

IP

Defi

= I  
i

nen

r

ifia

19

M1

15

6.

7.

II

X

$$(f(x_i)) \text{ decreasing} \Rightarrow f(x_K) - f(x^*) \leq f(x_i) - f(x^*) \quad \forall i \in K$$

$$\therefore K(f(x_K) - f(x^*)) = \sum_{i=1}^K (f(x_K) - f(x^*))$$

$$\leq \sum_{i=1}^K (f(x_i) - f(x^*))$$

$$\leq \frac{1}{2\alpha} \|x_0 - x^*\|^2$$

### \* Lecture 9

#### 4.4 Constrained Optimization

Consider  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  and  $S \subseteq \mathbb{R}^n$  and find

$$x^* = \underset{\substack{x \in S}}{\operatorname{argmin}} f(\vec{x}) \quad \text{or} \quad f_{\min} = \underset{\substack{x \in S}}{\min} f(\vec{x}) \quad (\text{or max})$$

same problem

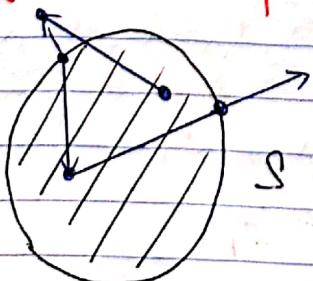
E.g.  $A \in \mathbb{R}^{m \times n}$  and consider

$$f(\vec{x}) = \|A\vec{x}\|_2 \quad f: \mathbb{R}^n \rightarrow \mathbb{R}$$

$$\max_{\substack{\vec{x}: \|A\vec{x}\|_2 = 1}} f(\vec{x}) = \max_{\substack{\vec{x}: \|A\vec{x}\|_2 = 1}} \|A\vec{x}\|_2 \\ = \|A\|_2.$$

• one approach (modification of gradient descent)

Idea: At each step, "project" back into  $S$ .



• More general

The Karush-Kuhn-Tucker approach (generalizes Lagrange multipliers)



扫描全能王 创建

negation:  $\vec{x} \notin S \Rightarrow g_i(\vec{x}) \neq 0$  or  $h_j(\vec{x}) > 0$  for some  $i$  and  $j$

number of conditions you have

$1 \leq i \leq n$

- \* If we can write  $S = \{ \vec{x} | g_i(\vec{x}) = 0 \text{ and } h_j(\vec{x}) \leq 0 \quad 1 \leq j \leq m \}$  for some functions  $g_i, h_j: \mathbb{R}^n \rightarrow \mathbb{R}$ .

Define the Lagrangian

$$\vec{x}^T \vec{g}(\vec{x}) = \sum_i \lambda_i g_i(\vec{x}) L(\vec{x}; \vec{\lambda}, \vec{\alpha}) = f(\vec{x}) + \sum_i \lambda_i g_i(\vec{x}) + \sum_j \alpha_j h_j(\vec{x})$$

$$\vec{x}^T \vec{h}(\vec{x}) \text{ then min}_{\vec{x} \in \mathbb{R}^n} \max_{\vec{\lambda}} \max_{\vec{\alpha}: \vec{\alpha} \geq 0} L(\vec{x}; \vec{\lambda}, \vec{\alpha}) \quad (*)$$

has the same solutions  $\vec{x}$  as  $\min_{\vec{x} \in S} f(\vec{x})$ .

Note  
use

Justification Anytime  $\vec{x} \notin S \Rightarrow g_i(\vec{x}) \neq 0$

$$\max_{\vec{\lambda}} \max_{\vec{\alpha}: \vec{\alpha} \geq 0} L(\vec{x}; \vec{\lambda}, \vec{\alpha}) = \infty$$

unbounded.  $\therefore$  can make  $\vec{\lambda}$  as large as we want.

and whenever  $\vec{x} \in S$

$$\max_{\vec{\lambda}} \max_{\vec{\alpha}: \vec{\alpha} \geq 0} L(\vec{x}; \vec{\lambda}, \vec{\alpha}) = f(\vec{x}) \quad \because g_i(\vec{x}) = 0 \quad \text{and} \quad h_j(\vec{x}) \leq 0.$$

Theorem [Korush, Kuhn, Tucker]: Let  $f$  be differentiable.  $\exists \vec{\lambda} \geq 0$  such that  $\vec{x} \in S$  and  $\nabla_{\vec{x}} L(\vec{x}; \vec{\lambda}, \vec{\alpha}) = 0$

$$\textcircled{1} \quad \nabla_{\vec{x}} L(\vec{x}; \vec{\lambda}, \vec{\alpha}) = 0$$

$$\textcircled{2} \quad \vec{x} \in S$$

gradients

could be a set of possible slopes.

(i.e. set of subderivatives).

$$\textcircled{3} \quad \sum \alpha_i h_i(\vec{x}) = 0$$

The following conditions are necessary for  $\vec{x}$  to be a solution to  $(*)$

Eg. Let  $A \in \mathbb{R}^{n \times n}$  be symmetric.

$$f(\vec{x}) = -\vec{x}^T A \vec{x} \quad \text{find}$$

$$\min_{\vec{x}: \| \vec{x} \|_2 \leq 1} f(\vec{x})$$

$$h(\vec{x}) = \| \vec{x} \|_2^2 - 1 = \vec{x}^T \vec{x} - 1 \leq 0 \quad \text{How to get?}$$

$$L(\vec{x}; \vec{\alpha}) = -\vec{x}^T A \vec{x} + \alpha(\vec{x}^T \vec{x} - 1)$$

usually the

most significant one

$$\textcircled{1} \Rightarrow \nabla_{\vec{x}} L(\vec{x}; \vec{\alpha}) = 0$$

$$\Rightarrow -\vec{x}^T (A + A^T) + 2\alpha \vec{x}^T = 0$$

$$\Rightarrow -2\vec{x}^T A = 2\alpha \vec{x}^T \quad \text{taking transpose on both sides.}$$

$$A = B \Rightarrow A^T = B^T$$

$$\Rightarrow A \vec{x} = \alpha \vec{x} \quad \text{so } \vec{x} \text{ is an eigenvector of } A$$

e.g. if  $AB = 0$

$$\Rightarrow B^T A^T = E^T D^T C^T$$

$$\therefore A^T = A \quad \therefore (2A)^T = 2A^T = 2A \quad \alpha \text{ is a constant } / 1 \times 1 \text{ matrix}$$



扫描全能王 创建

$$\begin{aligned} \vec{x}^T \vec{y} &= \vec{y}^T \vec{x} \\ \| \vec{x} \|_2 &= \| \vec{y} \|_2 \end{aligned}$$

$$\begin{aligned} \sum_i x_i y_i &= \sum_i y_i x_i \\ \vec{x}^T \vec{y} &= \vec{y}^T \vec{x} \\ \textcircled{2} \quad \| \vec{x} \|_2 \leq 1 &\Rightarrow \vec{y}^T \vec{x} = \vec{y}^T \vec{x} \| \vec{x} \|_2 - 1 = 0 \Rightarrow \| \vec{x} \|_2 - 1 = 0 \Rightarrow \| \vec{x} \|_2 = 1 \\ \textcircled{3} \quad \| \vec{x} \|_2 = 1 &\Rightarrow \vec{y}^T \vec{x} = 0 \Rightarrow \vec{y}^T (\vec{x} - \vec{x}) = 0 \quad \text{either } \vec{x} = 0 \text{ or } \vec{x} = 0 \\ \text{solution: } \vec{x} &\text{ is the (normalized) eigenvector of } A \text{ with the largest eigenvalue.} \\ f(\vec{x}) &= -\lambda \end{aligned}$$

$f(\vec{x}) = -\lambda$   
when  $\vec{x}$  is an EV w/  
EV  $\lambda$ .

## \* Chapter 5 Basics of ML

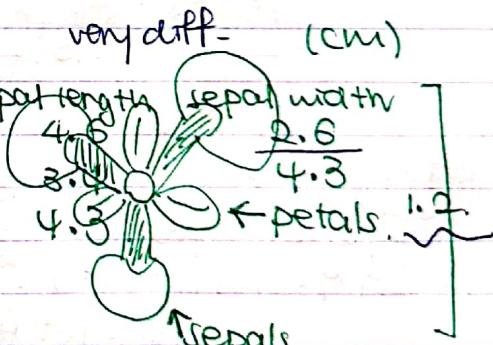
"Definition" A learning algorithm is a computer program which accomplishes some task  $T$  based on some experience  $E$  with respect to some performance measure  $P$ , such that performance  $P$  increases as we add experience to  $E$ .

(not formal) understand what is meant by E.T.P by example.

• Experience  $E$ : usually a dataset  $X \in \mathbb{R}^{m \times n}$  where each column is called a feature and each row is called an example.

E.g. Suppose we have 150 Iris plants.

	petal length	petal width	sepal length	sepal width
Plant 1	5.1	1.3	1.7	0.5
Plant 2	4.9	1.4	4.0	1.4
Plant 3	4.7	1.8	4.4	1.4
⋮	⋮	⋮	⋮	⋮



3 different species known 1, 2, 3.

$E$  may also consist of a vector  $y$  of labels.

$$E.g. \quad y = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$$

A learning algorithm with  $X, y$  is called supervised.

" " " " " " unsupervised.

Unsupervised Goal Determine a probability distribution

$p(x)$  where  $x \in \mathbb{R}^n \leftrightarrow x_i, *$ .

Supervised usually want to determine  $p(y|x)$ .

white logistic regression  
is a classification  
also. If  $y$  takes discrete values (e.g.  $1, \dots, k$ )  $\rightarrow$  classification.  
If  $y$  takes continuous " (e.g.  $\mathbb{R}$  or  $[a, b] \subseteq \mathbb{R}$ ) \rightarrow regression



扫描全能王 创建

## Lecture 10

(E.T.P, program).

- Task I: ① Classification Given  $(x, y) \in \mathbb{R}^{m \times n} \times \mathbb{R}^m$  (Examples) produce  $f: \mathbb{R}^n \rightarrow \{1, \dots, k\}$  s.t.  $f(x) = y$ . performs well as a prediction.

E.g. object recognition  $x_i, * = \text{pixel data of an image}$   
 $y_i = \text{cat or dog}$

- ② Regression Given  $(x, y)$ ,  $y \in \mathbb{R}$

task: produce  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  s.t.  $f(x) = y$  performs well as a prediction

E.g. (Predicting Housing Prices)

$x_i, * = \text{data about a house.}$

$y_i = \text{price of house}$

each task  
would require  
a very different  
type of  
algorithm.

- ③ Given an image of text, produce text in Unicode (digital).

- ④ Translate from one language to another.

- ⑤ Generate new examples similar to those in E.

- ⑥ Given X with some entries missing, fill in entries (as reliably as possible).

- ⑦ Given X, remove noise (SVD often useful)

## Performance Measure P

Depends on T (the specific task you are doing).

- ① (classification) Accuracy/Precision/Recall

different types of performance measure for classification



扫描全能王 创建

② (Regression) Average (squared) error in prediction.

★ Performance is not measured using examples from E but from a different set of examples.

$$x^{(test)}, y^{(test)}$$

### • Program

Examples: ② (Regression) Linear Regression

using same idea gives us E: (supervised) Dataset  $x \in \mathbb{R}^{m \times n}$  labels  $y \in \mathbb{R}^m$ .  
Polynomial regression.  $y_i \in \mathbb{R}$

T: Produce  $f: \mathbb{R}^n \rightarrow \mathbb{R}$

P: Given test set  $x^{(test)}, y^{(test)}$   $m = \# \text{examples}$

want to minimize

case ①: Expect  $y_i$  depends (linearly) on  $x_{i,*}$ ,

so  $\exists w$  s.t.  $y_i \approx x_{i,*} w + b$ .

then, we want to minimize Introduce  $\hat{x} = \begin{bmatrix} 1 & x \end{bmatrix}$

$$\frac{1}{m} \sum_{i=1}^m \| \hat{x}_{i,*} w - y_i \|_2^2$$

so no need to have  $+b$

$$= \frac{1}{m} \| \hat{x} w - y^{(test)} \|_2^2 \quad \text{should always be evaluated on the test sets.}$$

Define the cost function (Loss, Error, ...)

$$J(w) = \frac{1}{m} \| \hat{x} w - y^{(test)} \|_2^2 \quad \leftarrow \begin{array}{l} \text{cost function} \\ \text{defined using } E. \end{array}$$

Find  $w^* = \arg \min_w J(w)$

More justification  
in 3.5.1.

solution 1: Use gradient descent to approximate  $w^*$ .

proposition:  $J$  is twice differentiable, convex, and has Lipschitz continuous gradient. Furthermore

$$\textcircled{1} \quad \nabla_w J(w) = \frac{2}{m} [\hat{x}^T \hat{x} w - \hat{x}^T y]$$

$$\textcircled{2} \quad H(J)(w) = \frac{2}{m} \hat{x}^T \hat{x}$$

proof: HW



扫描全能王 创建

solution 2: proposition If  $\hat{X}$  has LI columns and  $J$  has a global minimum,  $w^*$ , then  $w^* = (\hat{X}^\top \hat{X})^{-1} \hat{X}^\top y$

proof: Since  $J$  is differentiable,  $w^*$  global min  $\Rightarrow w^*$  local min  
 $\Rightarrow \nabla J(w^*) = 0$ .

$$\nabla J(w^*) = \frac{2}{m} [\hat{X}^\top \hat{X} w - \hat{X}^\top y] = 0$$

$$\Rightarrow \hat{X}^\top \hat{X} w = \hat{X}^\top y$$

If  $\hat{X}^\top \hat{X}$  not

some columns of  $\hat{X}$  are LI,  $\hat{X}^\top \hat{X}$  is pos def  $\Rightarrow \hat{X}^\top \hat{X}$  is invertible.

Note: invertible, then

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

use  $\exists w \neq 0$  s.t.  $\hat{X}^\top \hat{X} w = 0$  hard to invert a matrix if the

by rank nullity.

matrix is large

then  $w^\top \hat{X}^\top \hat{X} w = 0$

pos def.

case ② we expect  $y_i$  is a polynomial function  $x_{i,t}$  of maximum total degree  $n$ .

$$\text{so } f(x_1, \dots, x_n) = f(\vec{x}) = \sum_{\alpha \in \mathbb{Z}_{\geq 0}^n} a_\alpha x^\alpha \quad \text{E.g. } \alpha = (0, 1, 2, 3) \\ x^\alpha = x_1 x_2^2 x_3^3 x_4^3$$

solution: modify  $X$  by adding new columns, one for each  $\alpha \in \mathbb{Z}_{\geq 0}^n$  s.t.  $| \alpha | \leq n$   $| \alpha | = \sum_i \alpha_i$   
 total degree

$$\text{E.g. } X = \begin{bmatrix} 1 & 2 \\ 2 & 1 \\ 3 & 4 \end{bmatrix} \xrightarrow{\substack{\text{max total deg} \geq 3}} \begin{bmatrix} x_1 & x_2 & x_1^2 & x_1 x_2 & x_1^3 \\ 1 & 2 & 1 & 2 & 4 \\ 2 & 1 & 4 & 2 & 1 \\ 3 & 4 & 9 & 12 & 16 \end{bmatrix}$$

again, form  $\hat{X}$  and proceed as in Case 1.

final coefficient  $w$  that minimize prediction error.



扫描全能王 创建

## \*Lecture 11

Example ① Logistic Regression

E:  $x \in \mathbb{R}^{m \times n}$ ,  $y \in \{0, 1\}^m$  (supervised)

I: produce a function  $f: \mathbb{R}^n \rightarrow [0, 1]$

from this, we define  $\text{pred } p: \mathbb{R}^n \rightarrow \{0, 1\}$ .

$$\text{pred } p(\vec{x}) = \begin{cases} 1 & f(\vec{x}) > 0.5 \\ 0 & f(\vec{x}) \leq 0.5 \end{cases}$$

P: Many choices exist. Let's go with accuracy.

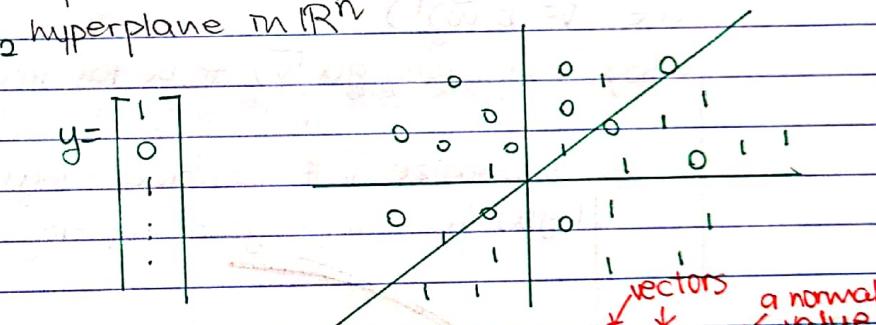
Note: Also could use precision or recall.

$$\text{Acc} = \frac{\#\{i | \text{pred } p(x_i) = y_i\}}{\#\text{rows in } X^{(test)}}$$

A solution: case ① we expect  $y=1$  examples are reliably separated from  $y=0$  examples by a hyperplane in  $\mathbb{R}^n$

E.g.  $x \in \mathbb{R}^{m \times 2}$

$$x = \begin{bmatrix} 1 & 2 \\ 2 & 1 \\ 4 & 3 \\ \vdots & \vdots \end{bmatrix} \quad y = \begin{bmatrix} 1 \\ 0 \\ 1 \\ \vdots \end{bmatrix}$$



Goal: Find  $\vec{w} \in \mathbb{R}^n$  b  $\in \mathbb{R}$  s.t. the hyperplane  $\vec{w}^\top \vec{x} + b = 0$  separates 1's from 0's reliably.

we need the

b, so the

equation

for the hyperplane

does not go

through the

origin.

separates 1's from 0's reliably.

equation for a hyperplane

algorithm?

want function  $J(\vec{w}, b)$  which is minimized when accuracy is maximized. (Cost function)

$$\vec{x} \cdot \vec{w} + b > 0 \Rightarrow \text{strong predict for 1.}$$

$$\vec{x} \cdot \vec{w} + b < 0 \Rightarrow \text{strong predict for 0.}$$

we want these.

This suggests we use a function g with the following properties

$$\lim_{z \rightarrow \infty} g(z) = 1.$$

$$g(0) = 1/2.$$

$$\lim_{z \rightarrow -\infty} g(z) = 0.$$



扫描全能王 创建

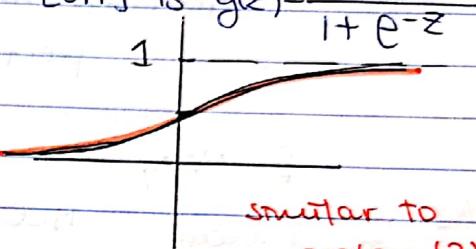
With this, we can compute  $g(\vec{x} \cdot \vec{w} + b)$  and predict 1 if  $> 0.5$   
 0 if  $< 0.5$

Defn: The sigmoid function  $g: \mathbb{R} \rightarrow [0,1]$  is  $g(z) = \frac{1}{1 + e^{-z}}$

prop: ①  $g$  satisfies (+)

$$\textcircled{2} \quad g'(z) = g(z)(1-g(z))$$

proof: Exercise



similar to  
 $\arctan(?)$

Again, we use  $\hat{x} = \begin{bmatrix} \vdots \\ x \\ \vdots \end{bmatrix}$  so  $\hat{x}^T v$  instead of  $x^T w + \begin{bmatrix} b \\ \vdots \\ b \end{bmatrix}$  (i.e.  $v = (b \ \vec{w})^T$ )

Idea: encourage  $\vec{g}(\vec{x}^T \vec{v})$  to be far from 0 if  $y=1$ .

← " " " " " 1 " y=0.  
 penalize  $\vec{v}$  if  $y=1$  and  $-\log(g(\vec{x}^T \vec{v}))$  is large  
 | log(x) " " " y=0 " "  $-\log(1-g(\vec{x}^T \vec{v}))$  is large

$$\log(x) = \ln(x)$$

$$\log = \ln$$

combine the two cases into one big function

solution: define the cost function

$$J(v) = \frac{1}{m} \sum_i$$

Given dataset

$$x, y = x^{(\text{train})}, y^{(\text{train})}$$

\*  $x^{(test)}, y^{(test)}$

$$-(-y_i) \log(1 - g(\hat{X}_{i,*}^T \vec{v})) - y_i \log(g(\hat{X}_{i,*}^T \vec{v}))$$

the test data

$$= \frac{1}{m} \left[ -(\vec{t} - \vec{y})^T \log (\vec{t} - G(\hat{\vec{x}}_V)) - \vec{y}^T \log (G(\hat{\vec{x}}_V)) \right]$$

where  $\text{Log}: \mathbb{R}^m \rightarrow \mathbb{R}^m$   $\text{Log}(x_1, \dots, x_m) = (\log(x_1), \dots, \log(x_m))$ .



扫描全能王 创建

$$G: \mathbb{R}^m \rightarrow \mathbb{R}^m \quad G(x_1, \dots, x_m) = (g(x_1), \dots, g(x_m))$$

and find  $v^* = \arg \min_v J(v)$  using some form of gradient descent

proposition: ①  $\nabla_v J(v) = \frac{1}{m} (G(\hat{x}v) - \hat{y})^T \hat{x}$

$$H(J)(v) = \frac{1}{m} \hat{x}^T D(v) \hat{x}$$

$$D(v) = \text{diag}(G'(\hat{x}v))$$

$$G'(x_1, \dots, x_m) = (g'(x_1), \dots, g'(x_m))$$

②  $J$  is twice differentiable convex, and has Lipschitz continuous gradient.

Proof: HW2

∴ Gradient descent reliably finds the global min (if it exists)

case ②: Expect a nonlinear "decision boundary"

solution: modify  $x$  by adding new columns corresponding to polynomial or other functions of current columns.

### 9.5.2 Capacity, overfitting, underfitting

Recall: Learning algorithms (E.T.P. program)

the program is "trained" using  $E$  (e.g.  $x^{(\text{train})}, y^{(\text{train})}$ )

but  $P$  is determined by "test" data (not in  $E$ ).

(e.g.  $x^{(\text{test})}, y^{(\text{test})}$ ) for supervised algo.

usually have some "cost function" and we minimize them w.r.t training set to determine program.

That is, minimize "test error" by minimizing "training error".

Question: How can we statistically guarantee this approach works?



扫描全能王 创建

independent identically  
distributed

Need to make assumptions (i.r.o. assumptions):

Both the training set and test set are gotten by independently randomly sampling the same probability distribution.  
(call this  $P_{data}$ ).

Usually our task  $T$  is accomplished by choosing some set  $\vec{w}$  of parameters found by minimizing some cost function  $J(\vec{w})$  on training set.

E.g.



扫描全能王 创建