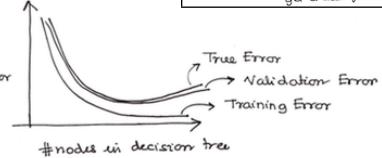


With Pruning:



As tree grows more complicated, training error decreases. At some point, true error stops improving and may even get worse. This will be reflected in the validation error, provided we have enough validation data.

Properties of Entropy:

1. $H(X)$ does not depend on the exact values taken by X , only on the probabilities of distinct values.

2. $H(X) \geq 0$ (always). Why?

$$H(X) = -\sum_{i=1}^k p_i \log p_i, \text{ where } \log p_i \leq 0, \text{ as } p_i \leq 1.$$

3. If X takes k values, $H(X) \leq \log k$. This maximum is achieved only when each value occurs w.p. $1/k$.

Properties of Conditional Entropy:

1. Suppose $X = Z$. What is $H(X|Z)$?

$$H(X|Z) = \sum_Z \Pr(Z=z) H(X|Z=z).$$

Now, given that $Z=z$, we know that $X=z$ w.p. 1.

So, $H(X|Z=z) = 0$.

Thus, $H(X|Z) = \sum_Z \Pr(Z=z) \cdot 0 = 0$.

(This may hold even if $H(X)$ is very large!)

2. Suppose X, Z are independent. What is $H(X|Z)$?

Since X, Z are independent, $\Pr(X|Z=z)$ for any z and x , is equal to $\Pr(X=x)$.

i.e. $X|Z=z$ has exactly the same distribution as X .

So, $H(X|Z=z) = H(X)$

$$H(X|Z) = \sum_Z \Pr(Z=z) H(X|Z=z) = H(X) \sum_Z \Pr(Z=z) = H(X)$$

How to Avoid Overfitting in Decision Trees? By Pruning.

- By pruning the tree using a validation dataset.

1. Split the training ~~data~~ data into training set S and validation set V .

2. Build ID3 tree T using training sets

3. Prune using \vee :

Repeat:

For each node v in T :

Replace subtree rooted at v by single node that predicts majority label in v to get tree T'

If $\text{error}(T')$ on $V \leq \text{error}(T)$ on V , then $T = T'$

Continue till there is no such node v .

Geometric interpretation of: $w_{vt} = w_t + \gamma_t x_t$.

Case 1: $\gamma_t = 1$; $w_{vt} = w_t + x_t$

w_{vt} moves "closer to" x_t
or
 x_t moves towards the + side of the decision boundary (hyperplane)

Case 2: $\gamma_t = -1$; $w_{vt} = w_t - x_t$.

w_{vt} moves "away from" x_t
or
 x_t moves towards the - side of the decision boundary hyperplane.

In both cases, we are moving towards the "correct solution".

- Note: 1. Lower margin \Rightarrow more mistakes.
2. May need > 1 pass over training data to get a classifier with no mistakes.

Measure of Separability: Margin.

For a vector w , and training set S , margin of w wrt S is:

$$y = \min_{(x, y) \in S} \frac{|w \cdot x|}{\|w\|}$$

Example: $S = \{(1, -1), 1\}, \{(-1, 1), -1\}, \{(0, 0), 1\}, \{(-1, 0), -1\}$
 $w = (1, 0)$

What is the margin of w wrt S ? $\|w\| = 1$.

$$\frac{|w \cdot x_1|}{\|w\|} = 1 \quad \text{So, margin} = 0.01$$

$$\frac{|w \cdot x_2|}{\|w\|} = 1$$

$$\frac{|w \cdot x_3|}{\|w\|} = 1$$

$$\frac{|w \cdot x_4|}{\|w\|} = 1$$

Voted Perceptron:

Initially, $m = 1$, $w_1 = 0$, $c_m = 1$.

For $t = 1, 2, 3, \dots$

If $y_t < w_{mt} \cdot x_t + c_m$ then:

$$w_{mt+1} = w_{mt} + y_t x_t$$

$$m = m+1$$

$$c_m = 1$$

$$\text{Else: } c_m = c_m + 1.$$

$$\text{Output: } (w_1, c_1), (w_2, c_2), \dots, (w_m, c_m)$$

On to account of how long w_m survived

How to classify test example x ?

$$\text{Output: } \text{sign}\left(\sum_{i=1}^m c_i \text{sign}(w_i \cdot x)\right)$$

$$[2] K(x, z) = \langle x, z \rangle^2$$

What is the feature space?

* Suppose x, z are 2-d vectors, $x = [x_1, x_2]$, $z = [z_1, z_2]$

$$K(x, z) = (x_1 z_1 + x_2 z_2)^2$$

$$= x_1^2 z_1^2 + 2 x_1 x_2 z_1 z_2 + x_2^2 z_2^2$$

If $\Phi(x) = [x_1^2, x_2^2, \sqrt{x_1} x_2]$, then,

$$K(x, z) = \langle \Phi(x), \Phi(z) \rangle$$

If x, z are d -dimensional vectors, $x = [x_1, \dots, x_d]$, $z = [z_1, \dots, z_d]$,

$$\text{then, } K(x, z) = (x_1 z_1 + \dots + x_d z_d)^2$$

$$= \sum_{i=1}^d x_i^2 z_i^2 + 2 \sum_{i < j} x_i z_i x_j z_j + \dots + x_d^2 z_d^2$$

If $\Phi(x) = [x_1^2, \dots, x_d^2, \sqrt{x_1} x_2, \dots, \sqrt{x_d} x_d]$, then

feature map: from example 2, $\Phi(x)$ and if

$$\Phi(x) = [\Phi(x_1), \sqrt{x_1} x_2, \dots, \Phi(x_d)]$$

$K(x, z) = \langle \Phi(x), \Phi(z) \rangle$

String Kernels.

Let s and t be strings over an alphabet Z , p an integer > 0 .

$K(s, t) = \#\text{of common substrings of length } p \text{ in } s \text{ and } t$.

e.g. for $p=1$, $K(s, t) = \#\text{common letters in } s \text{ and } t$.

so, if $s = axax^2$, $t = aayx^2$, then $K(s, t) = 2$ for $p=1$.

Kernel Properties:

Not all functions $K(x, z)$ are kernels.

Conditions for a function to be a kernel:

1. Symmetry: For all x, z , $K(x, z) = K(z, x)$

2. Positive Semi Definiteness: For a set of

points x_1, \dots, x_m , define kernel matrix as:

$$K_{m \times m}, K_{ij} = K(x_i, x_j)$$

For all x_1, \dots, x_m , the kernel matrix is PSD.

These are necessary and sufficient conditions.

How to show a function $K(x, z)$ is a kernel?

1. Either find a feature map ϕ s.t. $K(x, z) = \langle \phi(x), \phi(z) \rangle$

2. Or show conditions (1) and (2) hold (usually harder).

Fact:

Proof:

$$\|x-y\|^2 = \|x\|^2 + \|y\|^2 - 2 \langle x, y \rangle$$

$$\|x-y\|^2 = \langle x-y, x-y \rangle \quad (\text{From (1)})$$

$$= \langle x, x \rangle - 2 \langle x, y \rangle + \langle y, y \rangle = 2 \langle x, y \rangle$$

(From properties of dot products)

$$= \|x\|^2 + \|y\|^2 - 2 \langle x, y \rangle$$

Training error of a classifier:

$$\text{err}(h) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(h(x_i) \neq y_i)$$

True error of a classifier:

$$\text{err}(h) = \Pr(h(x) \neq y)$$

(x, y i.i.d.)

For a fixed classifier, with more data, training error should approach true error.

Example:

Training Data: $\{(4, 0), 1\}, \{(1, 1), -1\}, \{(0, 1), -1\}, \{(-2, -2), 1\}$

Round 1:

$$\star w_1 = 0$$

$\star y_1 < w_1 \cdot x_1$ for $t=1$

$$= 0 \text{ as well.}$$

$$\star w_2 = w_1 + y_1 x_1$$

$$= (4, 0)$$

Round 2:

$$\star y_2 < w_2 \cdot x_2 \Rightarrow 0 < 0$$

\star So $w_3 = w_2 + y_2 x_2$

$$= (4, 0) - (1, 1) = (3, -1)$$

Training error of h classifier:

$$\text{err}(h) = \frac{1}{4} \sum_{i=1}^4 \mathbb{1}(h(x_i) \neq y_i)$$

Intuitively, average entropy given that we know z .

Entropy: Let X be a random variable that takes values x_1, \dots, x_n with probabilities p_1, \dots, p_n . Then, entropy of X , denoted by $H(X)$ is defined as:

$$H(X) = -p_1 \log p_1 - p_2 \log p_2 - \dots - p_n \log p_n$$

(Note: Use the convention that $0 \cdot \log 0 = 0$)

Linear classification by Hyperplane not thru Origin

We can transform this problem to linear classification by a hyperplane through the origin.

Original Problem:

Training data: $(x_i, y_i), i=1, \dots, n, x_i \in \mathbb{R}^d, y_i \in \{-1, 1\}$

Separating hyperplane: $\tilde{w} \cdot x + b = 0$

Transform it to:

Training data: $(x_i, y_i), i=1, \dots, n, x_i \in \mathbb{R}^{d+1}, y_i \in \{-1, 1\}$

Separating hyperplane: $w \cdot x + b = 0$

Projection of x along w is the vector $\frac{\langle x, w \rangle}{\|w\|} u$.

This is also called "component" of x along w .

Component of x perpendicular to w is $x - \frac{\langle x, w \rangle}{\|w\|^2} u$.

Examples of Kernels:

$$[1] K(x, z) = \langle x, z \rangle^2$$

What is the feature space?

* Suppose x, z are 2-d vectors, $x = [x_1, x_2], z = [z_1, z_2]$

$$K(x, z) = (x_1 z_1 + x_2 z_2)^2$$

$$= x_1^2 z_1^2 + 2 x_1 x_2 z_1 z_2 + x_2^2 z_2^2$$

If $\Phi(x) = [x_1^2, x_2^2, \sqrt{x_1} x_2]$, then,

$$K(x, z) = \langle \Phi(x), \Phi(z) \rangle$$

If x, z are d -dimensional vectors, $x = [x_1, \dots, x_d]$, $z = [z_1, \dots, z_d]$,

$$\text{then, } K(x, z) = (x_1 z_1 + \dots + x_d z_d)^2$$

$$= \sum_{i=1}^d x_i^2 z_i^2 + 2 \sum_{i < j} x_i z_i x_j z_j + \dots + x_d^2 z_d^2$$

then $K(x, z) = \langle \Phi(x), \Phi(z) \rangle$

$$[2] K(x, z) = e^{-\|x-z\|^2/\sigma^2}$$

(could be Gaussian Kernel)

* Time taken to write down $\Phi(x) = \Phi(d^2)$

* Space taken to store $\Phi(x) = \Phi(d^2)$

* Time taken to compute $K(x, z) = \langle \Phi(x), \Phi(z) \rangle$

Corresponds to Φ on \mathbb{R}^d to \mathbb{R}^{d^2}

Graph Kernels.

Let G_1, G_2 be two graphs, $\#G$ an integer.

$$K(G_1, G_2) = \langle \Phi(G_1), \Phi(G_2) \rangle$$

Time to compute $\Phi(G)$ is $O(d^3)$

However, perceptron will never converge to a single w if the data is not linearly separable, as we make more passes over training data.

Margins:

* Data linearly separable with a margin

* Perceptron stops once a separator is found, but we might want to compute the max margin separator

Max margin separator:

- Max margin separator

- Separator output by perceptron

IG(Z) = $H(X) - H(X|Z)$

Information Gain(Z) = $H(X) - H(X|Z)$

(essentially, how much entropy of X is reduced because we know Z).

Conditional Entropy:

Let X, Z be two r.v.s. The conditional entropy of X given Z is defined as:

$$H(X|Z) = \sum_Z \Pr(Z=z) H(X|Z=z)$$

(Intuitively, average entropy given that we know Z).

Entropy: Let X be a random variable that takes values x_1, \dots, x_n with probabilities p_1, \dots, p_n . Then, entropy of X , denoted by $H(X)$ is defined as:

$$H(X) = -p_1 \log p_1 - p_2 \log p_2 - \dots - p_n \log p_n$$

(Note: Use the convention that $0 \cdot \log 0 = 0$)

Boosting Algorithm:

Input: Training set $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$, $y_i = \pm 1$
 $D_t(t) = 1/n$ for all $t = 1, \dots, n$

For $t = 1, 2, 3, \dots$

$h_t = \text{weak-learner wrt } D_t$. (so, $\text{err}_{D_t}(h_t) < 0.5$)

$E_t = \text{err}_{D_t}(h_t)$

$$\alpha_t = \frac{1}{2} \ln \frac{1-E_t}{E_t} \quad (\text{so, } E_t \text{ is high when } h_t \text{ is low, and almost 0 when } h_t \text{ is close to 0.5})$$

$$D_{t+1}(t) = \frac{D_t(t) e^{-\alpha_t y_t h_t(x_t)}}{\sum_i D_t(i)} \quad (\text{D}_t \text{ goes up if } i \text{ is misclassified by } h_t; \text{ so higher } D_t \text{ means harder example.})$$

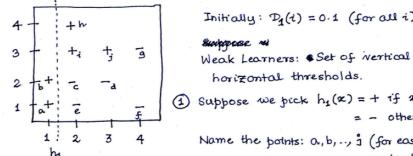
where Z_t is a normalization constant to ensure that

$$\sum_i D_{t+1}(i) = 1.$$

Final classifier: $H(x) = \text{sign} \left(\sum_{t=1}^T \alpha_t h_t(x) \right) \quad (\text{weighted majority of } h_t(x_t)'s)$

Boosting Algorithm Example:

Training data: $(1, 1, +)$ $(2, 1, -)$ $(4, 1, -)$
 $(1, 2, +)$ $(2, 2, -)$ $(3, 2, -)$
 $(2, 3, +)$ $(3, 3, +)$ $(4, 3, -)$
 $(2, 4, +)$



Initially: $D_1(t) = 0.1$ (for all i)

Suppose: Weak Learners: Set of vertical and horizontal thresholds.

- ① Suppose we pick $h_1(x) = +$ if $x_1 \leq 1.5$
 $= -$ otherwise

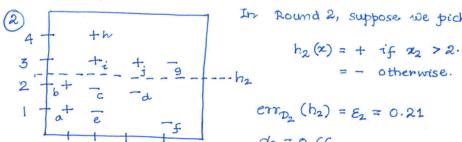
Name the points: a, b, \dots, f (for ease of understanding)

Then:
 $\text{err}_{D_1}(h_1) = E_1 = 0.3$ $\alpha_1 = 0.42$

Weights of a, b, c, d, e, f, g : $D_2 = 0.07$

Weights of h_1, i, j : $D_2 = 0.17$

$$Z_2 = 7 \cdot e^{-0.42} \cdot 0.1 + 3 \cdot 0.1 \cdot e^{-0.42} \\ = 0.92$$



In Round 2, suppose we pick

$$h_2(x) = + \text{ if } x_2 > 2.5 \\ = - \text{ otherwise.}$$

$$\text{err}_{D_2}(h_2) = E_2 = 0.21$$

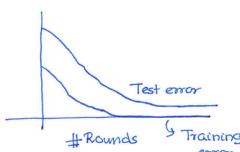
$$\alpha_2 = 0.66$$

$$\begin{aligned} \text{Weights of } a, b: D_3 := 0.07 \times e^{0.66}/Z_3 = 0.17 \\ \text{Weights of } c, d, e, f: D_3 := 0.07 \times e^{0.66}/Z_3 = 0.04 \\ \text{Weights of } h_1, i, j: D_3 := 0.17 \times e^{0.66}/Z_3 = 0.11 \\ \text{Weight of } g: D_3 := 0.07 \times e^{0.66}/Z_3 = 0.17 \\ Z_3 = 0.81 \end{aligned}$$

Boosting and Overfitting:

Overfitting can happen with boosting, but often does not.

Typical boosting run:



Reason is that the margin of classification often increases with boosting.

ID3 Algorithm:

1. Initially, whole training data is at root.

2. While there is an impure node:

(a) Pick any impure node v

(b) Pick a feature f and threshold t along which to "split" the data at v . Done according to "Splitting Rule" (to be described later)

(c) Modify tree as:

$v: \text{Is feature } f \leq t?$ (all data in v)

Yes ↘

No ↘

(all data in v_1 for which $x_f < t$)

(all data in v_2 for which $x_f \geq t$)

If any of v_1 or v_2 is pure (i.e. has data of only one label), then ~~make it a leaf~~ make it a leaf that predicts this label.

Averaged Perceptron uses the same algorithm as voted perceptron, but the classification rule is different.

Averaged Perceptron Classification Rule for test example x :

$$\text{Output: } \text{sign} \left(\sum_{i=1}^m \alpha_i w_i \cdot x_i \right)$$

When do we have high variance?

High variance when there is a small amount of training data, and a very complicated concept class.

When you have high variance, it is called overfitting.

Variance decreases with larger training data, and increases with more complicated classifiers.

High bias \Rightarrow High training and test errors

High variance \Rightarrow Low training error, high test errors.

Bias-variance tradeoff: If we make the concept class more complicated, then, for the same training set size, bias decreases but variance increases. Thus there is a bias-variance tradeoff.

Bias is the true error of the best classifier in the concept class

e.g. best linear separator, best decision tree, etc.

When do we have high bias?

High bias when the concept class cannot model the true data distribution well; this doesn't depend on the training data size.

When you have high bias, it's called underfitting

Problem 1

	Cancer	No Cancer
1. Gene A	$\frac{1}{2}$	$\frac{1}{2}$
0. No Gene A	$\frac{1}{2}$	$\frac{1}{2}$

	Cancer	No Cancer
1. Gene B	$\frac{1}{3}$	$\frac{2}{3}$
0. No Gene B	$\frac{2}{3}$	$\frac{1}{3}$

1.1) conditional distribution for $X|Y = y$ for $y = 0, 1$

$$P(X|Y=0) = \frac{P(X=0|Y=0)}{P(Y=0)} = \frac{\frac{1}{2}}{\frac{1}{2} + \frac{1}{2}} = \frac{1}{2}$$

$$X|Y=1 : P_r(X=0|Y=1) = \frac{P_r(X=0 \cap Y=1)}{P_r(Y=1)} = \frac{\frac{1}{2} \cdot \frac{1}{3}}{\frac{1}{2} \cdot \frac{1}{3} + \frac{1}{2} \cdot \frac{2}{3}} = \frac{1}{3}$$

$$X|Y=1 : P_r(X=1|Y=1) = \frac{P_r(X=1 \cap Y=1)}{P_r(Y=1)} = \frac{\frac{1}{2} \cdot \frac{2}{3}}{\frac{1}{2} \cdot \frac{1}{3} + \frac{1}{2} \cdot \frac{2}{3}} = \frac{2}{3}$$

1.2) Calculate Conditional Entropies $H(X|Y)$ and $H(X|Z)$

$$H(X|Y) = P(Y=0) H(X|Y=0) + P(Y=1) H(X|Y=1)$$

$$H(X|Y=0) = -\frac{1}{2} \log(\frac{1}{2}) - \frac{1}{2} \log(\frac{1}{2}) = .3010299957$$

$$H(X|Y=1) = -\frac{1}{3} \log(\frac{1}{3}) - \frac{2}{3} \log(\frac{2}{3}) = .1956762468$$

$$H(X|Y) = \frac{1}{2}(.3010299957) + \frac{1}{2}(.1956762468) = .2487877464$$

Definitions:

1. Weak Learner: A simple rule of thumb that doesn't necessarily work very well.

2. Strong Learner: A good classifier (with high accuracy)

$$\text{err}_W(h) = \sum_{i=1}^n w_i \mathbb{1}(h(x_i) \neq y_i)$$

When to stop boosting? Use a validation dataset to find a stopping time.

Stop when validation error does not improve.

When x and y are orthogonal,

$$\langle x, y \rangle = 0$$

$$\|x - y\|^2 = \|x\|^2 + \|y\|^2$$

Cauchy-Schwarz Inequality:

For any vectors x and y ,

$$\langle x, y \rangle \leq \|x\| \cdot \|y\|$$

A matrix $A_{d \times d}$ is PSD if for all $d \times 1$ vectors x ,

$$x^T A x \geq 0$$

$$x^T A x = x^T z = \sum_{i=1}^d x_i z_i = \sum_{i=1}^d \sum_{j=1}^d A_{ij} x_i z_j$$

Kernels to Distances:

Given a kernel function $K(x, z)$, define a distance function:

$$D_K(x, z) = \sqrt{K(x, x) + K(z, z) - 2K(x, z)}$$

If $K(x, z) = \langle \Phi(x), \Phi(z) \rangle$ then:

$$D_K(x, z) = \langle \Phi(x), \Phi(x) \rangle + \langle \Phi(z), \Phi(z) \rangle - 2 \langle \Phi(x), \Phi(z) \rangle$$

$$= \|\Phi(x) - \Phi(z)\|^2$$