

Learning
Supervised: (labeled)

"right ans" given. //

- Regression. (continuous valued output)
- classifiers (discrete valued output) (0 or 1)
or more.
- Support vector machine (open with
as no. of features)

Unsupervised learning:
(Not labeled)

- clustering algs.
- coordinate parity algs.

cost fun"

$$h(\mathbf{x}) = \theta_0 + \theta_1 x$$

how to choose parameter

(optimize the prob.)

objective fun: $\min_{\theta_0, \theta_1} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$

$J(\theta_0, \theta_1)$

cost fun

or squared error fun

gradient descend algo to minimize cost fun $J(\theta_0, \theta_1)$

repeat until converging {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

(for $j=0$ & $j=1$)

learning rate.

$\alpha \uparrow \Rightarrow$ big step ↗ ↘

$$\left. \begin{array}{l} \theta_0 = \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) \\ \theta_1 = \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) \end{array} \right\}$$

$$\theta_0 = \theta_0' ;$$

$$\theta_1 = \theta_1' ;$$

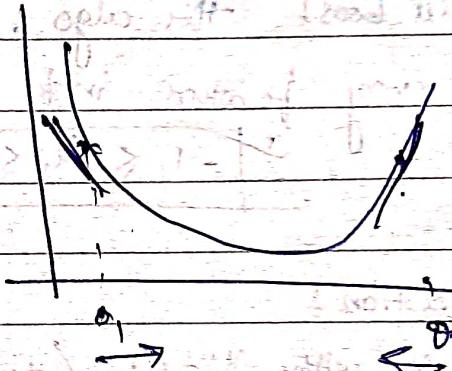
simultaneous update.
N.N. imp

$$\# \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

slope $\pm \nabla \times$ or $-\nabla \times$.

$$\theta_1 = \theta_1 - \lambda (\pm \nabla \times)$$

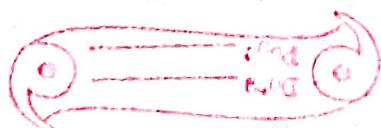
$$\Rightarrow \quad \text{always inc.}$$



$$\theta_1 = \theta_1 - \lambda (-\nabla \times)$$

$$= \theta_1 + \lambda ()$$

"Batch" gradient descent



LR (with more var) (Multivariate)

$$h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$\boxed{g_{\theta}(x) = 1}$$

$$\Rightarrow h_{\theta}(x) = \theta^T x$$

↑
↑

$$x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \quad \text{if } \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$\theta^T = [\theta_0, \theta_1, \dots, \theta_n]$$

gradient descent algo for this ($n > 1$)

Repeat {

$$\theta_j := \theta_j - \alpha \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

$$\frac{\partial}{\partial \theta_j} J(\theta)$$

$$(j = 0, \dots, n)$$

feature scaling :

(features are in same scale)

it will boost the algo.

get every feature into approximately a range,

$$\boxed{-1 \leq x_i \leq 1}$$

mean normalization :

Replace x_i with $\frac{x_i - \bar{x}_i}{\text{range}}$ (to make features have approx '0' mean)

$$x_i \leftarrow \frac{x_i - \bar{x}_i}{\text{range}} \rightarrow \text{avg. of } x$$

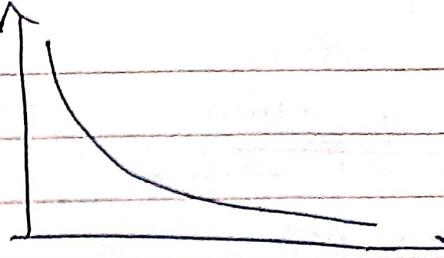
$$\underbrace{\text{range}}_{\text{max} - \text{min}}$$

$$\text{range} = (\text{max} - \text{min})$$

No. of iterations \uparrow $J(\theta) \downarrow$

\Rightarrow Gradient descent algo working properly.

Min J(θ)



No. of iteration.

Automatic Convergence:

Polynomial regression: $\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \dots$

Normal Eqn: solve θ analytically.

$$J(\theta) = \text{Cost}$$

$$J(\theta_0, \theta_1, \dots, \theta_m) = \frac{1}{2m} \sum \left(h(x^{(i)}) - y^{(i)} \right)^2$$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \dots = 0. \quad (\text{for every } j)$$

\Rightarrow solve for $\theta_0, \theta_1, \theta_2, \dots, \theta_m$

Normal eqn

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

This θ minimizes the Cost func.

(feature scaling is

not required in Normal Eqn method)

$(x^T x)^{-1}$ is non-invertible \Leftrightarrow

caz: (i) Redundant feature (i.e. LD)

(ii) Too many features (i.e. $m \leq n$)
(regularization)



Logistic Regression :

Classification prob:

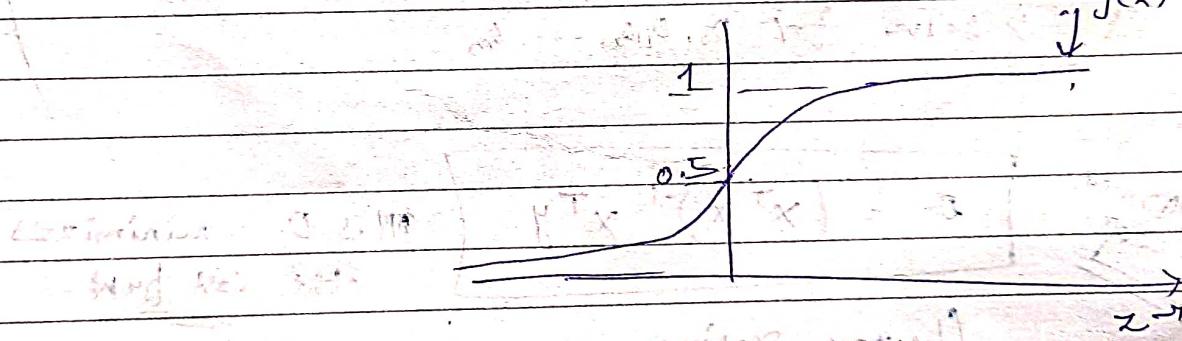
$$\text{Binary: } y = \{0, 1\} \quad 0 = \text{-ve class} \\ \quad \quad \quad 1 = \text{+ve class}$$

multiclass C.P. + g = {0, 1, 2, 3}

Logistic Regression : $0 \leq h_{\theta}(x) \leq 1$
(Classification algo)

Hypothesis Representation!

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



$$h_{\theta}(x) = P(y=1 | x; \theta)$$

Prob that $y=1$, given λ , parameterized by α .

$$P(y=0 | u; \theta) + P(y=1 | u; \theta) = 1$$

decision boundary:

$$\rightarrow y = 1 \text{ if } h_\theta(x) > 0.5$$

$\theta^T x > 0$

$$\frac{1}{1 + e^{-\theta^T x}} > \frac{1}{2}$$

$$\rightarrow y = 0 \text{ if } h_\theta(x) < 0.5$$

$$\theta^T x < 0$$

Linear DB: $h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$ $\rightarrow \theta^T x$

Non-linear decision boundary:

$$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$$

(decision boundary is a property of the hypothesis set,
not a property of the training set.)

Cost function:

(finding θ)

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}(h_\theta(x), y)$$

$$\text{cost}(h_\theta(x), y) = \frac{1}{2} (h_\theta(x) - y)^2 \rightarrow \text{Linear Regression}$$

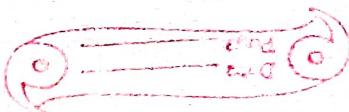
Logistic R: $h_\theta(x)$ is non-convex func.

$$\text{cost}(h_\theta(x), y) = \begin{cases} -\log(h_\theta(x)) & \text{if } y = 1 \\ -\log(1 - h_\theta(x)) & \text{if } y = 0 \end{cases}$$

Simplified C.F:

$$\text{cost}(h_\theta(x), y) = -y \log(h_\theta(x)) - (1-y) \log(1 - h_\theta(x))$$

(derived using Max Likelihood)



$$G.D \quad J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}$$

Repeat {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

}

$$\frac{\partial J(\theta)}{\partial \theta_j} = \left(h_\theta(x^{(i)}) - y^{(i)} \right) x_j^{(i)} \quad | \quad \text{same as linear,}$$

Advanced optimization :

optimization algs : \rightarrow gradient descent

$\left. \begin{array}{l} \rightarrow \text{conjugate gradient} \\ \rightarrow \text{BFGS} \\ \rightarrow \text{L-BFGS} \end{array} \right\} \text{for large problems}$

Multiclass classification : One vs all : (or one vs rest)

$$h_\theta^{(i)}(x) = P(y=i|x; \theta) \quad (i=1,2,3)$$

$$\max_i h_\theta^{(i)}(x)$$

Regularization

Prob of overfitting : (or high variance)

(underfit or high bias)

Addressing overfitting : \rightarrow more no. of features.

\rightarrow regularization.

reduce mag/val of parameters θ_j)

lest function:

Small values of parameters $\theta_0, \theta_1, \dots, \theta_n$

⇒ simpler hypothesis.

⇒ Less power to overfitting.

$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

Regularization term.
Regularization parameter.

If λ is too large,

$$\Rightarrow h_\theta(x) = \theta_0 \quad (\text{as } \theta_1 \approx \theta_2 \approx \dots \approx \theta_n \approx 0)$$

(or of underfitting)

Regularized Linear Regression

Repeat {

$$\theta_0 = \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\theta_j = \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right] \quad (j=1, 2, 3, \dots)$$

$$\theta_j = \underbrace{\theta_j \left(1 - \alpha \frac{\lambda}{m}\right)}_{\text{regularization step}} - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

$$\checkmark \left(1 - \alpha \frac{\lambda}{m}\right) < 1$$

grad descent -



$$\text{Normal Eqn} \rightarrow \mathbf{x} = \begin{bmatrix} (\mathbf{x}^{(1)})^T \\ \vdots \\ (\mathbf{x}^{(n)})^T \end{bmatrix}_{m \times (n+1)} \quad \mathbf{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix} \in \mathbb{R}^n$$

$$\min_{\theta} J(\theta) \\ \Rightarrow \theta = \left(\mathbf{x}^T \mathbf{x} + \lambda \begin{bmatrix} 0 & & & \\ 1 & 1 & & \\ & & \ddots & \\ & & & 1 \end{bmatrix} \right)^{-1} \mathbf{x}^T \mathbf{y}$$

Regularized Logistic Regression:

$$\text{Same as linear, but } h_{\theta}(x) = \frac{1}{e^{-\theta^T x} + 1}$$

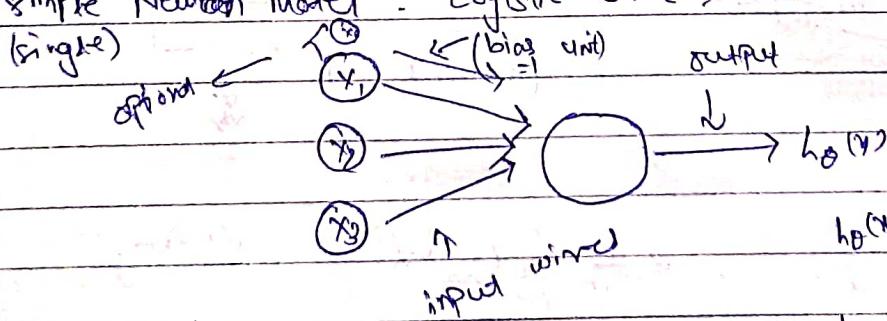
Neural Networks Representation

Non-Linear Hypothesis:

Model Representation:

In a neuron / dendrite - input wires.
Axon - output wire.

Simple Neural Model : Logistic Unit :



(sigmoid (logistic) activation function)

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

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}} \quad \mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

In Neural Network k -
(multi layers)

Weights
or parameters,

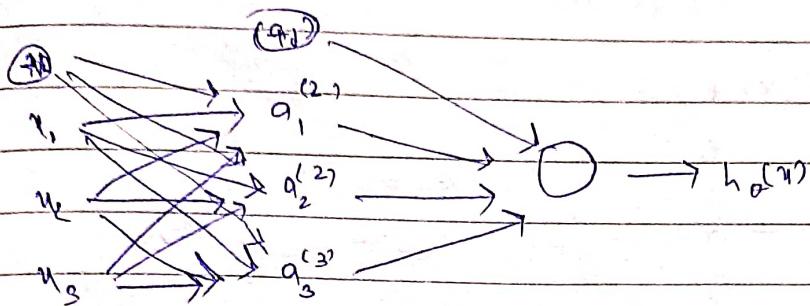
layers : Input layer : last layer : -> output layer.

Intermediate layers : hidden layer.)

$a_i^{(j)}$ = "activation" of unit "i" in layer "j".

$\theta^{(j)}$ = Matrix of weights controlling function mapping from layer j to $j+1$.

If N.N has s_j units in layer j , s_{j+1} units in layer $j+1$,
then $\theta^{(j)}$ will be of dimension $s_{j+1} \times (s_j + 1)$



$$a_1^{(2)} = g\left(\theta_{10}^{(1)} u_0 + \theta_{11}^{(1)} u_1 + \theta_{12}^{(1)} u_2 + \theta_{13}^{(1)} u_3\right) \rightarrow z_1^{(2)}$$

$$a_2^{(2)} = g\left(\theta_{20}^{(1)} u_0 + \theta_{21}^{(1)} u_1 + \theta_{22}^{(1)} u_2 + \theta_{23}^{(1)} u_3\right) \rightarrow z_2^{(2)}$$

$$a_3^{(2)} = g\left(\theta_{30}^{(1)} u_0 + \theta_{31}^{(1)} u_1 + \theta_{32}^{(1)} u_2 + \theta_{33}^{(1)} u_3\right) \rightarrow z_3^{(2)}$$

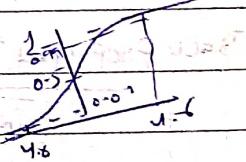
$$h_\theta(u) = g\left(\theta_{10}^{(2)} a_0^{(2)} + \theta_{11}^{(2)} a_1^{(2)} + \theta_{12}^{(2)} a_2^{(2)} + \theta_{13}^{(2)} a_3^{(2)}\right)$$

$$u = \begin{bmatrix} u_0 \\ u_1 \\ u_2 \\ u_3 \end{bmatrix} \quad X^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix}$$

$$a^{(2)} = \theta^{(1)} a^{(1)}$$

$$a^{(2)} = g(X^{(2)})$$

(forward propagation)
(multiclass classification)

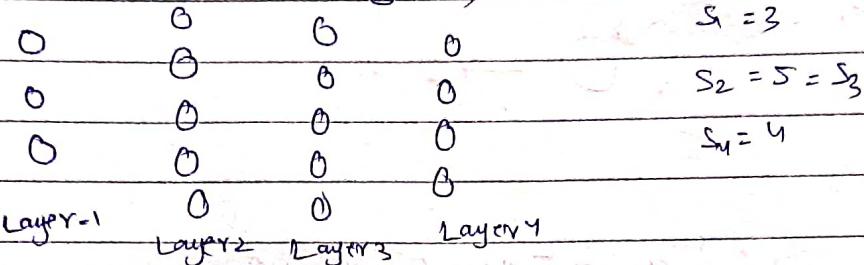


Neural Network Learning

Cost funcⁿ: (application in classification Problems)

($L = \text{No. of layers}$) $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

($S_L = \text{No. of units in layer } L$)



Binary classificⁿ:

$$y = 0 \text{ or } 1$$

(1 output unit)

$$\Rightarrow S_L = 1 \Rightarrow k = 1$$

$$h_{\theta}(x) \in \mathbb{R}$$

$$\begin{matrix} y \\ \rightarrow 0 \\ \rightarrow 1 \end{matrix}$$

Multiclay

$$y = \overline{R^k} \rightarrow$$

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 1 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

k output units

$$h_{\theta}(x) \in \mathbb{R}^k$$

$$\left| \begin{array}{c} (k=3) \\ S_L = k \end{array} \right.$$

$$h_{\theta}(x) \in \mathbb{R}^k \quad (h_{\theta}(x))_i = i^{\text{th}} \text{ output.}$$

$$J(\theta) = \frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_{\theta}(x^{(i)}))_k + (1 - y_k^{(i)}) \log(1 - h_{\theta}(x^{(i)}))_k \right] - \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{S_l} \sum_{j=1}^{S_{l+1}} (\theta_{ji}^{(l)})^2$$

Back propagation algo:

$$\begin{array}{c} \text{minimize the cost func} \\ \min J(\theta) \\ \frac{\partial}{\partial \theta_j^{(l)}} [J(\theta)] \end{array}$$

$\delta_j^{(1)}$ = "error" of node j in layer 1.

$$(L=4 \text{ in the diagram})$$

$$\delta_j^{(4)} = \hat{y}_j^{(4)} - y_j$$

$$(\hat{h}_{\theta}(x))_j$$

$$\delta^{(3)} = (\theta^{(2)})^T \delta^{(4)} + g'(z^{(3)})$$

$$\delta^{(2)} = (\theta^{(1)})^T \delta^{(3)} + \underbrace{g'(z^{(2)})}_{\text{No } \delta^{(1)}} \rightarrow \alpha^{(2)} \cdot \alpha^{(1)} (1 - \alpha^{(2)})$$

~~$\delta^{(1)}$~~ (as input layer)

$$\frac{\partial}{\partial \theta^{(l)}} J(\theta) = \alpha_j^{(l)} \delta_i^{(l+1)}$$

(ignoring λ ; if the regularization term)
 $\lambda = 0$

$\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$

→ set $\Delta_{ij}^{(1)} = 0$ (for all i, j, l)

→ for $i = 1 \dots m$

 set $a^{(1)} = x^{(i)}$

 → Perform forward propagation to compute $a^{(l)}$ for $l = 2, 3, \dots, L$

 → Using $y^{(i)}$, compute $\delta^{(L)} = a^{(L)} - y^{(i)}$

 → Compute $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$

 → $\Delta_{ij}^{(1)} = \Delta_{ij}^{(1)} + \alpha_j^{(1)} \delta_i^{(2)}$

 or $\Delta_{ij}^{(1)} = \Delta_{ij}^{(1)} + \delta^{(2)} (\alpha^{(1)})^T$

$$\Rightarrow D_{ij}^{(1)} = \frac{1}{m} \Delta_{ij}^{(1)} + \lambda \theta_{ij}^{(1)} \quad \text{if } j \neq 0$$

$$\Rightarrow D_{ij}^{(1)} = \frac{1}{m} \Delta_{ij}^{(1)} \quad \text{if } j = 0$$

$$\frac{\partial}{\partial \theta_{ij}^{(1)}} J(\theta) = D_{ij}^{(1)}$$

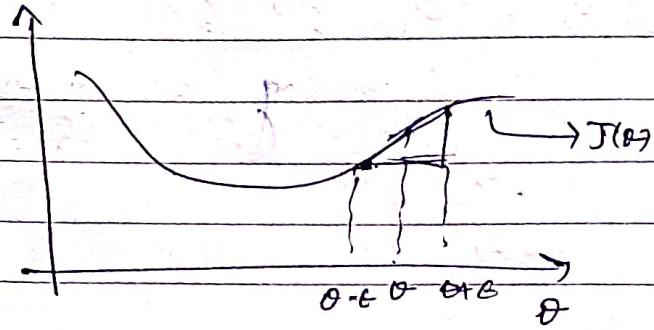


Backpropagation Intuition:

$$\delta_i^{(l)} = y_i^{(l)} - a_i^{(l)}$$

Gradient Checking:

Numerical "estimate" of gradients



$$\frac{d}{d\theta} [J(\theta)] \approx \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}$$

Random initialization:

(solving symmetric breaking issue)

Initialize $\theta_j^{(l)}$ to a random value in $[-\epsilon, \epsilon]$.

Advantages of Applying M.L.:

Machine learning diagnostic:

↓
a test

(deeper above)

Evaluating hypothesis:

→ Dataset: divide into training set & test set

↪ 1. Train on training set (70%) (30%) (around 15%)

↪ 2. Learn parameter θ from training set. (Min training error $J(\theta)$)

↪ 3. Compute test set error

4. Cross validation set.
(or validation set) (15%)

$$\rightarrow J_{\text{test}}(\theta) = \frac{1}{2M_{\text{test}}} \sum_{i=1}^{M_{\text{test}}} (h_{\theta}(x_{\text{test}}^{(i)}) - y_{\text{test}}^{(i)})^2$$

test error (for linear regression)

Similar for logistic regression (Classification)

✓ test set error using
misclassification error (0/1 misclassification error)

Model selection:

at training/Validation/Test set

Train/validation/Test error:

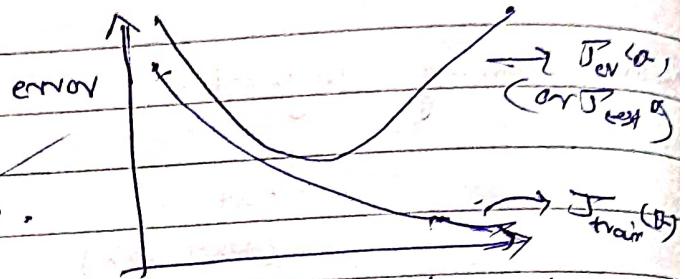
$$J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2$$

similarly $J_{\text{cv}}(\theta)$, M_{cv} , X_{cv} , y_{cv}
& $J_{\text{test}}(\theta)$, M_{test} , X_{test} , y_{test} .

J_{cv} - training error
 J_{cv} = cross validation error.

Diagnosing Bias vs Variance

Left if $J_{cv}(\theta)$
 \Rightarrow high bias (underfit)
 right \Rightarrow high var (overfit)



Bias (underfit) $\Rightarrow J_{train}(\theta) \approx J_{cv}(\theta)$
 (both high)

degree of polynomial (d).

Nonsense (overfit) $\Rightarrow J_{train}(\theta)$ is less but
 $J_{cv}(\theta)$ is high.

($J_{cv}(\theta) > J_{train}(\theta)$)

Regularization of Bias/Variance
 (help prevent overfitting)

Reg parameter λ : very high \Rightarrow underfit.
 intermediate \Rightarrow good fit.
 very low \Rightarrow overfit.

Linear
 model

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

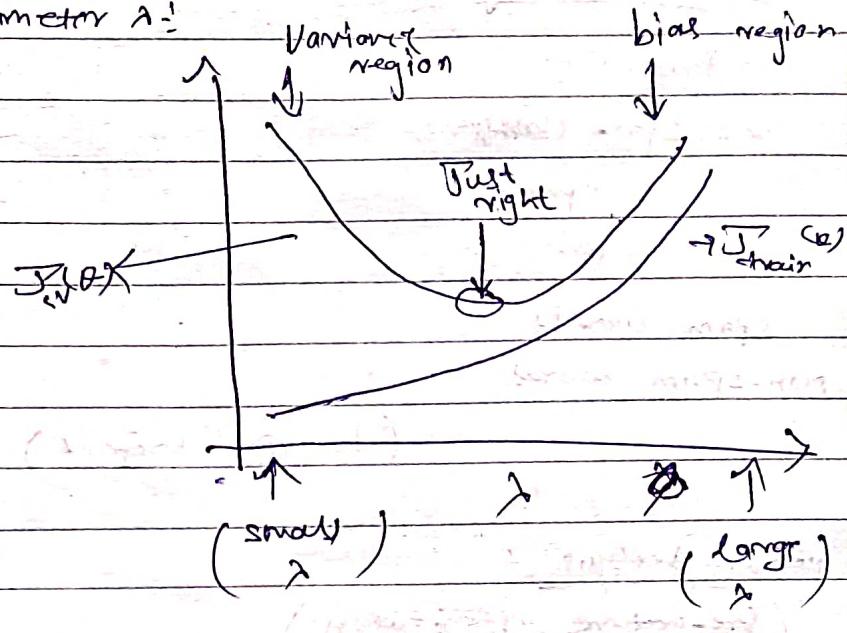
$$J_{train}(\theta) = \frac{1}{2m_{train}} \sum_{i=1}^{m_{train}} h_\theta(x^{(i)}_{train}) - y^{(i)}_{train})^2$$

$J_{cv}(\theta)$

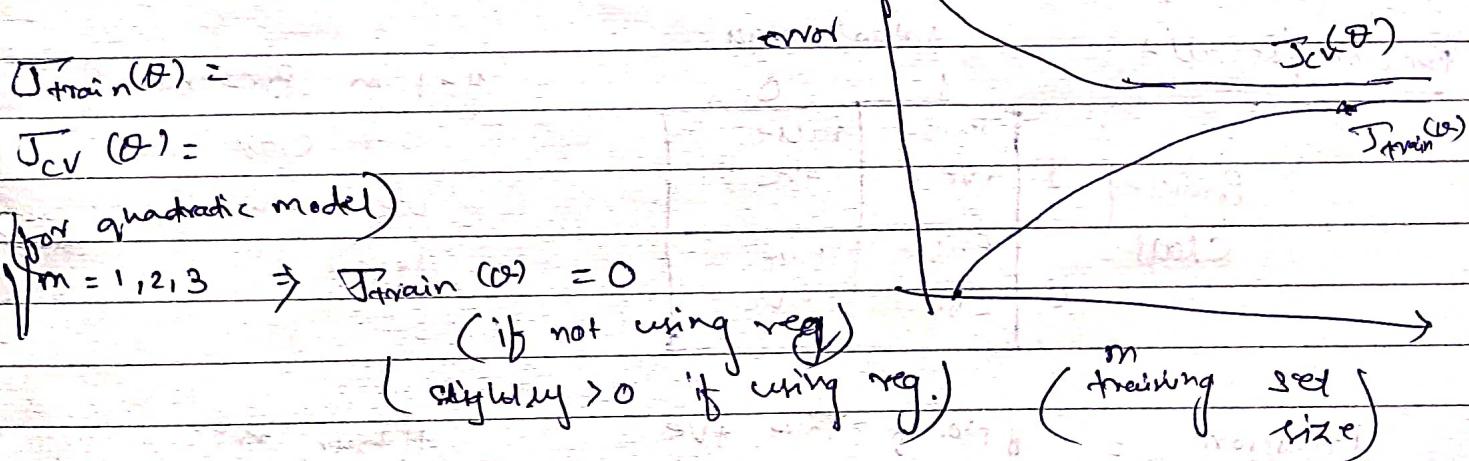
\hookrightarrow without regularization term

Try diff $\lambda = 0$	Put these	get	Put other	... in point
0.01	$\lambda = m$	value of θ_m in $J_{cv}(\theta)$	for which $J_{cv}(\theta)$ is min	for which $J_{cv}(\theta)$ is min
0.02	$J_{train}(\theta)$	for each		
1				
10				

Bias/Varr as a funcⁿ of reg. parameter λ :



Learning Curves



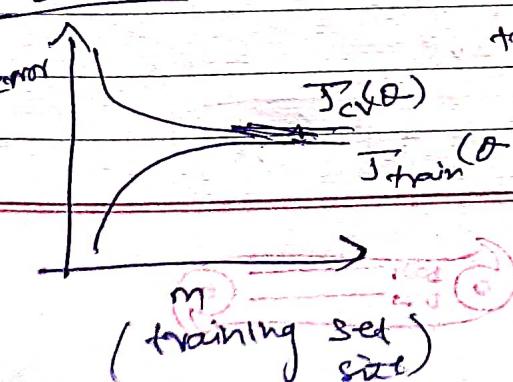
high bias:

getting more training data → high variance

would help

if algo is suffering

from
High bias



$J_{cv}(\theta)$

$J_{train}(\theta)$

gap ↑

getting more data may help in this case

small N.N \rightarrow fewer parameters \Rightarrow Underfitting
 Large N.N \rightarrow More " " \Rightarrow Overfitting.

Machine Learning System Design

Prioritizing

spam classification:

spam - 1

non-spam - 0

spam words:

non-spam words:

("honeyPot" Project)

Error Analysis

(pre-mature optimization)

Error Metrics for Skewed Class

Precision / Recall =

Actual class

		1	0
		True +ve	False +ve
Predicted Class	1	+ve	-ve
	0	-ve	+ve

$y=1$ in presence of

more class that we want to detect -

$$\text{Precision} = \frac{\# \text{No. of True +ve}}{\# \text{Predicted +ve}} = \frac{\# \text{True +ve}}{\# \text{True +ve} + \# \text{False +ve}}$$

high precision \Rightarrow good.

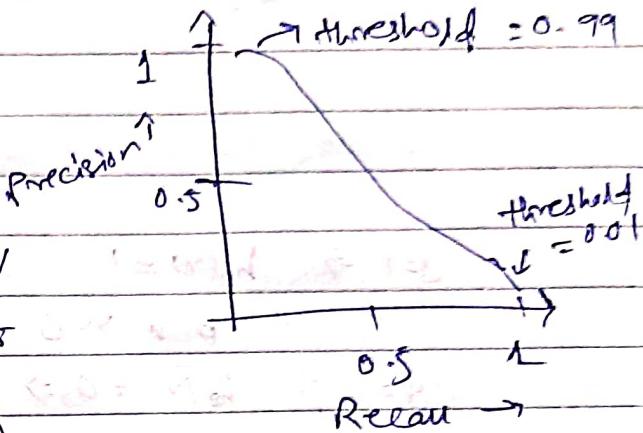
$$\text{Recall} = \frac{\# \text{True +ve}}{\# \text{actual +ve}} = \frac{\# \text{True +ve}}{\# \text{True +ve} + \# \text{False -ve}}$$

high recall \Rightarrow good.

Trading off Precision & Recall:

Higher precision, lower recall

generally, Predict 1 if $h_\theta(x) >$ threshold,



F₁ score (f-score)

	Precision (P)	Recall (R)	Avg
Alg 1	0.5	0.7	0.44
2	0.7	0.1	0.175
<u>3</u>	<u>0.0291</u>	2.0	<u>0.392</u>

$$f_1 \text{ score} = 2 \frac{P \cdot R}{P + R} \quad (\in [0, 1])$$

not a good way.

1 $\rightarrow 0.444$
2 $\rightarrow 0.175$
3 $\rightarrow 0.0392$

Data for ML:

$\rightarrow J_{\text{train}}(\theta)$ will be small

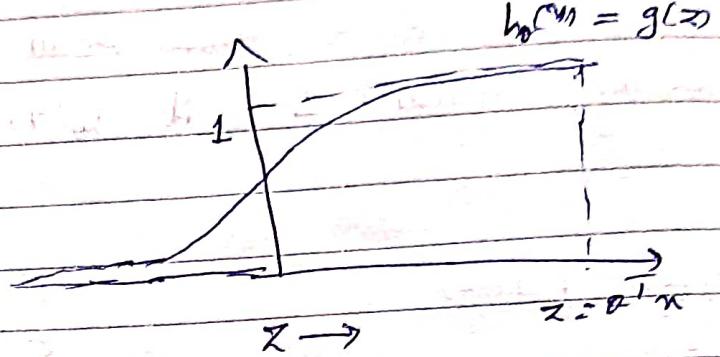
Using a v. large dataset $\rightarrow J_{\text{train}}(\theta) \approx J_{\text{test}}(\theta)$
(some cases)

$\rightarrow J_{\text{test}}(\theta)$ will also be small.

Support Vector Machine : (SVM)

optimization objective :

$$L.R : h_{\theta}(y) = \frac{1}{1 + e^{-\theta^T x}}$$



$$y=1 \Rightarrow h_{\theta}(y) = 1$$

$$\Rightarrow \theta^T x > 0$$

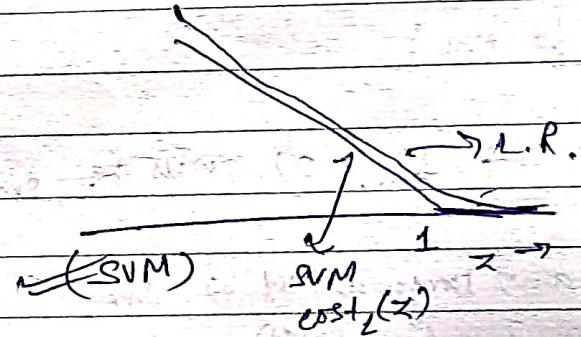
$$y=0 \Rightarrow h_{\theta}(y) = 0 \Rightarrow \theta^T x < 0$$

$$\text{cost func} = - (y \log h_{\theta}(y) + (1-y) \log h_{\theta}(y))$$

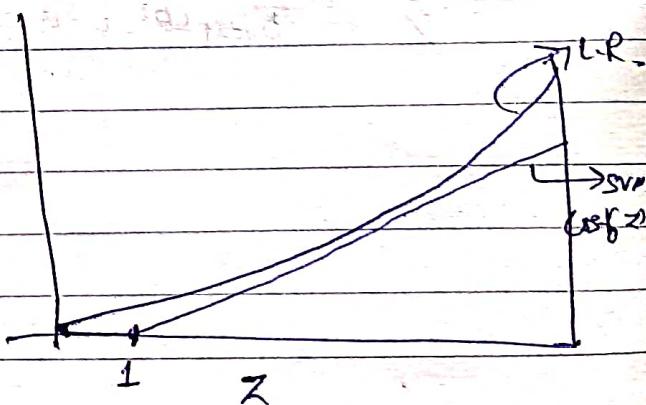
$$= -y \log \frac{1}{1+e^{-\theta^T x}} + -(1-y) \log \frac{1}{1+e^{-\theta^T x}}$$

$$\text{if } y=1 \quad (\text{and } \theta^T x > 0)$$

$$\Rightarrow \text{L.F.} = -y \log \frac{1}{1+e^{-\theta^T x}}$$



$$\text{If } y=0 \quad (\text{and } \theta^T x < 0)$$



$$\text{Cost func} = \underset{\theta}{\text{Min}} \frac{1}{m} \sum_{i=1}^m y^{(i)} \text{Cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{Cost}_0(\theta^T x^{(i)}) \\ (\text{SVM}) + \frac{\lambda}{2m} \sum_{j=0}^n \theta_j^2$$

2

$$L.R : A + \lambda B$$

$$\text{SVM: } C A + B ; \boxed{\epsilon = \frac{1}{2}}$$

then LR & SVM give same result

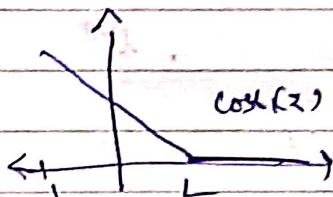
$$\text{Hypothesis : } h_0(n) = \begin{cases} 1 & , \text{ if } D_n^T X > 10 \\ 0 & , \text{ otherwise} \end{cases}$$

Large Margin Intuition - ?
Large Margin classifier)

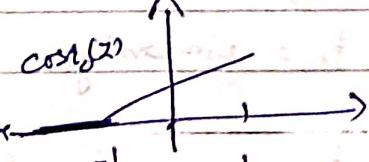
$$S V M = \underset{\text{cost}}{\underset{\theta}{\min}} \sum_{i=1}^m \left[y^{(i)} \text{Cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{Cost}_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

very large $C \Rightarrow$ to minimize, we have to put this 0.

when $y^{(i)} = 1$, $\Theta^T x^{(i)} \geq 1$



when $y^{(i)} = 0$, $\vec{w}^T \vec{x}^{(i)} \leq -1$



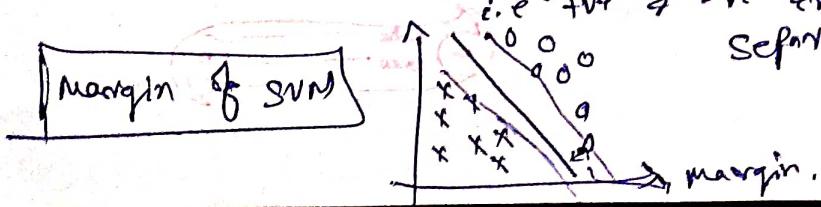
$$\therefore \min_{\theta} C(x_0) + \frac{1}{2} \sum_{j=1}^n \theta_j^2 = \min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

$$\text{st } \theta^T x^{(i)} > 1 \quad \text{if } y^{(i)} = 1$$

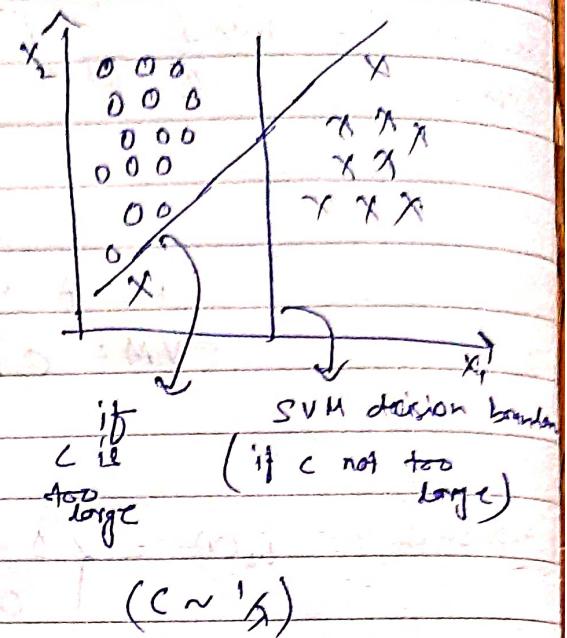
$$4 \cdot \theta^7 x^{(i)} \leq -1 \text{ if } g^{(i)} = 0$$

SVM Decision Boundary : (Linearly Separable Case)

We examples are
separated by a st. line.



Large Margin classifier in presence of outliers



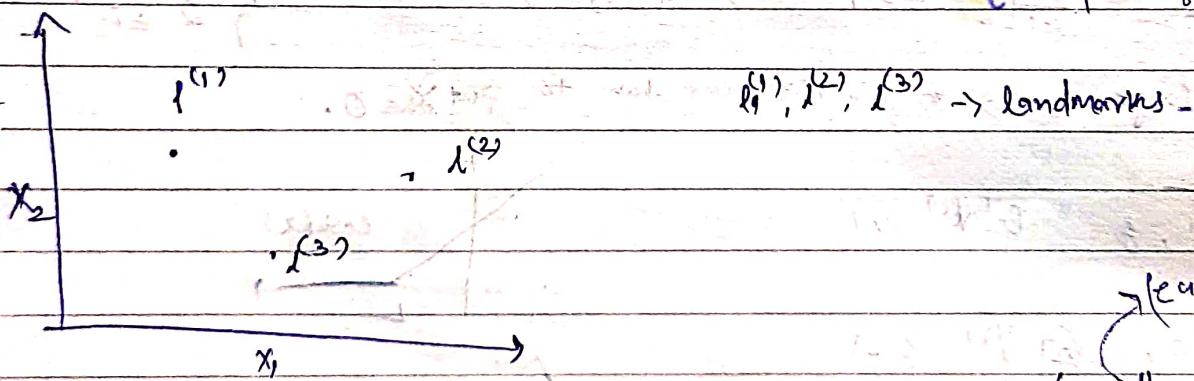
Math behind Large Margin classification

Kernels!

non-linear decision boundary

$$\rightarrow \theta_0 + \theta_1 f_1 + \theta_2 f_2 + \dots$$

$$h_\theta(x) = \begin{cases} 1 & \text{if } \theta_0 + \theta_1 f_1 + \dots > 0 \\ 0 & \text{otherwise} \end{cases}$$



$$\text{Given } x : f_1 = \text{similarity } (x, l^{(1)}) = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right)$$

$$f_2 = \text{similarity } (x, l^{(2)}) = \exp\left(-\frac{\|x - l^{(2)}\|^2}{2\sigma^2}\right)$$

$$f_3 = \text{similarity } (x, l^{(3)}) = \exp\left(-\frac{\|x - l^{(3)}\|^2}{2\sigma^2}\right)$$



similarity is kernel func.
 (Gaussian Kernel)
 is used here.

$$t_i = \exp\left(-\frac{\|x - x^{(i)}\|^2}{2\sigma^2}\right) = \exp\left(-\frac{\sum_{j=1}^n (x_j - x_j^{(i)})^2}{2\sigma^2}\right)$$

$$\text{If } x \approx x^{(i)} \Rightarrow t_i \approx \exp\left(-\frac{0^2}{2\sigma^2}\right) \approx 1$$

$$\text{If } x \text{ is far from } x^{(i)} \Rightarrow t_i = \exp\left(-\frac{\text{large no.}}{2\sigma^2}\right)$$

for practice

Landmarks will be at the same location of the training examples.

~~parametr~~ $\rightarrow C \approx \lambda^{-1}$ large $C \Rightarrow$ lower bias,
 (. small λ) high variance.

\rightarrow small $C \Rightarrow$ lower var
 (large λ) high bias

σ^2 : Large \rightarrow features fit very more smoothly
 (higher bias, low var)

small \rightarrow less smoothly -

(lower bias, higher var)

SVM with kernel's flags \rightarrow watch again

(refer to any internet stuff)

Using A SVM

Most used kernels in SVM +
Linear Kernel (no kernel),
Gaussian Kernel.

→ SVM Software Packages :- liblinear, libsvm etc.

→ choosing Parameter γ .

→ Choice of Parameter C

→ Choice of Kernel (sim func)

Gaussian Kernel \Rightarrow need to choose γ^2 .

Mercer's Theorem

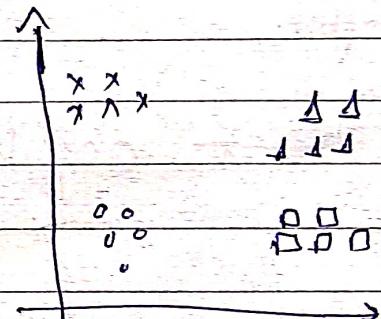
This is a tech condⁿ that should be satisfied by kernels.

• Multiclass Classification \Rightarrow

Use built in functionality

or

Use one-vs-all method.



Logistic Regression vs SVM:

$n = \text{no. of features}$, $m = \text{no. of training ex.}$

→ n is large \Rightarrow LR or SVM w/o kernel.
 $(n \gg m)$ (linear kernel)

→ m is small, n is intermediate \Rightarrow SVM with gaussian kernel.

→ n , m is large \Rightarrow create / add more features,
then LR or SVM w/o kernel.

unsupervised Learning

Clustering

(in Labeled data)

K-Means Algo :

'Input : $K = (\text{no. of clusters})$

training set $\{u^{(1)}, u^{(2)}, \dots, u^{(m)}\}$

$u^{(i)} \in \mathbb{R}^n$ (drop $u_0 = 1$ convention)

Randomly initialize K cluster centroids $u_1, u_2, \dots, u_K \in \mathbb{R}^n$

Repeat

cluster
assignment
step

for $i=1$ to m

$c^{(i)} = \text{index (from 1 to } K \text{) of cluster centroid}$
closest to $u^{(i)}$ $\min \|u^{(i)} - u_k\|^2$.

more
centroid
step

for $k=1$ to K

$u_k = \text{avg (mean) of points assigned to}$
cluster k .

K-means for non-separated clusters \rightarrow

Optimization objective \vdash

$c^{(i)}$ = index of cluster to which $u^{(i)}$ is assigned.

($1, 2, \dots, K$)

u_k = cluster centroid in ($u_k \in \mathbb{R}^n$)

$u_{c^{(i)}}$ = u_k " if the cluster to which $u^{(i)}$ had been
assigned.

obj. J :

$$J(c^{(1)}, \dots, c^{(m)}, u_1, \dots, u_K) = \frac{1}{m} \sum_{i=1}^m \|u^{(i)} - u_{c^{(i)}}\|^2$$

$$\min_{u^{(1)}, \dots, u^{(m)}} J$$

u_1, \dots, u_K (Distortion)

(Cost function)

\rightarrow cluster assignment step minimizes $J(\cdot)$
wrt $c^{(1)}, c^{(2)}, \dots, c^{(m)}$

(holding u_1, u_2, \dots, u_K fixed)

\rightarrow more centroid step minimizes $J(\cdot)$
wrt u_1, \dots, u_K .

Random Initialization

Randomly initialize cluster centroids +

(initialize k -means multiple times)

for $i=1$ to 100 {

Randomly initialize k -means.

Run k -means : get $c^1, \dots, c^m, \mu_1, \dots, \mu_k$

compute cost func (distortion) $J(\dots)$

Pick lowest $J(\dots)$.

→ (It works ^{good} for small no. of clusters!)

($n, k = 2$ to 10)

Choosing No. of Clusters

Elbow Method:

Dimensionality Reduction

Application :

Data compression :

Data visualization,

Principal Component Analysis

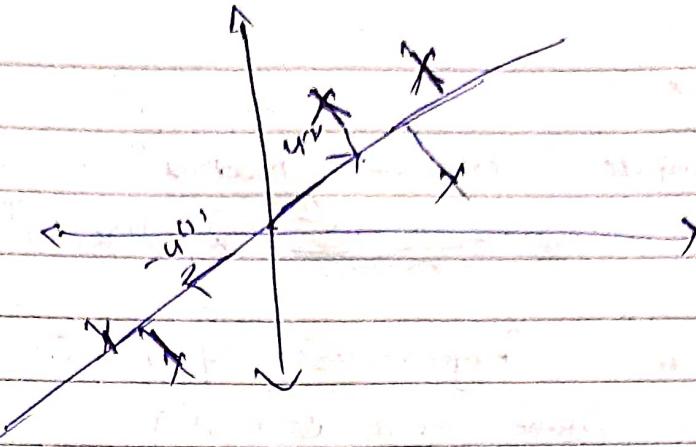
(PCA)
(algo for dimensionality reduction)

Projection error → Minimize.

✓ Reduce from 2-D to 1-D + By finding a dirⁿ/a vector)

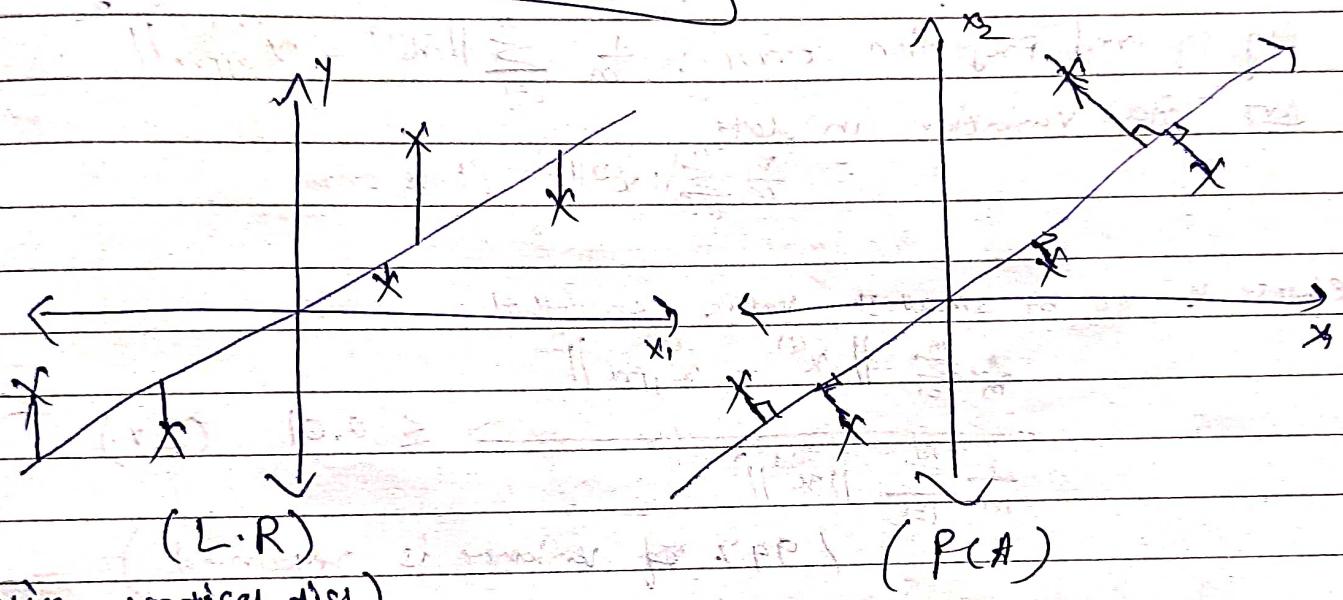
onto which to project the data so as to minimize
the projection error.





Reduce n-D to k-D \vdash find k-vectors u^1, u^2, \dots, u^k

PCA is not Linear Regression.



We want to predict the
special var/distinguish var y
using x_i in L.P.

(there is no sp var/
distinguish var in fcn)

x_1, x_2, \dots, x_n → are features
(treated equally)

Data Preprocessing (Before PCA):

(feature scaling / mean normalization)

Mean = 0

$$M_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$$

Replace $x_j^{(i)}$ with $x_j^{(i)} - M_j$ in data

If $x_j^{(i)}$ have diff scale?

$$x_j^{(i)} \leftarrow \frac{x_j^{(i)} - M_j}{S}$$

$S \rightarrow$ (max-min)
range.

Algo (PCA)

→ Complete "covariance matrix":

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}^{(i)}) (\mathbf{x}^{(i)})^T$$

covariance matrix

satisfies

"Symmetric Positive Definite."

→ "eigenvectors" of Matrix Σ

(take 1st K vectors)
for u^1, \dots, u^K

Principal
choosing no. of components K

Avg squared projection error: $\frac{1}{m} \sum_{i=1}^m \| \mathbf{x}^{(i)} - \mathbf{x}_{\text{approx}}^{(i)} \|_2^2$

Total variation in data

$$= \frac{1}{m} \sum_{i=1}^m \| \mathbf{x}^{(i)} \|_2^2 \quad (\text{wrt origin})$$

Choose K to be smallest value so that

$$\frac{\frac{1}{m} \sum_{i=1}^m \| \mathbf{x}^{(i)} - \mathbf{x}_{\text{approx}}^{(i)} \|_2^2}{\frac{1}{m} \sum_{i=1}^m \| \mathbf{x}^{(i)} \|_2^2} \leq 0.01 \quad (L7)$$

($\approx 99\%$ variance is retained)

Reconstruction from compressed Representation:

$$\mathbf{x} = \text{Ureduce } \mathbf{z}$$

$$\mathbf{x}_{\text{approx}}^{(i)} = \text{Ureduce } \mathbf{z}^{(i)}$$

$$\begin{pmatrix} \mathbf{I}_n \\ \mathbf{R} \end{pmatrix} \approx \mathbf{x}$$

Supervised Learning Pipeline

$(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$ $x^{(i)} \in \mathbb{R}^{10,000}$

labeled inputs

unlabeled dataset: $x^{(1)}, x^{(2)}, \dots, x^{(n)} \in \mathbb{R}^{10,000}$



PCA

$x^{(1)}, x^{(2)}, \dots, x^{(m)} \in \mathbb{R}^{10,000}$

New training set: $(z^{(1)}, y^{(1)}), (z^{(2)}, y^{(2)}), \dots, (z^{(m)}, y^{(m)})$

$$h_{\theta}(z) = \frac{1}{1 + e^{-\theta z}} \quad (\text{for L.R.})$$

~~Same $x \rightarrow z$ (for test set)~~

PCR

Mapping $x^{(i)} \rightarrow z^{(i)}$ should be defined by running PCA
only on the training set

find the mapping $x \rightarrow z$

then apply this mapping to x_{test} & x_{cv}

→ Do NOT Use PCA to Prevent Overfitting

→ II to Design of ML System

Anomaly Detection

Density Estimation

→ Dataset $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

→ Is x_{test} anomalous?

model - $P(x)$

$P(x_{\text{test}}) < \epsilon \rightarrow$ flag anomaly

$P(x_{\text{test}}) > \epsilon \rightarrow \text{ok}$

appln: fraud detector

(ii) Fraud Detection

→ $x^{(i)}$ = features of user i's activities.

→ Model $P(x)$ from data.

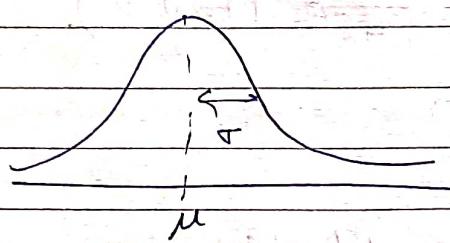
Identify unusual users by checking which have $P(x) \ll e^{-\infty}$

(iii) Manufacturing

→ Monitoring computers in a data center.

Gaussian Distribution

(Normal)



$$P(x; \mu, \sigma^2) = \frac{1}{(\sqrt{2\pi})\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

Parameter estimation

Data set: $\{x^{(1)}, \dots, x^{(m)}\}, x^{(i)} \in \mathbb{R}^n$

$$x^{(i)} \sim N(\mu, \Sigma)$$

↑ unknown

finding μ, Σ is aim of Parameter estimation problem

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)} \quad \Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)^2$$

algo

Training: $\{x^{(1)}, \dots, x^{(m)}\}$

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

$x_i \sim N(\mu_i, \Sigma_i)$

$$P(x) = P(x_1; \mu_1, \Sigma_1) \cdots P(x_n; \mu_n, \Sigma_n)$$

(events are closed to independent)

$$\prod_{j=1}^n P(x_j; \mu_j, \Sigma_j)$$

Anomaly Detection vs supervised Learning

A.D

→ small no. of +ve examples. ($y=1$)
Large " " " " . ($y=0$)
(Normal)

↓ → fraud detection.

→ Manufacturing

→ Monitoring machines in
doctor centre

etc.

S.L

→ Large +ve & -ve examples.

etc. → email spam classifier.

→ weather Predictor.

→ cancer classification

#

Try to make the data Gaussian.

If not gaussian \rightarrow $\log(n) \rightarrow$ gaussian.
or $\log(n+c)$
or $\log(\sqrt{n})$

Error Analysis:

(Miss some features to create a new feature & examine)

Multivariate Gaussian Dist.

\mathbb{R}^n

Don't model $P(n_1), P(n_2) \dots$ etc separately.

Model $P(n)$ all in one go.

Parameters: $\mu \in \mathbb{R}^n$
 $\Sigma \in \mathbb{R}^{n \times n}$ (co-variance matrix)

$P(x)$

$$P(x; \mu, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

$$\Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}$$

mean) x_1 has var = 2

x_2 has var = 2

Correlation (x_1, x_2) = 0

Recommender System & (app) & ML)

Adv. formulae!

Product recommendation by Amazon, Netflix etc.

Content based recommendation:

- (1) Predicting movie rating:
 $n_u = \text{No. of users} \quad \| \quad n_m = \text{no. of movies}$
- $r(i,j) = 1$ if user j has rated movie i (0 otherwise)
- $y^{(i,j)} = \text{rating by user } j \text{ on movie } i$ (if defined)
- $\theta^{(i)}$ = parameter vector for user j
- $x^{(i)}$ = feature vector for movie i
- for user j , movie i , predicted rating: $(\theta^{(j)})^T x^{(i)}$
- $m(j) = \text{no. of movies rated by user } j$.

$$\min_{\theta^{(1)}, \dots, \theta^{(m)}} \frac{1}{2} \sum_{j=1}^{n_u} \sum_{i: r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \lambda \sum_{j=1}^{n_u} (\theta^{(j)})^2$$

Optimization algo: (min J)

To learn $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(m)}$

$$\min_{\theta^{(1)}, \dots, \theta^{(m)}} \frac{1}{2} \sum_{j=1}^{n_u} \sum_{i: r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \lambda \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2$$

grad. descent update:

$$\theta_k^{(j)} = \theta_k^{(j)} - \alpha \sum_{i: r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} \quad (\text{for } k=0)$$

(regularization term)

$$\theta_k^{(j)} = \theta_k^{(j)} - \alpha \left(\sum_{i: r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} + \frac{\lambda}{2} \theta_k^{(j)} \right) \quad (\text{for } k \neq 0)$$

(regularization term)

for linear R:

We have $\frac{1}{m} \sum_{i=1}^m$ here, so we can ignore it.

Collaborative filtering: (with η)

(uses feature learning)

opt. algo.:

$$\text{given } \theta^{(1)}, \dots, \theta^{(n)}, \text{ to learn } \eta^{(i)}: \\ \min_{\eta^{(i)}} \frac{1}{2} \sum_{j: r(i,j)=1} ((\theta^{(j)})^T \eta^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{k=1}^n (\eta_k^{(i)})^2$$

'to learn' $\eta^{(i)}, \dots, \eta^{(m)}$:

$$\min_{\eta^{(1)}, \dots, \eta^{(m)}} \frac{1}{2} \sum_{i=1}^{n_m} \sum_{j: r(i,j)=1} ((\theta^{(j)})^T \eta^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (\eta_k^{(i)})^2$$

so $\theta \rightarrow \eta \rightarrow \theta \rightarrow \eta \dots$

Collaborative filtering Algo: (improved)

Putting this f. this together:

$$J(\theta^{(1)}, \dots, \theta^{(n)}, \eta^{(1)}, \dots, \eta^{(m)}) = \frac{1}{2} \sum_{(i,j): r(i,j)=1} ((\theta^{(j)})^T \eta^{(i)} - y^{(i,j)})^2$$

$$\begin{aligned} \min_{\theta^{(1)}, \dots, \theta^{(n)}, \eta^{(1)}, \dots, \eta^{(m)}} & \frac{1}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\eta_k^{(j)})^2 \\ & + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{i=1}^{n_m} ((\theta^{(j)})^T \eta^{(i)} - y^{(i,j)})^2 \end{aligned}$$

$\rightarrow \exists x \in \mathbb{R}^n$

\rightarrow Initialize $\eta^{(1)}, \dots, \eta^{(m)}, \theta^{(1)}, \dots, \theta^{(n)}$ to small random values

$\rightarrow \min J(\theta, \eta)$

\rightarrow for a user with parameter θ & a movie with feature x , then predict a star rating of $\theta^T x$

Vectorization! Low-Rank Matrix Factorization!

↙(LRMF)

Predicted ratings +

$$\left[\begin{array}{c} (\theta^{(1)})^T u^{(1)} \\ | \\ (\theta^{(4)})^T u^{(4)} \end{array} \right] - \dots - \left[\begin{array}{c} (\theta^{(n_u)})^T u^{(1)} \\ | \\ (\theta^{(n_u)})^T u^{(n_u)} \end{array} \right] = \underline{\underline{x}} \theta^T$$

$$x = \begin{bmatrix} - (u^{(1)})^T - \\ - (u^{(2)})^T - \\ | \\ - (u^{(n_u)})^T - \end{bmatrix} \quad \theta = \begin{bmatrix} - (\theta^{(1)})^T - \\ | \\ - (\theta^{(n_u)})^T - \end{bmatrix}$$

This has a Mathematical Property Related →
 (Low-Rank Matrix)
 (.Prop.)

(Rank deficient - If Matrix doesn't have the full rank)

How to find movie j related to movie i?

$$\Rightarrow \|u^{(i)} - x^{(j)}\| = \text{small} \quad (\text{feature vectors of movies})$$

then movies j & i are "similar"

Implementation Detail: Mean Normalization:

Take Mean,

Normalize i.e. $(x_i - \bar{x})$.

Net it to learn $\theta^{(1)}, \underline{\underline{x}}^{(1)}$

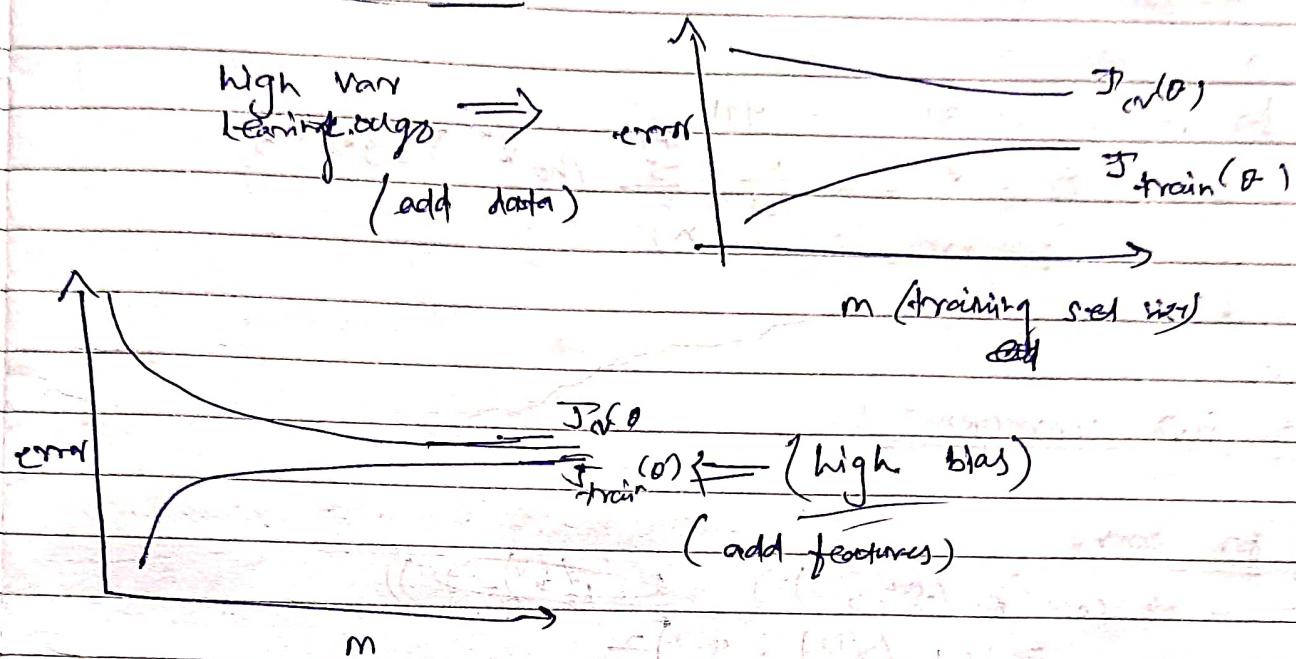


For user j on movie i predict:

$$\rightarrow (\theta^{(i)})^T (x^{(i)}) + b_i$$

Large Scale M.L.

Learning with large datasets:



Stochastic Gradient descent:

(batch G.D.)

(gradient descent = Batch G.D.)

$$\text{cost}(\theta, (x^{(i)}, y^{(i)})) = \frac{1}{2} (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$J_{\text{train}}(\theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}(\theta, (x^{(i)}, y^{(i)}))$$

Just 1 random
training $\{m\}$,
instead of
all.

→ Randomly shuffle dataset.

→ Repeat of for $i=1, \dots, m$

$$\theta_j = \theta_j - \alpha (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

$$(\text{for } j=0, \dots, n) \quad \hookrightarrow \frac{\partial}{\partial \theta_j} \text{cost}(\theta, (x^{(i)}, y^{(i)}))$$

4 3

→ $(x^{(i)}, y^{(i)}) \rightarrow$ progress → $(x^{(j)}, y^{(j)}) \rightarrow$ progress → - - -

Progress towards global minima

Batch G.D. : Use all m obs. in each interval.

Stochastic G.D. : Use 1 " " "

Mini Batch grad. descent :

VR

here - Use b obs. in each iter.

b = Mini-batch size.

Say + b = 10, m = 1000 :

Repeat {

for i = 1, 11, 21, 31, ... - 991

$$\theta_j = \theta_j - \alpha \frac{1}{10} \sum_{k=1}^{10} (h_\theta(x^{(k)}) - y^{(k)}) \alpha_j^{(k)}$$

(for every j = 0, ..., n)

}

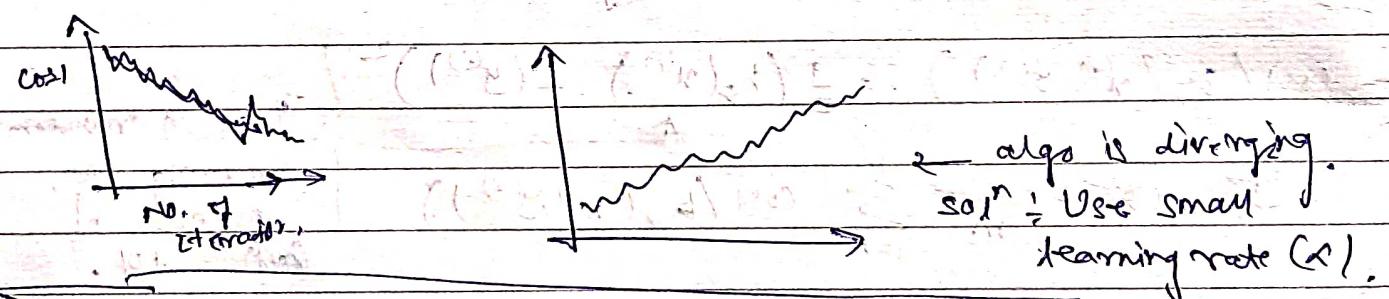
Stochastic G.D. convergence :

Checking for conv.

$$\rightarrow \text{cost}(\theta, (x^{(i)}, y^{(i)})) \quad \cancel{(h_\theta(x^{(i)}) - y^{(i)})}$$
$$= \frac{1}{2} (h_\theta(x^{(i)}) - y^{(i)})^2$$

→ During learn. compute $\text{cost}(\theta, (x^{(i)}, y^{(i)}))$ before updating θ . using $(x^{(i)}, y^{(i)})$

→ Every 1000 iter (say), plot this averaged over the last 1000 obs. processed by algo.



L.R. \Rightarrow cost. & over time if our θ starts to converge,

$$\text{Can start w/ decresc. } d = \frac{\text{cost 1}}{\text{iter no. 1} + \text{cost 2}}$$

cos 2 iter no. $\Rightarrow \alpha$



Training batch, learning rate α .

Online Learning: ~~also called Stochastic Gradient Descent~~

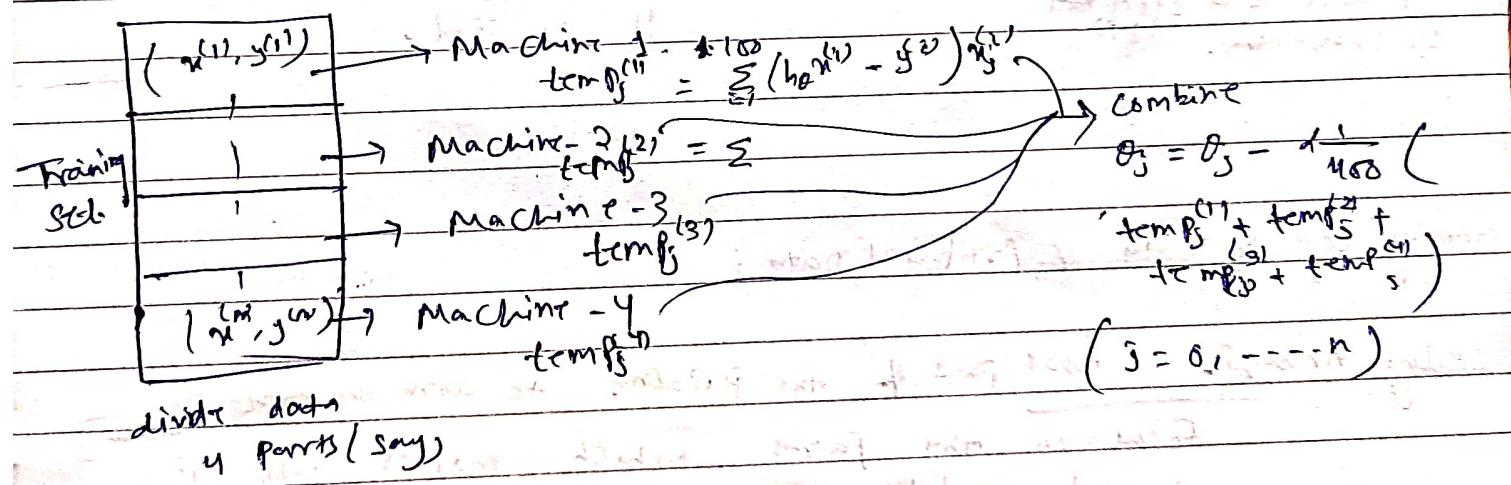
(Logistic regression or M.M.)

{ Repeat forever <
 get (x, y) corresponding to user.
 Update θ using (x, y)
 $\theta_j = \theta_j - \alpha (h_\theta(x) - y) x_j \quad (j=0, \dots, n)$

(It can adapt to changing user preference)

Pred (Predicted CTR Prob) $P(y=1 | x; \theta)$

Map Reduce & Data Parallelism



Application Example: Photo OCR

(Photo Optical Character Recognition)

Problem description & Pipeline:

Pipeline of Photo OCR

- Text Detection.
- Character Segmentation.
- Character Classification.

Sliding windows

(supervised learning classifier)

Step-size / stride

"expansion" operator

for char. segmentation!

(run supervised learning model?)

-ve ex + (y=1)

If there is a split between
(gap)

2. character.

-ve ex / (y=0)

If no split ->

! Getting lots of Data & Artificial Data!

Ceiling Analysis: what part of the pipeline to work on more

Focus on the parts which makes diff in overall model.