

Intro to ML by Tom Mitchell

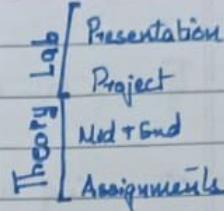
ML Murphy

Page No.:

Date: / /

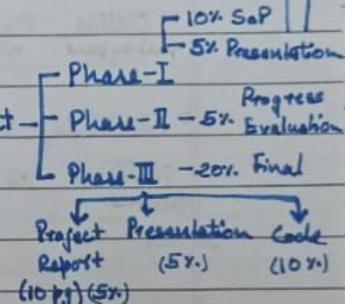
→ Classical ML - Statistical
→ DL

Gen AI



QJ Journal for CSE

A/A* main conference → article related to opted project
L article per member (at least)



Theory (60%)

- Midsem - 20%
- Endsem - 30%
- Assignments - 10%
(Best 2/3)

Lab (40%)

- New/ Unsolved Problem
- Refractoring the code ← (already solved but no code)
- Build new architecture

Task T

Performance P
Data/Esp E

⇒ Learner System

- Learner - Takes esp and bg → builds model
- Reasoner - works with model
 - provides sol \approx
 - also perf. measure

Creating Learner

- choose training experience — Features
- choose target function — Class of functions
- choose how to repr. target fx \approx — Designing of learning algo
- choose a learning algo to infer target fx \approx

choose esp fx etc.

target fx \approx → problem type (classification, regression etc.)

$$\hookrightarrow y = f(x) \quad x \rightarrow \text{feature}$$

$y \rightarrow$ actual expected op

Unsupervised Learning

→ No examples (only)

Semi-supervised Learning

→ few labelled many unlabelled

Reinforcement Learning

→ policy, agent, environment, reward-penalty

Online v Batch

↳ use all to train in one go (offline)

↳ (new+new) + new + ...

↳ incremental

Instance based v Model based

(categorization technique)

↳ measure of similarity

↳ Build a model for prediction

fitness or
Utility metrics — How good

Cost function — How bad

Types of ML Systems:

→ On basis of human supervision

Supervised, Unsupervised, Semi-supervised, Reinforcement
 (labeled) unlabeled most unlabeled reward-punishment

→ incremental or learning on fly

Online, Batch
 part by part all at once

→ comparing to known data pts or by detecting patterns

instance vs model based
 distance \Rightarrow model

Supervised Learning

Given:
 Set of input features x_i
 target features y_i
 set of input-target features
 New examples

Output

Prediction

Classification - discrete categories

Regression - continuous

Features : attr.

Feature vector : attr vec for numericals

Instance Space X : set of all possible objs that can be defined by features

Example (x_i, y_i)

Concept (c) : subset of X

Target fct \rightarrow : Map each $x \in X$ to $y \in Y$

Training Data : observed egs

Challenges

Insufficient training Data

- overfitting
- capturing noise rather than features
- would perform poorly on new data

Non representative Training Data

- inaccurate reflection of distribution and density
- poor generalization

Poor Quality Data

- Data contains errors outliers and missing values
- inaccurate predictions

Irrelevant Features

- including irrelevant features increases complexity and can lead to overfitting
- fits noise rather than meaningful patterns
- more resources required

Overfitting Training Data

- captures noise and random fluctuations too well that do not generalize to new data

Underfitting

- too simplistic to capture the patterns

$$E(\gamma - \hat{\gamma})^2 = E[(f(x) - \hat{f}(x) + \epsilon)^2] = E[f(x) - \hat{f}(x)]^2 + \text{Var}(\epsilon)$$

Page No.:

Date: / /

Bias

$$\sum_{i=1}^N f(x_i) - \hat{y}_i = \text{EPE} \quad (\text{Expected prediction error})$$

True function ↗ prediction function

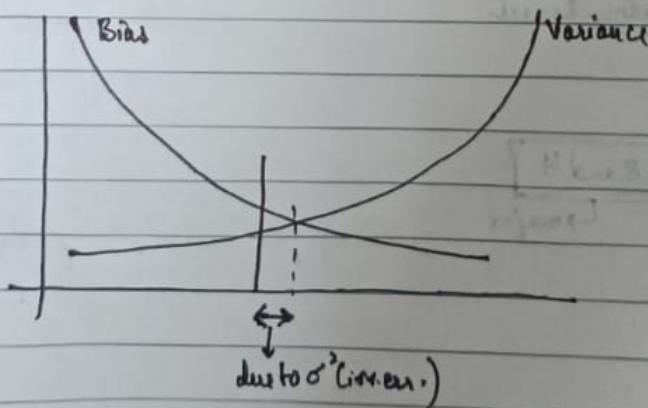
$$\begin{aligned} &= E[(f(x_i) - \hat{y}_i)^2] \\ &= E[f(x_i)^2] + E[\hat{y}_i^2] - 2 E[f(x_i) \hat{y}_i] \\ &= f(x_i)^2 + E[\hat{f}(x_i)^2] - 2 f(x_i) E[\hat{f}(x_i)] \\ &= [f(x) - E[\hat{f}(x)]]^2 + E[(\hat{f}(x) - E[\hat{f}(x)])^2] + \sigma^2 \end{aligned}$$

Bias
deviation from avg. prediction (from train) of actual values

Variance
deviation from avg. prediction (from test/pred.) of pred. val.

Irreducible Error
Error that can not be avoided
due to faults in data collection (e.g.: faulty/noisy sensor)

Bias Variance Tradeoff



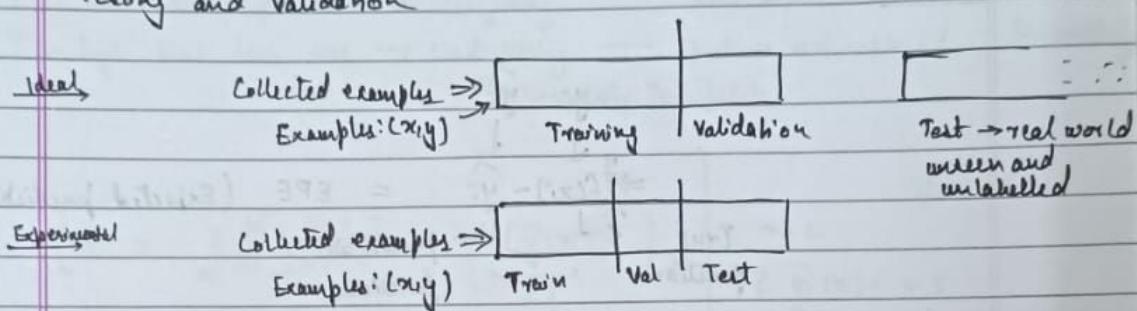
point chosen depends on metric observed (prec. recall acc.)

same \Rightarrow model type (only parameter/hyp diff) \rightarrow pool
diff model type \rightarrow acc, prec, recall

Page No.:

Date: / /

Testing and Validation



$$\text{Generalization error} = \text{Error}_{\text{test}} - \text{Error}_{\text{train}}$$

Gen err less \rightarrow accurate
more \rightarrow over/under

Hyper-parameter tuning

Importance of features \rightarrow parameters
(epochs, batchsize, ...) Training \rightarrow hyperparameters
important for initiating process

- Manual Search
- Grid Search
- Random Search

1, 2, 3 and 4
↳ major

No Free Lunch Theorem | Inductive Bias | PAC Learning

Page No.:

Date: / /

(PAC - Probability Abstraction Current)

\Rightarrow Regularization, Feature Selection, Dropout

- low bias high var \rightarrow Overfitting — noise captured — high complexity
- high bias low var \rightarrow Underfitting — patterns not captured — low complexity

\Rightarrow Complexity, Data points

$$\text{Var} = \frac{1}{N} \sum_{i=1}^N \frac{1}{M} \sum_{j=1}^M \left(\hat{y}_j(x_i) - \underbrace{\mathbb{E}[\hat{y}(x_i)]}_{\text{mean model}} \right)^2$$

multiple instances \Rightarrow multiple model \Rightarrow mean model pred.

 $\text{Var} = 0.167$
 $\text{Bias}^2 = 0$

$$f(x) = 5$$

$$n = 2$$

$$\beta \hat{y}_1(x) = 4 \cdot 5$$

$$\hat{y}_2(x) = 5$$

$$\hat{y}_3(x) = 6 \cdot 5$$

~~cont~~ can't say over/under without var range
what if var is considered not that high?

Hypothesis Space

\rightarrow all possible hypotheses

VC dimensions

\rightarrow how effectively we can select hypothesis based on data.

generalization bound

~~tree exp~~

Theoretical framework for evaluation

~~tree exp~~

Inductive bias — Assuming (assumptions) \rightarrow Effectiveness \rightarrow PAC Learning

Deductive bias — Inferring (deductions/inferences)

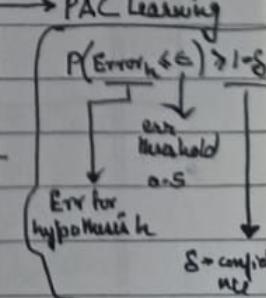
not accepting Shattering
at: if null is rejected

\leftarrow (Vapnik-Chervonenkis)

VC dimension
If hyp is complex enough but a simple
hyp. has also the same result, pick
simple one.

Theory

VC dimension



Tom Mitchell book

Page No.:

Date: / /

consistent hypothesis - remains same

Eg for Find-s

$f_1 \quad f_2 \quad f_3$

red round small + \rightarrow add to hypothesis

Yellow long medium - \rightarrow check and modify so that - we are not
in hypothesis

Green round medium + \rightarrow add to hypothesis

Initially

$s_0 \quad \{ \phi, \phi, \phi \}$

$s_1 \quad \{ \text{red, round, small} \}$

$s_2 \quad \{ \text{SL} \}$

$s_3 \quad \{ \text{red, green} \} \{ ?, ?, ?, ? \}$

$\Rightarrow s_3$ is $c(x)$

Voron space?

+ w sample \rightarrow minimal generalization from specific boundary

Eg for Candidate Elimination

- w sample \rightarrow minimal specialization from General boundary

check consistency

Eg

Sky	AirTemp	Humidity	Wind	Water	Forecast	PlaySport
Sunny	Warm	Normal	Strong	Warm	Same	Yes
"	"	High	"	"	"	"
Rainy	Cold	"	"	"	Change	No
Sunny	Warm	"	"	Weak	"	Yes

$s_0 < \phi \phi \phi \phi \phi \phi >$

$g_0 < ??? ??? >$

$s_1 < \text{Sunny, warm, normal, strong, warm, same} > g_1 = g_0$

$s_2 < \text{Sunny, Warm, ?, strong, weak, same} > g_{1,2} = g_1$

$s_3 = s_2$

$g_3 = \left\langle \begin{array}{l} \text{Sunny} \\ \text{Same} \end{array} \right\rangle g_3 = \text{Sunny}$

(rule-taking v rule-making) and context-based validity

Add more questions here

Page No.

7/22

Date:

Today

Justify role in SoP

Generalization

$n \rightarrow$ no. of data points.

$VC(H) \rightarrow$ VC dim. of
hyp. space

⇒ Generalization Bound

$$\text{Err}_{\text{true}} \leq \text{Error}_{\text{Emp.}} + \sqrt{\frac{VC(H) \log(n)}{n}}$$

If $VC(H) \rightarrow$ high $\rightarrow \text{Error}_{\text{true}} \rightarrow$ high

If $n \rightarrow$ high $\rightarrow \text{Error}_{\text{true}} \rightarrow$ low

SoP

- PS

- Related Work

- Solution - Flow diagram with phase divisions and work-member div

Concept Learning

→ Binary \rightarrow one = concept

→ Find-S algorithm

→ Candidate Elimination Algorithm

Dataset $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ $n \in \{1, 1000\}$
 $y_i \in \{\text{True, False}\}$

$$y_i = c(x_i)$$

$$\Rightarrow D = \{(x_1, c(x_1)), (x_2, c(x_2)), \dots\}$$

Goal to find hypothesis $h(x) = (x_1 \wedge x_2 \wedge \dots \wedge x_n)$

If all have true for an attr - keep
else not-consistent

Page No.: _____

$$S_3 = S_2 \Rightarrow G_{s_3} = \langle \text{Sunny}, ??? ? \rangle / \langle ?, \text{warm}, ??? ? \rangle / \langle ??? \text{ Normal } ?? \rangle \\ \langle ??? \text{ Cool } ? \rangle / \langle ??? ? \text{ Same } \rangle$$

$$\Rightarrow S_4 = \langle \text{sunny}, \text{warm}, ?, \text{strong}, ?, ? \rangle \quad G_4 = G_3$$

Version space \rightarrow space having all possible hypotheses

$\langle \text{Sunny ? ? strong ??} \rangle$	
$\langle \text{Sunny warm ??? ?} \rangle$	$\rightarrow V =3$
$\langle ? \text{ warm ? strong ??} \rangle$	

PS - why are we not taking 3 and
↳ in each hyp.

Drawbacks → not scalable for larger datasets
→ order of hypothesis matters

Decision Tree

- root
 - internal node
 - leaves

\Rightarrow for supervised learning

→ recursively partition the training data into homogeneous regions

for classification

→ less time taking and less resource heavy in test time

→ Size and shape

- internal and leaf node \neq
- branching factor
- depth

→ cross validation for best size & shape

→ At leaf

- Gini constant δ/p
- use KNN
- use other supervised learning model

⇒ How to split at Internal Nodes

- split should result in as pure groups as possible
- for classification, entropy is a measure
 - ↳ low entropy \rightarrow high purity

Technique to split

- optimality depends on test / case-by-case basis
- greedy approach is good in some cases
- information gain → informative rules?

Leaves should give homogeneous/informative δ/p

→ highest information gain \rightarrow root

Entropy and Information Gain

$S \rightarrow$ set of labelled inputs S from C classes, p_c as fraction of c class inputs

$$\text{Entropy } H(S) = - \sum_{c \in C} p_c \log p_c$$

$S \rightarrow$ split $\longrightarrow S_1, S_2$ (disjoint)

$$\begin{aligned} \text{Reduction in Entropy after split} & \uparrow I_G = H(S) - \frac{|S_1|}{|S|} H(S) - \frac{|S_2|}{|S|} H(S_2) \\ & \downarrow \text{Information gain} \end{aligned}$$

$H_S \rightarrow$ true, -ve classes in C

$$\Rightarrow H(S) = (\text{true } p_c) \log(\text{true } p_c) + (-\text{ve } p_c) \log(-\text{ve } p_c)$$

~~outlook~~

$$\cancel{\text{outlook}} \quad 0.94$$

$$S_3 = \begin{array}{ccccc} \text{sunny} & & \text{cloudy} & & \text{rainy} \\ \downarrow & & \downarrow & & \downarrow \\ 3- & 2+ & 4+ & 0- & 3+ \\ & & & & 2- \end{array}$$

$$H(S_3) = - \left(\frac{3}{5} \log \frac{3}{5} + \frac{2}{5} \log \frac{2}{5} \right) = \frac{1}{5} \left(3 \log 3 + 2 \log 2 - 5 \log 5 \right) = \left(3 \log \frac{3}{5} + 2 - 5 \log 5 \right)$$

$$H(S_0) = - \left(\log 1 \right) = 0$$

$$H(S_{\text{rainy}}) = - \left(\frac{3}{3} \log \frac{3}{3} + \frac{2}{3} \log \frac{2}{3} \right) = - (3 \log 3 + 2 - 5 \log 5)$$

$$I_G = 0.94 - \frac{5}{14} () - \frac{5}{14} = 0.94 - \frac{5}{14} \left[- (6 \log 3 + 2 - 10 \log 5) \right]$$

Brennan Leo Friedman, JH, OMSher

Temperature 4 h 6 m 4 c

$$\begin{array}{c}
 \text{28} \downarrow \quad \downarrow \quad \downarrow \\
 \text{2-} \quad \text{2+} \quad \text{4+} \quad \text{2-} \quad \text{3+} \quad \text{1-} \\
 \hline
 \downarrow \quad \downarrow \quad \downarrow \quad \downarrow \\
 \text{-} \left(\frac{1}{3} \log \frac{1}{3} + \frac{2}{3} \log \frac{2}{3} \right)
 \end{array}
 \quad \frac{1}{4} \log \frac{1}{4} + \frac{3}{4} \log$$

When to stop

- overfitting
- max entropy
- covered all features

避免过拟合 Avoiding Overfitting in DTs

- decision stump
- pre pruning and post pruning → based on purity / entropy
- feature selection

Notes

- Gini index instead of Ig
- variance can be used to assess purity
- for real valued features → use tests based on thresholding feat. values

Decision Tree for Regression

Instance Based Learning
+ Lazy / Memory based

Tom Mitchell

⇒ K-Nearst Neighbour

$$d(x_i, x_j) = \sqrt{\sum_{i=1}^n (a_i(x_i) - a_i(x_j))^2} \quad a_i \rightarrow \text{1}^m \text{ features}$$

Voronoi Diagram -

⇒ Training Algorithm

For each training algo $\langle x, f(x) \rangle$ add the example to the list of training

⇒ Classification Algorithm

$$\hat{f}(x_j) \leftarrow \arg\max_{v \in V} \sum_{i=1}^K \delta(v, f(x_i))$$

$$\delta(a, b) = \begin{cases} 1 & \text{if } a=b \\ 0 & \text{otherwise} \end{cases}$$

⇒ For regression

$$\hat{f}(x_j) = \frac{\sum_{i=1}^K f(x_i)}{K}$$

→ for distance weighted, mult $w_i = \frac{1}{d(x_j, x_i)^2}$ to δ value of $f(x_i)$

left null space?

Linear Regression

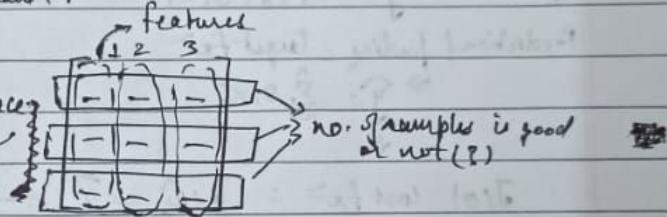
Stat → ISL

Prob → Murphy

Interp → Hands-On

PREREQUISITES

- Determinant
- Inversion → inverse opp. vector
- Vector Space →
- Linear Transformation
- Null space (Rank Nullity Thm) (relate with $Ax=0$) (cone w/o relation)
- Metric Space and Normal Space
- Norm $\| \cdot \|_1, \| \cdot \|_2, \dots, \| \cdot \|_\infty$
- Convex function \cup and Concave \cap
- Multicollinearity
- Inner product and Inner prod space
 ↓ dot b/w vectors



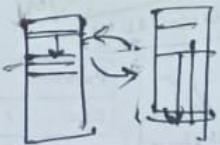
- Orthogonal → inner product is zero
- Orthogonal complement → If W is inner product space then W^\perp if elements of W

$$W^\perp = \{v \in V \mid v \perp w \ \forall w \in W\}$$

W^\perp are orthogonal to every element of W

- Eigenvalue, Eigenvector → search in limited space

- When to take mean, median, mode
- IQR, variance range → outlier detection
- Random Variable
- distribution
- pdf, pmf (interpretation of PDF($p(x=4)=10\%$))
- covariance



V1.

Page No.:

Date: / /

Linear Regression

Why linear

$$y = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n$$

↑
bias (inductive)

↑
Linearity

(for 1 ~~feature~~ feature)

Hypothesis $\rightarrow \hat{y} = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n$ (for n-features)

Prediction/finding target $f(x) \rightarrow$

$$\Rightarrow \hat{y} = \sum_{i=0}^n \theta_i x_i$$

↓

$J(\theta)$ cost function: $y - \hat{y} = \tilde{r}$ \leftarrow i (residuals)

$$\frac{1}{n} \sum_{i=0}^n (\hat{y}_i - y_i)^2 \quad \text{(ii)} \quad \equiv r = \arg \min_{\theta} J(\theta) = \frac{1}{n} \sum (\theta_i x_i - y_i)^2$$

$$\frac{1}{n} \sum (y_i - \hat{y}_i)^2 \quad \text{(iii)}$$

$$\frac{1}{n} \sum \sqrt{(\hat{y}_i - y_i)^2} \quad \text{(iv)}$$

$$h_{\theta}(x) = \sum_{i=0}^n \theta_i x_i$$

Goal: $\arg \min_{\theta} J(\theta) \rightarrow \hat{\theta}$ (optimal value of θ for our hypothesis)

If $x \in \mathbb{R}^d$ $\theta \in \mathbb{R}^{d+1}$

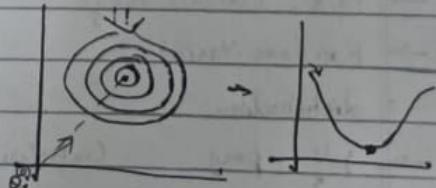
How to get the min?

↪ Gradient Descent Approach

$$\theta^{(0)} = \text{init}$$

for all $j \in \{1, \dots, d+1\}$

$$\theta_j^{(t+1)} = \theta_j^{(t)} - \alpha \frac{\partial J(\theta)}{\partial \theta_j} \quad i = \{0, d\}$$



$$J(\theta) = \frac{1}{2} \sum (\theta^T x - y)^T (\theta^T x - y)$$

$$\text{Goal } \nabla J(\theta) = 0$$

$$\Rightarrow \frac{1}{2} \sum (\theta^T x - y)^T (\theta^T x - y) = 0$$

$$\Rightarrow (\theta^T x - y)^T (\theta^T x - y) = 0$$

$$\Rightarrow (\theta^T x)^T \theta^T x - (\theta^T x)^T y - y^T (\theta^T x) + y^T y = 0$$

\Rightarrow

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n (\theta^T x_i - y_i)^2$$

$x = R^{n \times d+1}$ (considering θ_0, x_0)

$\theta = R^{d+1}$

$y = R^n$

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n (\theta^T x_i - y_i)^2$$

$$= \frac{1}{2} \left[(x_0^T \theta - y_0)^2 - (x_1^T \theta - y_1)^2 - \dots - (x_n^T \theta - y_n)^2 \right]$$

$$= \frac{1}{2} \left[\theta^T (x^T x) \theta - 2 \theta^T (x^T y) + y^T y \right]$$

$$\Rightarrow \frac{1}{2} \theta^T \left[\theta^T (x^T x) \theta - 2 \theta^T (x^T y) \right] = 0$$

$\Rightarrow \theta \cancel{\theta}$

$$\Rightarrow \theta = ((x^T x)^{-1} x^T y)$$

$$(x^T x) \theta = (x^T y) \Rightarrow \text{Normal Eq.}$$

Closed form - all given data
 GRD \rightarrow real time data.

Page No.:

Date: / /

repeat till convergence

$$\begin{aligned} &\rightarrow \text{No change in } \theta \Rightarrow \|\theta^{(t)} - \theta^{(t-1)}\| \\ &\rightarrow \text{No change in } J(\theta) \Rightarrow \|J(\theta^{(t)}) - J(\theta^{(t-1)})\| \\ &\rightarrow \text{If } \frac{\partial J(\theta)}{\partial \theta_j} \text{ is too small} \end{aligned}$$

$$\begin{aligned} \theta^{t+1} &= \theta^t - \alpha \nabla_{\theta} J(\theta^t) \\ &= \theta^t - \alpha \nabla_{\theta} \left[\frac{1}{2} \sum_{i=1}^n (\theta^T x^i - y^i)^2 \right] \\ &= \theta^t - \alpha \nabla_{\theta} \left[\frac{1}{2} \sum_{i=1}^n (\theta^T x^i - y^i)^2 \right] \\ \Rightarrow \theta^{t+1} &= \theta^t - \alpha \left[\sum_{i=1}^n (\theta^T x^i - y^i) x^i \right] \end{aligned}$$

$\theta = \theta_{\text{init}}$

\rightarrow BGD

for i in n -samples:

$\theta = \theta_{\text{init}}$ \rightarrow SGD

for j in d -features:

$$\theta_j \leftarrow \nabla_{\theta_j} J(\theta_j) \rightarrow \text{SGD}$$

$$\text{all}(\nabla_{\theta_j} J(\theta_j)) \rightarrow \text{accumulate BGD}$$

$$\theta \leftarrow \text{all}(\nabla_{\theta_j} J(\theta_j)) \rightarrow \text{BGD}$$

— BGD converges faster but consumes more memory

— SGD is less resource heavy but converges later

— miniBatch \rightarrow Best of both worlds \rightarrow sufficient resources and decent compute time required.

\rightarrow Closed Form Solution

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n (\theta^T x^i - y^i)^2$$

$$\theta \in \mathbb{R}^{nd+1}$$

$$x \in \mathbb{R}^{nd+1}$$

$$y \in \mathbb{R}^n$$

→ Assumption

$\Rightarrow \mathbf{x}^T \mathbf{x}$ is invertible \Rightarrow non-multicollinearity \Rightarrow linearly independent feature set
 \Rightarrow i.i.d

Probability \rightsquigarrow Likelihood

|
 ↳ known value, unknown parameter
 ↳ known parameters, known distribution

$$y^{(i)} = \boldsymbol{\theta}^T \mathbf{x}^{(i)} + \epsilon^{(i)} \quad \epsilon^{(i)} \sim N(0, \sigma^2) \quad i \in \text{iid}$$

$$\epsilon^{(i)} = y^{(i)} - \boldsymbol{\theta}^T \mathbf{x}^{(i)} \quad \sim N(0, \sigma^2)$$

$$\Rightarrow y^{(i)} \sim N(\mu(\boldsymbol{\theta}^T \mathbf{x}^{(i)}), \sigma^2)$$

$$p(y^{(i)} | \mathbf{x}^{(i)}, \boldsymbol{\theta}) \Rightarrow p_{\mathcal{Y}^N} = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y^{(i)} - \mu(\boldsymbol{\theta}^T \mathbf{x}^{(i)}))^2}{2\sigma^2}\right)$$

$$L(\boldsymbol{\theta}) \Rightarrow \prod_{i=1}^n p(y^{(i)} | \mathbf{x}^{(i)}, \boldsymbol{\theta})$$

$$= \frac{1}{(\sqrt{2\pi\sigma^2})^n} \exp\left[-\frac{1}{2\sigma^2} \sum_{i=1}^n (y^{(i)} - \boldsymbol{\theta}^T \mathbf{x}^{(i)})^2\right]$$

$$\log(L(\boldsymbol{\theta})) = L(\boldsymbol{\theta}) = \frac{-n}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (y^{(i)} - \boldsymbol{\theta}^T \mathbf{x}^{(i)})^2$$

$$L(\boldsymbol{\theta}) = n - \underbrace{\frac{1}{\sigma^2} \sum_{i=1}^n (y^{(i)} - \boldsymbol{\theta}^T \mathbf{x}^{(i)})^2}_{\text{loss function}}$$

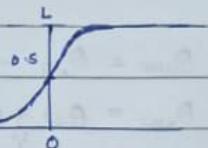
$$\Rightarrow \underset{\boldsymbol{\theta}}{\operatorname{argmax}} L(\boldsymbol{\theta}) = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} J(\boldsymbol{\theta})$$

New g choice

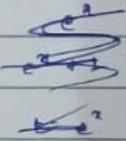
logistic / sigmoid

$$g(z) = \frac{1}{1 + e^{-z}} \quad z = \theta^T x$$

$$\lim_{z \rightarrow -\infty} = 0 \quad \lim_{z \rightarrow 0} = 0.5 \quad \lim_{z \rightarrow \infty} = 1$$



$$\begin{cases} 1 & \text{if } g(z) \geq 0.5 \\ 0 & \text{if } g(z) < 0.5 \end{cases}$$



$$h_\theta(x) = g(\theta^T x)$$

$$g(z) = \frac{1}{1 + e^{-z}}$$

$$\begin{aligned} & (1+e^{-z})^{-1} \\ & -1(1+e^{-z})(1+e^{-z}) \\ & -e^{-z} \\ & (1+e^{-z})^2 \end{aligned}$$

$$\begin{aligned} p(y^{(i)}=1|x, \theta) &= h_\theta(x) \\ p(y^{(i)}=0|x, \theta) &= 1 - h_\theta(x) \end{aligned}$$

$$p(y|x, \theta) = (h_\theta(x))^y \cdot (1 - h_\theta(x))^{1-y}$$

$$g(z)^y = g(z)(1-g(z))^{1-y}$$

$$L(\theta) = \prod_{i=1}^n p(y^{(i)}|x^{(i)}, \theta)$$

$$l(\theta) = \sum_{i=1}^n y^{(i)} \log(h_\theta(x^{(i)})) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)}))$$

Gradient Descent

$$\nabla_\theta l(\theta) = \sum_{i=1}^n \frac{y^{(i)}}{h_\theta(x^{(i)})} \cdot \frac{\partial h_\theta(x^{(i)})}{\partial \theta} \leftarrow \left(\frac{1-y^{(i)}}{1-h_\theta(x^{(i)})} \frac{\partial h_\theta(x^{(i)})}{\partial \theta} \right)$$

For single eq:-

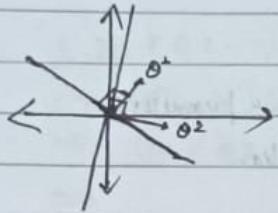
$$= \sum \frac{y^{(i)}}{g(\theta^T x)} g'(\theta^T x) x - \frac{(1-y^{(i)})}{1-g(\theta^T x)} g'(\theta^T x) x$$

$$\Rightarrow \frac{y}{g(\theta^T x)} g(\theta^T x) g(1-g(\theta^T x)) x + \frac{(1-y)(-1)}{1-g(\theta^T x)} g(\theta^T x) (1-g(\theta^T x)) x$$

$$\Rightarrow x [y(1-g(\theta^T x)) - (1-y)g(\theta^T x)] = x [y - y \cdot g(\theta^T x) - g(\theta^T x) + y g(\theta^T x)]$$

Perception Algorithm

$$y^{(i)} = g(\theta^T x^{(i)}) \quad y^{(i)} \in \{0, 1\}$$



$$g \rightarrow \text{originally} \quad \begin{cases} 1 & \theta^T x > 0 \\ 0 & \theta^T x \leq 0 \end{cases}$$

$$g \rightarrow \text{better choice} \quad \frac{1}{1 + e^{-\theta^T x}} \quad (\text{logistic})$$

Prereq.

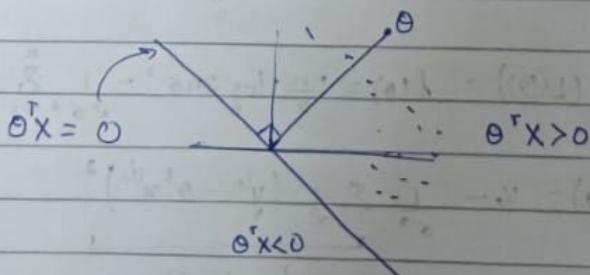
- Normal Vector
- Dataset

$$\begin{aligned} y &= \theta_0 + \theta_1 x \\ y &= \theta^T X \end{aligned} \quad \rightarrow \text{continuous valued regression}$$

$$Y = g(\theta^T X)$$

$\rightarrow g = \text{step function}$

$$\begin{cases} 1 & g(l) > 0 \\ 0 & g(l) \leq 0 \end{cases}$$



→ Step func. is very dependent on θ

Higher order norms are not that beneficial

Page No.:

Date: / /

$\Rightarrow \theta_0$ is not regularized

Ridge loss fn

$$J(\theta) = \text{MSE}(\theta) + \alpha \cdot \frac{1}{2} \cdot \sum_{i=1}^n \theta_i^2$$

Lasso Regression

Least Absolute Shrinkage and Selection Operator

$$J(\theta) = \text{MSE}(\theta) + \alpha \cdot \sum_{i=1}^n |\theta_i|$$

Elastic Net

$$J(\theta) = \text{MSE}(\theta) + r \cdot \alpha \cdot \sum_{i=1}^n |\theta_i| + \frac{1-r}{2} \alpha \sum_{i=1}^n \theta_i^2$$

lars ridge
hyperparameters

and Elastic net

→ How and Why does <which norm makes θ and lower wt> do what it does

→ Create the diag diagram for Elastic Net, L1 and L2

→ Why logistic function? (for classification) (why not others)

$$\frac{1}{1+e^{-x}}$$

(Hint: exponential eq \cong family)

→ What if No IID in loss functions

→ One SVML question

Newton method → more acc/optimal than GD

Swin-pres
Swin-Q/A
Tue/Fri
11:00

Page No.:

Date: / /

PS, ref article, flow diag/arch

$$\nabla_{\theta} = x(y - g(\theta^T x))$$

$$\nabla_{\theta} = \underbrace{x(y - g(\theta^T x))}_{\text{to modify impact on positive}}$$

$$\Rightarrow \theta_{\text{new}} = \theta_{\text{old}} - \eta \nabla_{\theta} l(\theta)$$

$$\Rightarrow \theta_{\text{new}} = \theta_{\text{old}} + \eta [x(y - g(\theta_{\text{old}}^T x))] \rightarrow \text{to modify impact on positive}$$

Newton-Raphson approach → Where benefit
→ Where loss

Regularizations

→ Good way to reduce everything (reduce no complexity)

f_1 dof - 3

$\Rightarrow f_1, f_1, f_1, f_1 \rightarrow 3$ features

Lin Reg (f_1, f_1^2, f_1^3)

1 2 1² 2² 12² 12³

Ridge Regression

→ reg term $\alpha \cdot \sum_{i=1}^n \theta_i^2$ added to cost function

→ forces the learning algo to keep min wts.

primal solution
dual solution

Page No.:

Date: / /

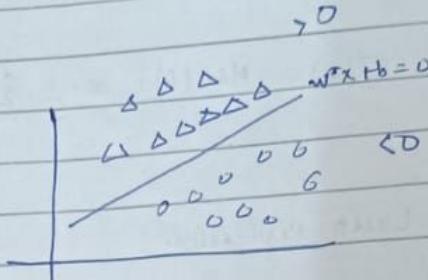
SVM

Confidence?

$$w^T x + b \geq 1$$

$$w^T x + b \leq -1$$

$$w^T x + b = 0$$



$$y_i(w^T x) \geq 1$$

$\frac{y_i}{\|w\|}$

$$\frac{|w^T x_i + b|}{\|w\|}$$

$1-y$
 y

$$\gamma = \min_{\text{pt}} (\text{Distance (pt, hyperplane)})$$

$$\underset{\mathcal{E}}{\text{obj}} \Rightarrow \max \gamma$$

$$\text{constraint: } \cancel{w^T x + b} \geq 1.$$

$$\Rightarrow \underset{w^T x + b}{\arg \max} \gamma$$

$$\text{constraint: } (w^T x + b)$$

KKT, Lagrangian multiplier

\Rightarrow GD won't work

one idea is \Rightarrow argmax $\gamma + (w^T x + b)$ \rightarrow computationally inefficient

Final obj \rightarrow maximize the min dist (pt, hyperplane / dec. surface)

$$\Rightarrow \text{dist}(x_i^+, w^T x_i + b) = \frac{|w^T x_i + b|}{\|w\|}$$

$$\Rightarrow \text{dist}(x_i^-, w^T x_i + b) = \frac{|w^T x_i^- + b|}{\|w\|}$$

$$\Rightarrow \text{dist}(x_i^-, w^T x_i + b) = \frac{|w^T x_i^- + b|}{\|w\|}$$

→ Lagrangian Method

$$L(w, b, \alpha) = \underbrace{\frac{1}{2} \|w\|^2}_{\text{min}} + \underbrace{\sum_{i=1}^n \alpha_i (1 - y_i(w^T x_i + b))}_{\max} \quad \rightarrow \text{primal form}$$

~~$$\frac{\partial L}{\partial w} = \frac{1}{2} \cdot \mathbb{E} \|w\|_2^2 + \sum_{i=1}^n \alpha_i (1 - y_i(w^T x_i + b))$$~~

$$\frac{\partial L}{\partial b} = \sum \alpha_i y_i = 0$$

$$\frac{\partial L}{\partial w} = \bar{w} - \sum \alpha_i y_i x_i = 0$$

max first min later or max later min first

difficult (?) to compute
(why?)

$$\max_{\alpha \geq 0} \left(\min_{w, b} \left(\frac{1}{2} \|w\|^2 + \sum_{i=1}^n \alpha_i (1 - y_i(w^T x_i + b)) \right) \right)$$

$$\frac{\partial L}{\partial w} = 0 \quad (\text{check})$$

$$\Rightarrow \frac{\partial L}{\partial w} \Rightarrow \cancel{\frac{1}{2} \cdot \mathbb{E} \|w\|_2^2} + \sum_{i=1}^n \alpha_i y_i x_i = w$$

$$\frac{\partial L}{\partial b} \Rightarrow \sum \alpha_i y_i = 0$$

$$\frac{\partial L}{\partial b} = 0 \quad (\text{check})$$

$$\mathcal{P}_{\text{max}} = \max_{\alpha \geq 0} \left(\frac{1}{2} \|w\|^2 + \sum \alpha_i (1 - y_i(w^T x_i + b)) \right)$$

$$= \max_{\alpha \geq 0} \left(\frac{1}{2} \|w\|^2 + \sum \alpha_i - (\sum \alpha_i y_i x_i) w^T - (\sum \alpha_i y_i) b \right)$$

=

$$L = \sum \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j (x_i^T x_j) = w(\alpha) (?)$$

→ verified

V	$\frac{\partial L}{\partial w} = w - \sum_{i=1}^n \alpha_i y_i x_i = 0$
E	$\frac{\partial L}{\partial b} = \sum_{i=1}^n \alpha_i y_i = 0$
R	.
I	.
F	.
Y	.

$$\begin{matrix} 1 & 2 & 3 \\ \frac{1}{1} & \frac{1}{2} & \frac{1}{3} \\ 0.1 & 0.2 & \underbrace{\frac{1}{0.1}}_{\frac{1}{0.2}} \end{matrix}$$

Page No.: / /
Date: / /

$$y_i(w^T x_i + b) \geq L$$

— hard margin
(no misclassification)

For +ve

$$y_i = 1$$

$$\Rightarrow w^T x_i + b \geq L$$

For -ve

$$y_i = -1$$

$$\Rightarrow w^T x_i + b \leq -L$$

$$\Rightarrow \text{dist}(x_i, \text{hyperplane}) = \frac{|y_i| w^T x_i + b|}{\|w\|}$$

— verify from
here —

$$\Rightarrow \min_i (\text{dist}(x_i, \text{hyperplane})) = \hat{\gamma}_i$$

$$\Rightarrow \gamma = \max_i (\hat{\gamma}_i) \quad \gamma \rightarrow \text{margin}$$

\downarrow functional margin

$$y_i(w^T x_i + b) = \text{dot}(\text{data}_i, \text{hyperplane})$$

\downarrow sensitive to vals of w and b

1 (? why?) [verify]

$$\hookrightarrow \gamma = \max_i \left(\frac{1}{\|w\|_2} \right)$$

\downarrow geometric margin

$$y_i(w^T x_i + b) \quad \|w\|_2$$

$$\Rightarrow \left[\gamma = \min_i \left(\frac{1}{\|w\|_2} \right); \quad y_i(w^T x_i + b) \geq L \quad (\text{for better computation}) \right]$$

Margin-width \rightarrow controlled by $C \rightarrow$ determines under/over

Page No.: / /
Date: / /

$$L = f(x) + \sum_i \lambda_i g_i(x) + \sum_j \mu_j h_j(x)$$

$$\text{New obj} \rightarrow \max_{\alpha} w(\alpha)$$

$$Q = y_i y_j (x_i^T x_j)$$

$$\left[\max_{\alpha} w(\alpha) = L - \frac{1}{2} \alpha^T Q \alpha \right] \Rightarrow \text{apply KKT if conditions are true}$$

KKT conditions

$$\hookrightarrow \text{Primal Feasibility} \quad y_i (w^T x_i + b) \geq L$$

$$\hookrightarrow \text{Dual Feasibility} \quad \alpha_i \geq 0$$

$$\hookrightarrow \text{Stationarity} \quad w = \sum_i \alpha_i y_i x_i \quad ; \quad \sum \alpha_i y_i = 0$$

$$\hookrightarrow \text{Slackness} \quad \alpha_i \cdot (y_i (w^T x_i + b) - L) = 0$$

\rightarrow applying KKT, we get

$$f(\alpha) = \text{sign} \left(\sum_i \alpha_i^* y_i (x_i^T x) + b \right)$$

For noisy data,

$$y_i (w^T x_i + b) \geq L - \epsilon \quad \epsilon > 0$$

Soft Margin
(may be misclassified)

\Rightarrow Final objective

$$\min_{w, b, \epsilon} \frac{1}{2} \|w\|^2 + C \sum_i \epsilon_i \quad \text{s.t.} \quad y_i (w^T x_i + b) \geq L - \epsilon_i$$

$$\min_{w, b, \epsilon, \mu, \alpha} \frac{1}{2} \|w\|^2 + \mu_j (C - \sum_i \epsilon_i) + \alpha_i (1 - \epsilon_i - y_i (w^T x_i + b))$$

$$\frac{\partial L}{\partial w} \quad \frac{\partial L}{\partial b} \quad \frac{\partial L}{\partial \epsilon}$$

Multiclass SVM

One v.	One vs all	K vs K
--------	------------	--------

Page No.:

Date: / /

SVM is by default only for binary classification

$$\max_w w(x) = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j (x_i^T x_j)$$

$$\text{st } 0 \leq \alpha_i \leq C, \sum_i \alpha_i y_i = 0$$

if $\alpha_i = 0 \Rightarrow$ data pts are not support vector

if $\alpha_i = C \Rightarrow$ data pts are inside the margin / misclassified

if $0 < \alpha_i < C \Rightarrow$ on the margin

SVM for Non-Linear data

Kernel trick \rightarrow project from lower to higher dimension

$$H \gg d$$

\Rightarrow polynomials in SVM $y_i(w^T x_i + b) \geq 1$

Very computationally expensive

\Rightarrow Kernel Trick \rightarrow assuming data is linearly separable in higher dimension

\hookrightarrow How to project in higher dimension

Kernel functions

- \hookrightarrow Gram Schmidt normalization
- \hookrightarrow Gaussian
- \hookrightarrow Exponential

\rightarrow Kernel functions must satisfy ~~condition~~ condition

$(x_i, x_j) \rightarrow$ positive semidefinite ~~A~~ $\Rightarrow A \geq 0$
for every datapoint

$$\text{poly} \rightarrow K(x_i, x_j) = (x_i^T x_j + c)^d$$

rbf / Gaussian / Sigmoid

Multiple correct

Fill in the blanks

numericals → scenario based

Small answers

⇒ Ensemble Learning

— Wisdom of the crowd

↳ aggregation of multiple outputs

— Different algorithms

→ Different hyper-parameters

Different representations (hypothesis)

→ Different training set

→ Why → balance bias-variance tradeoff

Bagging and Parting

五

Weak Model \rightarrow close to 50% probabilities

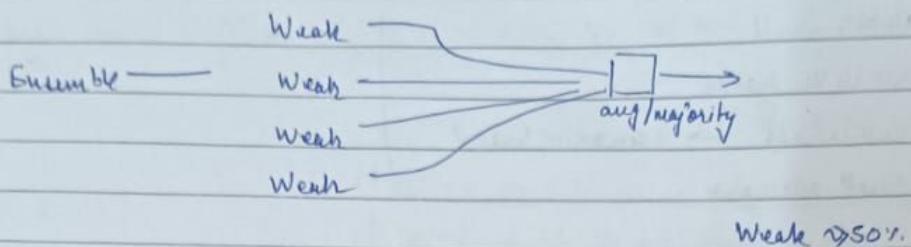
→ Bootstrapping → Picking random samples from data samples for training

→ Bagging - with replacement

→ Parting - without replacement

can be
in future
and also

→ how 63% → $n \rightarrow \infty$
 → Why not strong instead of weak



Bag

Dataset → Bagging with replacement (Bootstrapping)
 ↙ Parting without replacement

3G% → Y. Left
 DDB

63% = picked

Boosting → weighted sampling

In decision tree, ~~as~~ it is more prone to overfitting so bagging reduces variance

$$\text{Err} = (\text{Bias})^2 + \text{Var} + G$$

error $\rightarrow P < 0.5$
 Prob

$$\text{if } B \uparrow \Rightarrow P(\text{ensemble voting}) = \sum_{k=1}^B p^k (1-p)^{B-k}$$

↑
no. of learners

↓
Var ↓

↓
Lower to a point

Not effective on ~~models~~ models like linear regression, because var is already low.

For skewed datasets \rightarrow stratified Bootstrap \rightarrow ensure ~~the~~ correct ratio is picked

Random Forests

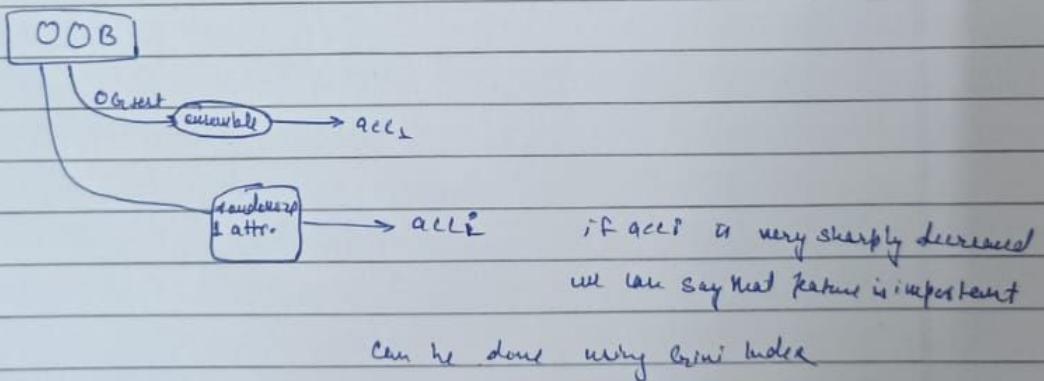
Bagging of data & features

base learners trained on a subset of features (5% is a good size)

\Rightarrow Tree is not optimal here, Ensemble is optimal

Variance Reduction split \rightarrow such that left subtree var + right subtree var is least

Feature importance



\rightarrow Slightly low interpretability of RF than DT

\rightarrow Hyperparameters $\rightarrow n, B, \text{depth}, \text{etc.}$