

## Decision Trees

Entropy:  $H(Y) = -\sum_{y \in \text{values}(Y)} P(Y=y) \log_2 P(Y=y)$

0 = pure, 1 = perfectly mixed

Mutual Information:  $I(X;Y) = H(Y) - \sum_{x \in \text{values}(X)} P(X=x) H(Y|X=x)$

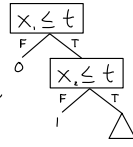
1 = 100% correlation between attribute & label

Memorizer: learns labels given attributes, MV otherwise

Decision stump:



Decision Tree



## K-Nearest Neighbor

Inductive Bias: Similar/Nearby points should have sim. labels  
All label dimensions are created equal

Euclidian Distance:  $d(p,q) = \sqrt{\sum_{i=1}^n (q_i - p_i)^2}$

\*feature scale could dramatically influence class results

## Model Selection

hyperparameters - tunable aspects of the model, that the learning algorithm does not select

K-fold cross-validation: create k partitions in D  
do k runs: train using k-1 partitions and calc validation error on remaining portion (rotating validation partition on each run)

report average validation error

leave-one-out validation: K = N partitions

train on N-1 samples; validate only one sample per run

\*model selection/hyperparameter optimization is just another form of learning

## Derivative Rules

$\frac{d}{dx}(\text{constant}) = 0$     $\frac{d}{dx}(x) = 1$     $\frac{d}{dx}(x^2) = 2x$     $\frac{d}{dx}(\sqrt{x}) = \frac{1}{2} \cdot x^{-\frac{1}{2}} = \frac{1}{2\sqrt{x}}$

$\frac{d}{dx}(e^x) = e^x$     $\frac{d}{dx}(a^x) = \ln(a) a^x$     $\frac{d}{dx}(\ln(x)) = \frac{1}{x}$     $\frac{d}{dx}(\log_a(x)) = \frac{1}{x \ln(a)}$

$\frac{d}{dx}(\sin(x)) = \cos(x)$     $\frac{d}{dx}(\cos(x)) = -\sin(x)$     $\frac{\sin(x)}{\cos(x)} = \tan(x)$     $\frac{d}{dx}(\tan(x)) = \sec^2(x)$

Multiplication by constant:  $\frac{d}{dx}(c \cdot f(x)) = c \cdot \frac{d}{dx}(f(x))$

Power Rule:  $\frac{d}{dx}(x^n) = nx^{n-1}$

Reciprocal Rule:

$\frac{d}{dx}(1/f) = -\frac{d}{dx}(f) / f^2$

Sum Rule:  $\frac{d}{dx}(f+g) = \frac{d}{dx}(f) + \frac{d}{dx}(g)$

Difference Rule:  $\frac{d}{dx}(f-g) = \frac{d}{dx}(f) - \frac{d}{dx}(g)$

Product Rule:  $\frac{d}{dx}(f \cdot g) = \frac{d}{dx}(f) \cdot g + f \cdot \frac{d}{dx}(g)$

Quotient Rule:  $\frac{d}{dx}(f/g) = \frac{\frac{d}{dx}(f) \cdot g - f \cdot \frac{d}{dx}(g)}{g^2}$

Chain Rule:

- as composition of functions:  $\frac{d}{dx}(f \circ g) = (f' \circ g) \cdot g'$

- using ' :  $\frac{d}{dx}(f(g(x))) = f'(g(x)) \cdot g'(x)$

- using  $\frac{d}{dx} \cdot \frac{dy}{dx} = \frac{dy}{du} \cdot \frac{du}{dx}$

$f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$

Inductive Bias:  
prefer the smallest tree consistent w/ the training data (i.e.) 0 error rate

## Perceptron

\* try to learn hyperplane directly

$h(x) = \text{sign}(\theta^T x)$  for  $y \in \{-1, +1\}$

First Perceptron: Linear Classifier

$\hat{y} = h(\vec{x}) = \text{sign}(\vec{w}^T x + b)$

$= \text{sign}(w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b)$

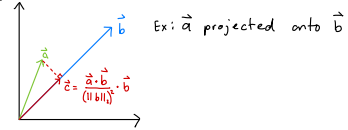
orthogonal:  $\vec{a}$  is orthogonal to  $\vec{b}$  iff  $\vec{a}^T \vec{b} = 0$

(right angle to each other)

dot product:  $\vec{a} \cdot \vec{b} = \vec{a}^T \vec{b} = \sum_i a_i b_i$

$l_2$  norm:  $\|\vec{u}\|_2 = \sqrt{\sum_i (u_i)^2}$  (length of vector / Euclidian Distance)

vector projection:



hyperplanes: 2D-line  $S = \{ \vec{x} : \vec{w}^T \vec{x} + b = 0 \}$

3D-plane  
4+D-hyperplane

halfspaces: all the points on one side of the hyperplane

$S_+ = \{ \vec{x} : \vec{w}^T \vec{x} + b > 0 \}$

$S_- = \{ \vec{x} : \vec{w}^T \vec{x} + b < 0 \}$

batch learning: learn all examples at once

online learning: gradually learn as example is recieved

initialize parameters  $\vec{w} = [w_1, w_2, \dots, w_n]^T = [0, 0, \dots, 0]^T$   
 $b = 0$  ← \*start at zero vector

for i in range(M):

- recieve instance  $\vec{x}^{(i)}$
- predict  $\hat{y} = h(\vec{x}) = \text{sign}(\vec{w}^T \vec{x} + b)$     $\text{sign}(a) = \begin{cases} +1 & \text{if } a \geq 0 \\ -1 & \text{otherwise} \end{cases}$
- recieve label  $y^{(i)}$
- update parameters: if  $y^{(i)} = +1$  and  $y^{(i)} \neq \hat{y}$ :  
 $\vec{w} \leftarrow \vec{w} + \vec{x}^{(i)}$   
 $b \leftarrow b + 1$    \*positive mistake

if  $y^{(i)} = -1$  and  $y^{(i)} \neq \hat{y}$ :  
 $\vec{w} \leftarrow \vec{w} - \vec{x}^{(i)}$   
 $b \leftarrow b - 1$    \*negative mistake

hyperplane w.  $\theta$ :  $H = \{ x : w^T x + b = 0 \} = \{ x : \theta^T x' = 0 \text{ and } x'_0 = 1 \}$

$\theta = [b, w_1, w_2, \dots, w_n]^T$

$x' = [1, x_1, x_2, \dots, x_n]^T$

halfspaces w.  $\theta$ :  $H_+ = \{ x : \theta^T x > 0 \text{ and } x'_0 = 1 \}$

$H_- = \{ x : \theta^T x < 0 \text{ and } x'_0 = 1 \}$

mistake bound: if some data has margin  $\gamma$  and all points lie inside a ball of radius R rooted at the origin, then the online

Perceptron algorithm makes  $\leq (R/\gamma)^2$  mistakes

$d = \frac{ax + by + c}{\sqrt{a^2 + b^2 + c^2}}$  where  $a = w_1, b = w_2, c = b, x = x_1, y = x_2$

Summary: -linear classifier

- simple: mistake = update params

- converge if linearly seperable

- can bound # of mistakes

## Linear Regression

assume  $y = w^T x$

parameters:  $w = [w_0, w_1, w_2, \dots, w_D]^T$

objective function: minimize the squared error

$$\ell_D(w) = \sum_{n=1}^N (w^T x^{(n)} - y^{(n)})^2$$

optimize objective: solve in closed form

- take partial derivatives  $\rightarrow$  gradient  $\nabla$

- set them equal to 0

suppose  $y \in \mathbb{R}$  and  $D$ -dimensional inputs  $x = [1, x_1, x_2, \dots, x_D]^T \in \mathbb{R}^{D+1}$

training data  $D = \{(x^{(n)}, y^{(n)})\}_{n=1}^N$

$$X = \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_D^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(N)} & \dots & x_D^{(N)} \end{bmatrix} \in \mathbb{R}^{N \times (D+1)}$$

is the design matrix

$y = [y^{(1)}, y^{(2)}, \dots, y^{(N)}]^T \in \mathbb{R}^N$  is the target vector

Algorithm:

Input:  $D = \{X, y\} = \{(x^{(n)}, y^{(n)})\}_{n=1}^N$

1) Compute the pseudoinverse of  $X$  either

a) Directly as  $X^\dagger = (X^T X)^{-1} X^T$  if  $X^T X$  is invertible or

b) Through singular value decomposition (SVD): if

$X = U \Sigma V^T$ , then  $X^\dagger = V \Sigma^\dagger U^T$  where  $\Sigma^\dagger$  inverts all non-zero elements of  $\Sigma$

2) Compute  $\hat{w} = X^\dagger y$

Output:  $\hat{w}$  weight vector

Inductive Bias:  
true relationship  
between inputs &  
outputs is linear

## Gradient Descent

current weight vector is  $\hat{w}$  (slope  $m$ )

move some distance  $\eta$  in most

downhill direction,  $\hat{v}$ :

$$w^{(t+1)} = w^{(t)} + \eta \hat{v}$$

gradient points in direction of

steepest increase,  $\hat{v}$  is opposite:

$$\hat{v}^{(t)} = - \frac{\nabla_w \ell_D(w^{(t)})}{\|\nabla_w \ell_D(w^{(t)})\|_2}$$

$$\eta^{(t)} = \eta^{(0)} \cdot \|\nabla_w \ell_D(w^{(t)})\|$$

step size init step size magnitude of grad

each iteration:

$$w^{(t+1)} = \underset{\substack{\uparrow \\ \text{next weight} \\ \text{vector}}}{w^{(t)}} + \underset{\substack{\uparrow \\ \text{curr} \\ \text{weight vec}}}{\eta^{(t)}} \underset{\substack{\uparrow \\ \text{step} \\ \text{size}}}{\hat{v}^{(t)}}$$

vector pointing in opposite dir from gradient

..... math

$$= \underset{\substack{\uparrow \\ \text{curr} \\ \text{weight vec}}}{w^{(t)}} - \underset{\substack{\uparrow \\ \text{learning} \\ \text{rate}}}{\eta^{(0)}} \underset{\substack{\uparrow \\ \text{gradient} \\ \text{(vect of} \\ \text{derivatives)} \\ \text{of loss func} \\ \text{wrt } w^{(t)}}}{\nabla_w \ell_D(w^{(t)})}$$

## Minimizing the Squared Error

$$\ell_D(w) = \sum_{n=1}^N (w^T x^{(n)} - y^{(n)})^2 = \sum_{n=1}^N (x^{(n)T} w - y^{(n)})^2$$

$$= \|Xw - y\|_2^2 \text{ where } \|z\|_2 = \sqrt{\sum_{d=1}^D z_d^2} = \sqrt{z^T z}$$

$$= (Xw - y)^T (Xw - y)$$

$$= (w^T X^T X w - 2w^T X^T y + y^T y)$$

$$\nabla_w \ell_D(\hat{w}) = (2X^T X \hat{w} - 2X^T y) = 0$$

$$\rightarrow X^T X \hat{w} = X^T y$$

$$\rightarrow \hat{w} = (X^T X)^{-1} X^T y$$

$$\ell_D(w) = \sum_{n=1}^N (w^T x^{(n)} - y^{(n)})^2 = \sum_{n=1}^N (x^{(n)T} w - y^{(n)})^2$$

$$= \|Xw - y\|_2^2 \text{ where } \|z\|_2 = \sqrt{\sum_{d=1}^D z_d^2} = \sqrt{z^T z}$$

$$= (Xw - y)^T (Xw - y)$$

$$= (w^T X^T X w - 2w^T X^T y + y^T y)$$

$$\nabla_w \ell_D(w) = (2X^T X w - 2X^T y)$$

$$H_w \ell_D(w) = 2X^T X$$

$H_w \ell_D(w)$  is positive semi-definite

## Matrix

\* for dot product,  $w_1 \cdot w_2 = h_2$

$$\begin{bmatrix} A & B \\ C & D \\ E & F \end{bmatrix} \times \begin{bmatrix} G \\ H \end{bmatrix} = \begin{bmatrix} AG + BH \\ CG + DH \\ EG + FH \end{bmatrix}$$

$$(AB)^T = B^T A^T$$

## Linear Regression via Gradient Descent

Input:  $D = \{X, y\} = \{(x^{(n)}, y^{(n)})\}_{n=1}^N$

design matrix  $\uparrow$  target vector  $\uparrow$  attributes  $\uparrow$  labels  $\uparrow$  learning rate  $\uparrow$  dataset size

1) initialize  $w^{(0)}$  to all zeros,  $w = [w_0, w_1, \dots, w_D]^T$  { init  $t=0$

2) while termination criteria is not satisfied (loop):

a) compute gradient  $\nabla_w \ell_D(w^{(t)}) \leftarrow$  gradient (vector of derivatives of loss function) wrt  $w^{(t)}$

b) update  $w$ :  $w^{(t+1)} = w^{(t)} - \eta^{(0)} \nabla_w \ell_D(w^{(t)})$

c) increment  $t$

Output:  $w^{(t)}$   $\leftarrow$  regression line