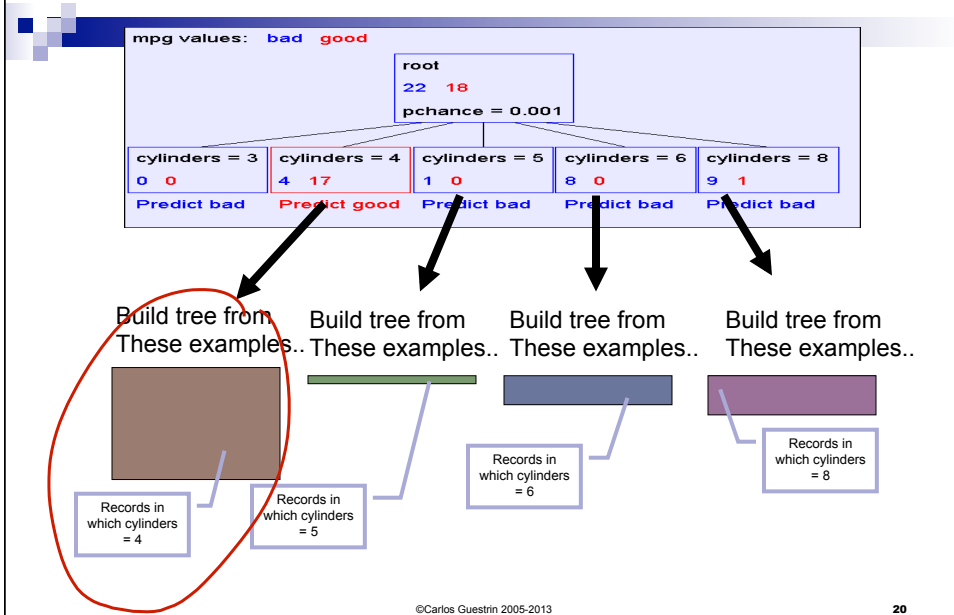
©Carlos Guestrin 2005-2013

18

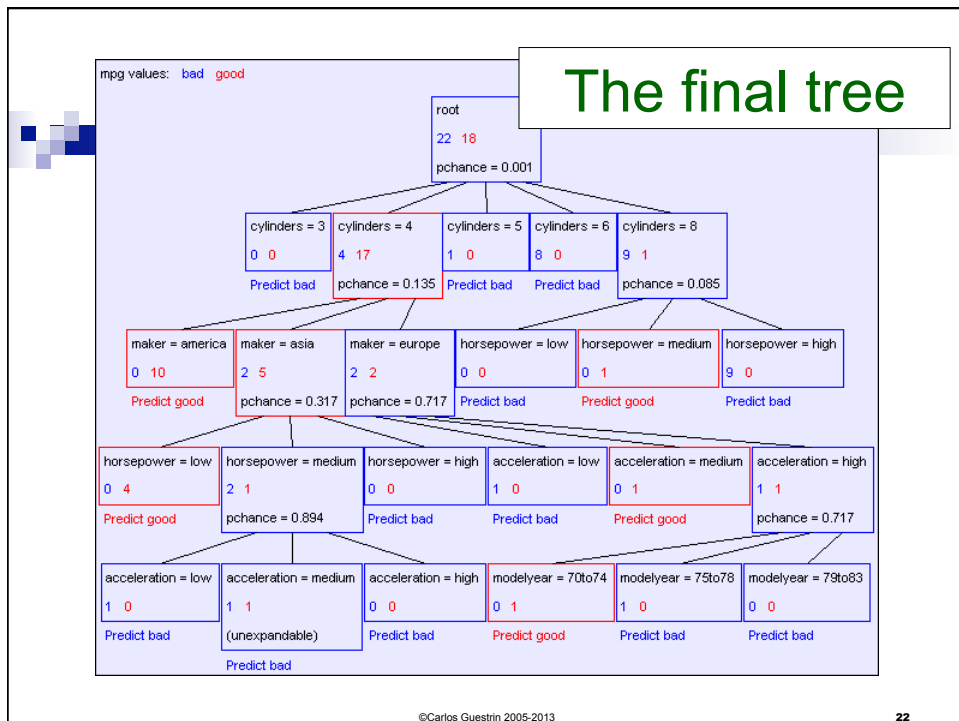
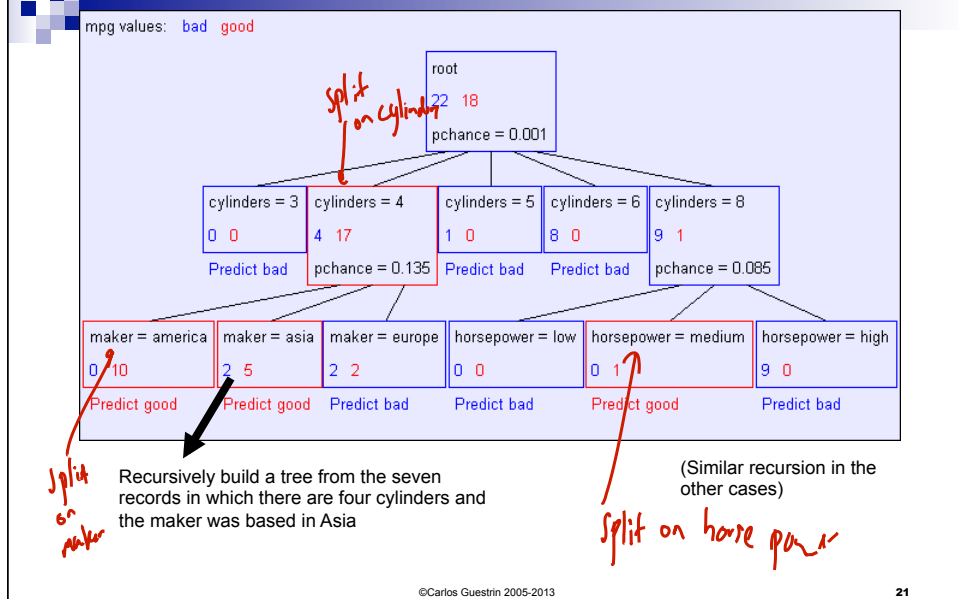
Recursion Step



Recursion Step



Second level of tree



Learning decision trees is hard!!!

- Learning the simplest (smallest) decision tree is an NP-complete problem [Hyafil & Rivest '76]
- Resort to a greedy heuristic:
 - Start from empty decision tree
 - Split on **next best attribute (feature)**
 - Recurse on subset of data associated with each value of the attribute

©Carlos Guestrin 2005-2013

25

Choosing a good attribute

X_1	X_2	Y
T	T	T
T	F	T
T	T	T
T	F	T
F	T	T
F	F	F
F	T	F
F	F	F

X_1

 t/f

 $t: 4$ $f: 0$ $Y=t: 1$ $Y=f: 3$

 absolutely sure

X_2

 t/f

 $t: 3$ $f: 1$ $Y=t: 2$ $Y=f: 2$

 Kind of sure

absolutely unsure / confused

X_1 is better, how now formalize??

©Carlos Guestrin 2005-2013

26

Measuring uncertainty

- Good split if we are more certain about classification after split

- Deterministic good (all true or all false)
- Uniform distribution bad

$P(Y=A) = 1/2$	$P(Y=B) = 1/4$	$P(Y=C) = 1/8$	$P(Y=D) = 1/8$
----------------	----------------	----------------	----------------

more certain

$P(Y=A) = 1/4$	$P(Y=B) = 1/4$	$P(Y=C) = 1/4$	$P(Y=D) = 1/4$
----------------	----------------	----------------	----------------

uniform
very uncertain

©Carlos Guestrin 2005-2013

27

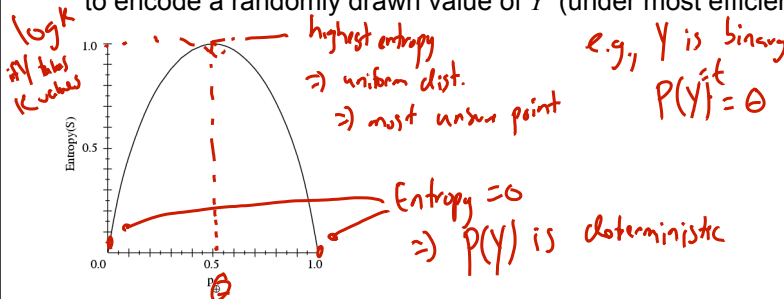
Entropy

Entropy $H(Y)$ of a random variable Y

$$H(Y) = - \sum_{i=1}^k P(Y = y_i) \log_2 P(Y = y_i)$$

More uncertainty, more entropy!

Information Theory interpretation: $H(Y)$ is the expected number of bits needed to encode a randomly drawn value of Y (under most efficient code)



©Carlos Guestrin 2005-2013

28

Andrew Moore's Entropy in a nutshell



Low Entropy



High Entropy

©Carlos Guestrin 2005-2013

29

Andrew Moore's Entropy in a nutshell



Low Entropy



High Entropy

...the values (locations of soup) sampled entirely from within the soup bowl

...the values (locations of soup) unpredictable... almost uniformly sampled throughout our dining room

©Carlos Guestrin 2005-2013

30

Information gain

← Reduction in uncertainty

X ₁	X ₂	Y
T	T	T
T	F	T
T	T	T
T	F	T
F	T	T
F	F	F

■ Advantage of attribute – decrease in uncertainty

□ Entropy of Y before you split

□ Entropy after split

- Weight by probability of following each branch, i.e., normalized number of records

$$H(Y | X) = - \sum_{j=1}^v P(X = x_j) \sum_{i=1}^k P(Y = y_i | X = x_j) \log_2 P(Y = y_i | X = x_j)$$

■ Information gain is difference $IG(X) = H(Y) - H(Y | X)$