

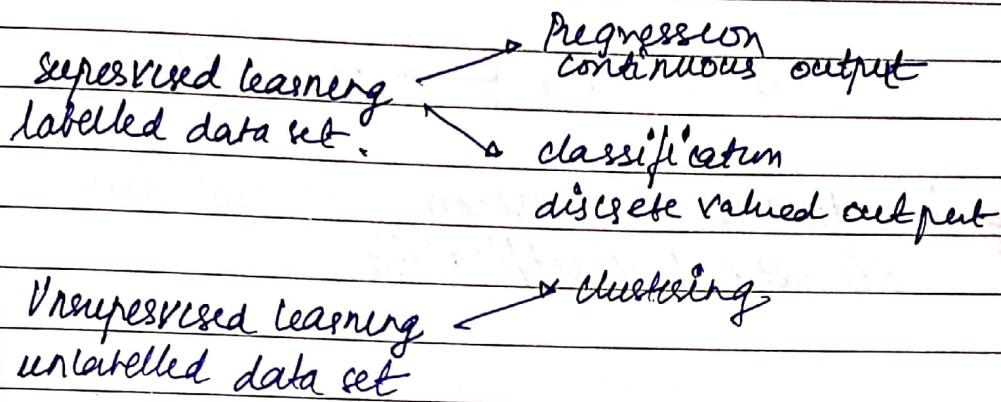
Machine learning coursera Andrew Ng

classmate

Date _____
Page _____

Tom Mitchell 1998

A well posed learning problem: A computer program is said to learn from experience E with respect to some task T and some performance measure P , if its performance ~~measures~~ P in T as measured by P improves with experience E .



Model representation. (Supervised learning)

$x^{(i)}$ input variables

$y^{(i)}$ output or target variable that we are trying to predict

$(x^{(i)}, y^{(i)})$ is called a training example.

$(x^{(i)}, y^{(i)})$ $i = 1, \dots, m$ is called a training set

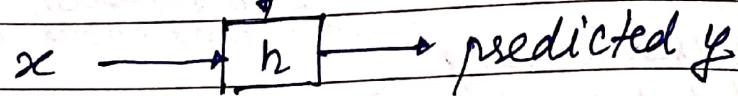
X space of input ~~variable~~ values

Y space of output ~~variable~~ values. called hypothesis

To learn a function $h: X \rightarrow Y$ so that $h(x)$ is a 'good' predictor for the corresponding value of y .

Training
set

Learning
Algorithm



If y is continuous : regression
discrete : classification

~~Classification~~
example

Univariate Linear regression

$$h(x) = \theta_0 + \theta_1 x$$

cost function

θ_i s parameters

how to choose θ_i ?

minimise error $\sum_{i=1}^n (h_\theta(x) - y)^2$

minimise $\frac{1}{2m} \sum_{i=1}^n (h_\theta(x) - y)^2$

lost function $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$

gradient descent

the gradient descent algorithm.

repeat until convergence

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

simultaneously update the parameters.

correct

$$\text{temp } 0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\text{temp } 1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\theta_0 := \text{temp } 0$$

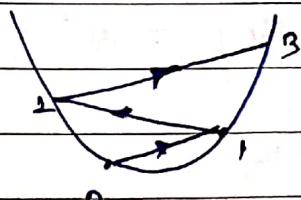
$$\theta_1 := \text{temp } 1$$

mathematically

$$\vec{\theta} := \vec{\theta} - \alpha \vec{\nabla} J$$

if α is too small gradient descent can be slow

if α is too large gradient descent can overshoot the min. It may fail to converge or even diverge



in case of univariate linear regression

$$J_{\theta} = \frac{1}{2m} \sum_{i=1}^m (\hat{y}_i - y_i)^2$$

$$\frac{\partial J}{\partial \theta_0} = \frac{1}{2m} \sum_{i=1}^m (\hat{y}_i + \theta_0 + \theta_1 x_i - y_i)^2$$

$$\frac{\partial J}{\partial \theta_1} = \frac{1}{2m} \sum_{i=1}^m (\hat{y}_i + \theta_0 + \theta_1 x_i - y_i)$$

$$\frac{\partial J}{\partial \theta_0} = \frac{1}{2m} \sum_{i=1}^m (\hat{y}_i + \theta_0 + \theta_1 x_i - y_i) x_i$$

Batch gradient descent.

this method looks at every example in the entire training set on every step

Multivariate linear regression

notation $x_j^{(i)}$ value of feature j in the i^{th} training eq
 $x^{(i)}$ input of the i^{th} training eq.

m the number of training eq.

n the number of features.

Recall

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_n x_n$$

$$h_{\theta}(x) = [\theta_0 \ \theta_1 \ \theta_2 \ \dots \ \theta_n] \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} = \theta^T x$$

$$x_0^{(i)} = 1$$

hypothesis $\theta^T x$

parameters θ n-dimensional vector

cost function $J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$

New algorithm

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

simultaneously update.

Feature scaling

To speed up gradient descent

have input values roughly in same range

$$x_i' := \frac{x_i - \mu_i}{s_i}$$

s_i : range of values or standard dev

Debugging gradient descent

- ① Make a plot $J(\theta)$ vs number of iterations
if $J(\theta)$ increases then you probably need to decrease α .
- ② Automatic convergence tests Declare convergence if $J(\theta)$ decreases by less than ϵ in one iteration, where ϵ is some small value such as 10^{-3} . In practice it is difficult to choose this threshold value.

To summarize

if α is too small: slow convergence

if α is too large: may not decrease on every iteration
and thus may not converge.

Features and Polynomial Regression

$$h_\theta(x) = \theta_0 + \theta_1 x_1$$

quadratic $h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2$

cubic $h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^3$

Created new features $x_2 = x_1^2$
 $x_3 = x_1^3$

scaling becomes impr

Normal equation

design matrix $X = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(n)} \end{bmatrix}$

solving $\bar{\nabla} J = 0$

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

gradient descent

Need to choose α

Needs many iterations

$O(n^2)$

Works well when n large

vs

Normal eqⁿ

No need to choose α

No need to iterate

$O(n^3)$

slow if n is very large

Normal equation Non invertibility

$X^T X$ is non invertible common causes

- ① Redundant features linearly dependent.
 - ② Too many features ($m \leq n$) delete some features or use regularization.

Logistic regression (classification algorithm)

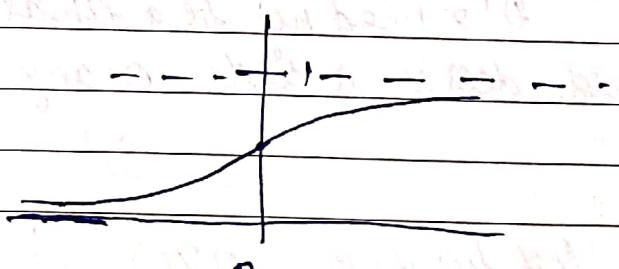
linear regression draw backs

logistic regression model $0 \leq h_0(x) \leq 1$

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

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

sigmoid func
logistic func



Interpretation of hypothesis output

$h_0(m)$ = estimated probability that $y=1$ on input x .

$$\partial_k \partial_{k\ell} = f$$

$$h_{\theta}(x) = P(Y=1 \mid x; \theta) \quad \text{prob that } y=1 \text{ given } x \\ \text{parameterized by } \theta.$$

logistic regression

$$h_{\theta}(x) = g(\theta^T x) = P(y=1/x; \theta)$$

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

predict $y=1$ $h_{\theta}(x) \geq 0.5$

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

$g(z) \geq 0.5$

when $z \geq 0$

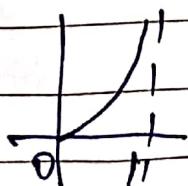
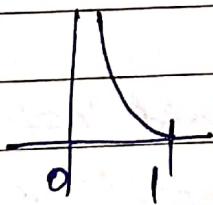
i.e. $\theta^T x \geq 0$

↑
decision boundary

$\theta^T x$ need not be a linear function (polynomial regression could describe a circle or any other shape to fit our data.)

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

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



guarantees convex

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

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h_\theta(x^{(i)})) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)}))$$

Vectorised implementation
 $h = g(X\theta)$

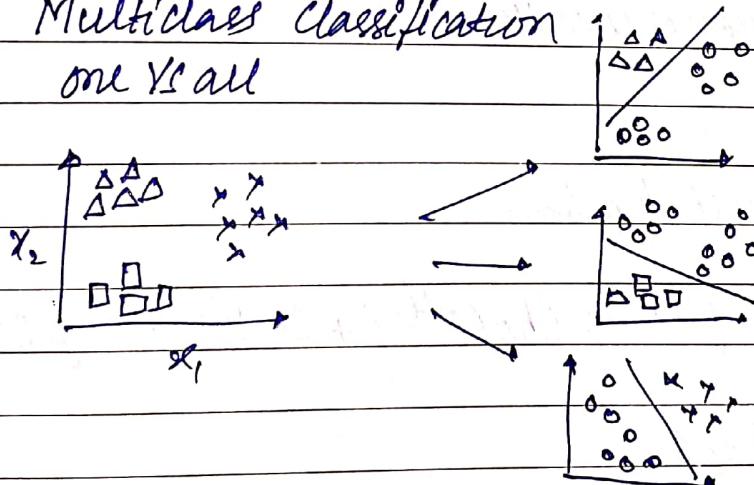
$$J(\theta) = -\frac{1}{m} (y^T \log(h) - (1-y)^T \log(1-h))$$

$$\frac{\partial J}{\partial \theta_j} = \frac{1}{m} \sum (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

$$\nabla = \frac{1}{m} X^T (g(X\theta) - y)$$

Multiclass classification

one vs all



To summarize

Train a logistic regression classifier for each class to predict the probability $y=i$

To make a new prediction pick the class that maximizes $h_\theta(x)$

The problem of overfitting

regularized cost function

linear regression $\min_{\theta} \frac{1}{2m} \left[\sum_{i=1}^m (\hat{y}_i - y_i)^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$

gradient

$$\frac{\partial J}{\partial \theta_j} = \frac{1}{m} \left(\sum_{i=1}^m (\hat{y}_i - y_i) x_i^{(i)} + \lambda \theta_j \right)$$

for $j \neq 0$

Normal equation

$$\theta = (X^T X + \lambda I)^{-1} X^T y$$

$$I = \begin{bmatrix} 1 & & & \\ & 1 & & \\ & & 1 & \\ & & & 1 \end{bmatrix}$$

if $m < n$, $X^T X$ is non-invertible
but $X^T X + \lambda I$ becomes invertible

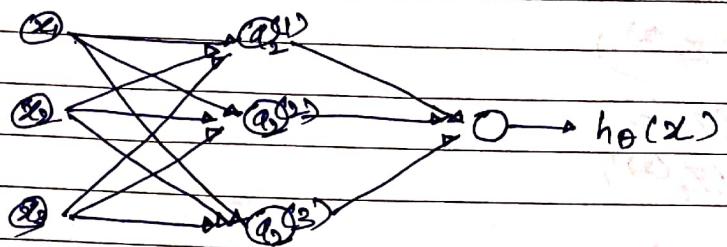
logistic regression regularized cost function

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(\hat{y}_i) + (1-y^{(i)}) \log(1-\hat{y}_i)$$

$$+ \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

* we don't regularize θ_0

Neural networks



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

$\Theta^{(j)}$ matrix of weights controlling function mapping from layer j to layer $j+1$

If we had one hidden layer

$$\begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \rightarrow \begin{bmatrix} a_1^{(2)} \\ a_2^{(2)} \\ a_3^{(2)} \end{bmatrix} \rightarrow h_\theta(x)$$

$$a_1^{(2)} = g(\Theta_{10}^{(1)} x_0 + \Theta_{11}^{(1)} x_1 + \Theta_{12}^{(1)} x_2 + \Theta_{13}^{(1)} x_3)$$

$$a_2^{(2)} = g(\Theta_{20}^{(1)} x_0 + \Theta_{21}^{(1)} x_1 + \Theta_{22}^{(1)} x_2 + \Theta_{23}^{(1)} x_3)$$

$$a_3^{(2)} = g(\Theta_{30}^{(1)} x_0 + \Theta_{31}^{(1)} x_1 + \Theta_{32}^{(1)} x_2 + \Theta_{33}^{(1)} x_3)$$

$$h_\theta(x) = a_1^{(3)} = g(\Theta_{10}^{(2)} a_0^{(2)} + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} a_2^{(2)} + \Theta_{13}^{(2)} a_3^{(2)})$$

If network has s_j units in layer j and s_{j+1} units in layer $j+1$
 then $\Theta^{(j)}$ will be of dimension $s_{j+1} \times (s_j + 1)$

↑
bias.

Vectorization.

$$q_1^{(2)} = g(Z_1^{(2)})$$

$$q_2^{(2)} = g(Z_2^{(2)})$$

$$q_3^{(2)} = g(Z_3^{(2)})$$

$$Z_K^{(2)} = \theta_{K,0}^{(1)} x_0 + \theta_{K,1}^{(1)} x_1 + \dots + \theta_{K,n}^{(1)} x_n$$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \quad Z^{(1)} = \begin{bmatrix} Z_1^{(1)} \\ Z_2^{(1)} \\ \vdots \\ Z_n^{(1)} \end{bmatrix}$$

$$Z^{(1)} = \theta^{j-1} x^{j-1}$$

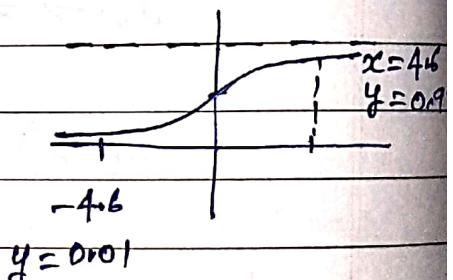
$$a^{(1)} = g(Z^{(1)})$$

last step

$$h_\theta(x) = a^{(n)} = g(z^{(n)})$$

Examples and intuitions

$$\begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} \rightarrow [g(Z^{(2)})] \rightarrow h_\theta(x)$$



x_0 bias variable

$$\theta^{(1)} = [-30 \ 20 \ 20]$$

$$h_\theta(x) = g(-30 + 20x_1 + 20x_2)$$

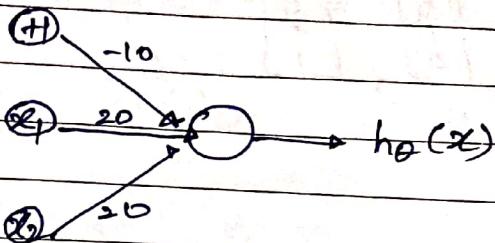
$$x_1 = 0 \text{ and } x_2 = 0 \quad g(-30) \approx 0$$

$$x_1 = 0 \text{ and } x_2 = 10 \quad g(-10) \approx 0$$

$$x_1 = 1 \text{ and } x_2 = 0 \quad g(10) \approx 0$$

$$x_1 = 1 \text{ and } x_2 = 1 \quad g(10) \approx 1$$

OR function



x_1	x_2	$h_\theta(x)$
0	0	$g(-10) = 0$
0	1	$g(10) = 1$
1	0	$g(10) = 1$
1	1	$g(30) = 1$

AND

$$\Theta^{(1)} = [-30 \quad 20 \quad 20]$$

NOR

$$\Theta^{(1)} = [10 \quad -20 \quad -20]$$

OR

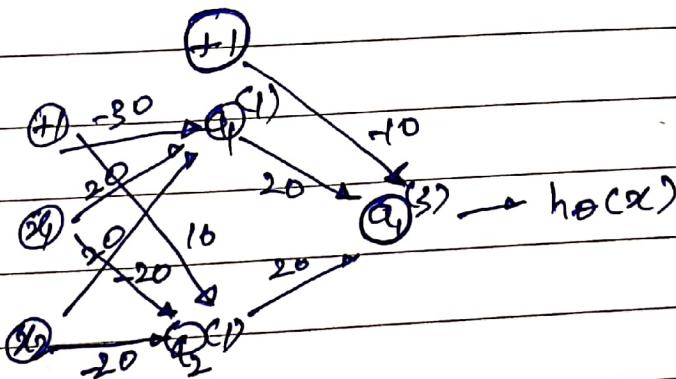
$$\Theta^{(1)} = [-10 \quad 20 \quad 20]$$

XNOR

$$\begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} \rightarrow \begin{bmatrix} af^{(2)} \\ qf^{(2)} \end{bmatrix} \rightarrow \begin{bmatrix} q^{(2)} \end{bmatrix} \rightarrow h_\theta(x)$$

$$\Theta^{(1)} = \begin{bmatrix} -30 & 20 & 20 \\ 10 & -20 & -20 \end{bmatrix}$$

$$\Theta^{(2)} = [-10 \quad 20 \quad 20]$$



multiclass classification using neural networks.

If 4 objects $y^{(i)} = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$

↓
set of resulting classes

cost function neural networks

L = total number of layers in the network.

s_l = number of units (not counting bias unit) in layer l .

K = number of output units/classes.

$h_{\theta}(x)_k$ hypothesis that results in K^{th} output

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_{\theta}(x^{(i)})_k) + (1-y_k^{(i)}) \log(1-h_{\theta}(x^{(i)})_k)$$

$$+ \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\theta_{j,i}^{(l)})^2$$

↓
excluding
the bias unit

Backpropagation

Training set of $(x^{(i)}, y^{(i)}) \quad i=1 \dots m$

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

For $i=1$ to 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 $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$

$$\delta^{(l)} = (\theta^{(l)})^T \delta^{(l+1)} + g'(z^{(l)}) \rightarrow a^{(l)} \cdot \star (1 - a^{(l)})$$

$$\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$$

Vectorized $\Delta^{(l)} := \Delta^{(l)} + \delta^{(l+1)} (a^{(l)})^T$

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

$$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} \quad j = 0$$

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

Implementation

- ① Unrolling parameters
- ② gradient checking
- ③ random initialisation

Evaluating a hypothesis

methods of troubleshooting

- ① getting more training examples.
- ② trying smaller sets of features.
- ③ trying additional features.
- ④ trying polynomial features.
- ⑤ increasing or decreasing λ .

dataset
/ \
training set test set

- ① learn $\hat{\theta}$ and minimize $J_{\text{train}}(\hat{\theta})$ using the training set.
- ② Compute the test set error $J_{\text{test}}(\hat{\theta})$

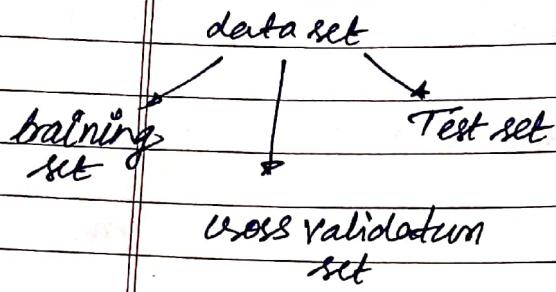
misclassification error

$$\text{err}(h_\theta(x), y) = 1 \quad \text{if } h_\theta(x) \geq 0.5 \text{ and } y=0 \\ \text{or } h_\theta(x) < 0.5 \text{ and } y=1$$

0 otherwise

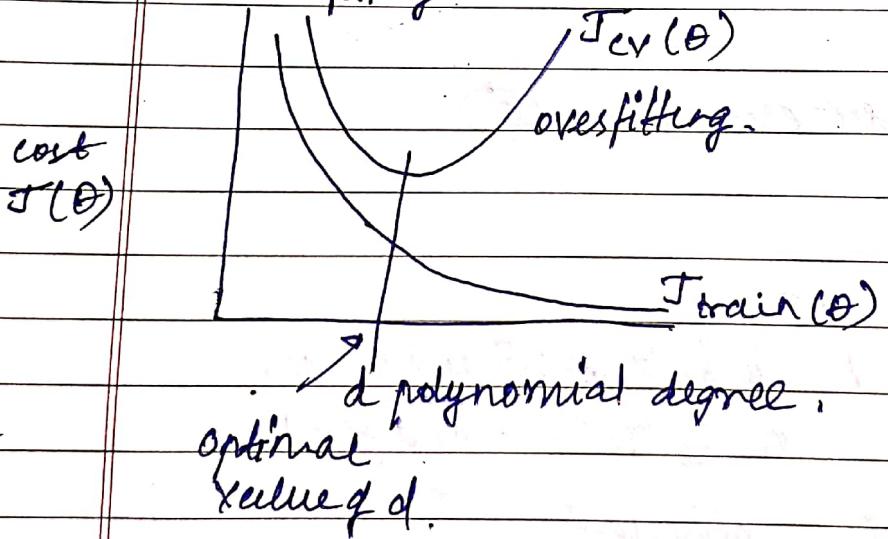
$$\text{Test error} = \frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} \text{err}(h_\theta(x_{\text{test}}^{(i)}), y_{\text{test}}^{(i)})$$

Model selection and train/validation/test sets.



- ① Optimise the parameters in Θ using the training set for each polynomial degree.
- ② Find the polynomial degree with least error using the cross validation set.
- ③ Estimate the generalization error using the test set with $J_{\text{test}}(\Theta^{(d)})$ $d = \theta$ from polynomial with lowest J .

Diagnosing Bias vs Variance
underfitting

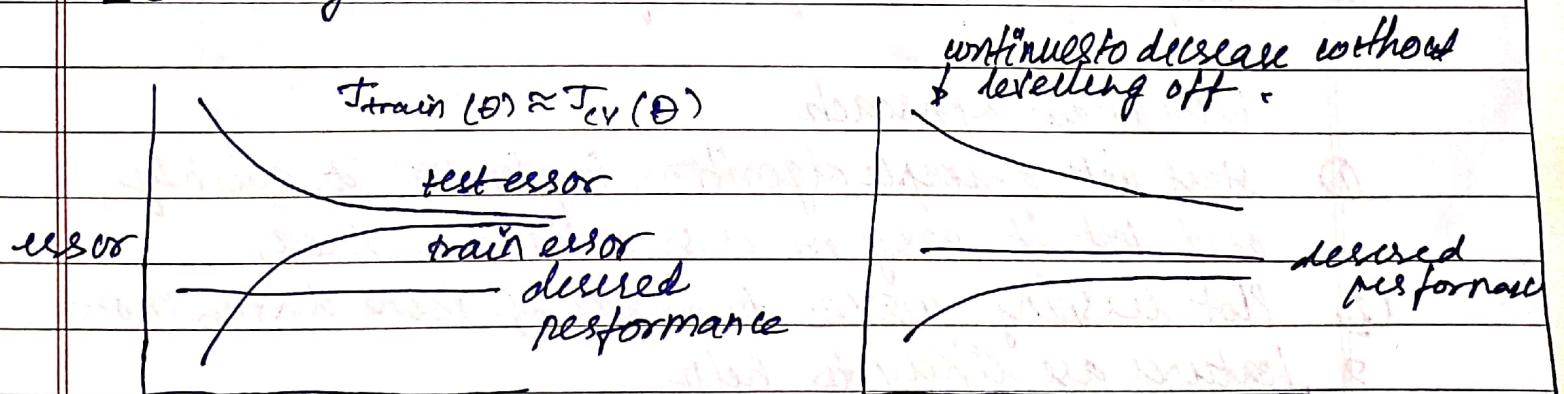


- (i) High bias (underfitting) both $J_{\text{train}}(\theta)$ and $J_{\text{cv}}(\theta)$ will be high AND $J_{\text{cv}}(\theta) \approx J_{\text{train}}(\theta)$
- (ii') High variance (overfitting) $J_{\text{train}}(\theta)$ will be low and $J_{\text{cv}}(\theta)$ will be much greater than $J_{\text{train}}(\theta)$.

Regularization and Bias/Variance

- ① Create a list of λ s $\lambda = 0 \ 0.01 \ 0.02 \ 0.03 \ \dots \ 10^{-2} 4$
- ② Create a set of models with different degrees or any other variants.
- ③ Iterate through λ s and for each λ go through all the models to learn some θ .
- ④ Compute the cross validation error using the learned θ (computed with λ) on the $J_{cv}(\theta)$ without regular
- ⑤ Select the best combo that produces the lowest error on the cross validation set.
- ⑥ Using the best combo θ and λ , apply it on $J_{test}(\theta)$ to see if it has a good generalization of the problem

Learning curves



high N training set size.

high bias.

more training data will
by itself not help much.

high variance

getting more data is
likely to help.

- ① getting more training examples (fixes high variance)
- ② trying smaller set of features (fixes high variance)
- ③ Adding features (Fixes high bias)
- ④ adding polynomial features (fixes high bias)
- ⑤ $\downarrow \lambda$ (fixes high bias)
- ⑥ $\uparrow \lambda$ (fixes high variance)

Machine Learning system design.

Prioritizing on what to work on.

- ① collect lots of data.
- ② Develop sophisticated features.
- ③ Develop algorithms to process your input in different ways.

Recommended approach

- ① Start with a simple algorithm, implement it quickly, and test it easily on cross validation data.
- ② Plot learning curves to decide if more data, more features are likely to help.
- ③ Manually examine the errors on the examples in the cross validation set and try to spot a trend where most of the errors were made.

We should try new things, get a numerical value of our error rate, and based on our result decide whether to keep the new feature or not.

Handling skewed ~~dataset~~ data

Precision / Recall

		Actual class	
		1	0
Predicted class	1	True positive	False positive
	0	False negative	True negative

Precision = $\frac{\text{True positives}}{\#\text{predicted positive}}$

Recall = $\frac{\text{True positives}}{\#\text{Actual positives}}$

Trading off precision and recall.

threshold = 0.7

1 if $h_0(x) \geq 0.7$
0 if $h_0(x) < 0.7$

higher precision lower recall.

threshold = 0.3

lower precision higher recall.

$$F_{\text{Score}} = \frac{2PR}{P+R}$$

Useful test: given the input x , can a human expert confidently predict y ?

Large data rationale

Use a learning algorithm with many parameters
low bias algo

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

Use a very large training set (unlikely to overfit)

$J_{\text{train}}(\theta) \approx J_{\text{test}}(\theta)$

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

Unsupervised learning

K-means algorithm

Input

→ K number of clusters

→ Training set $\{x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(m)}\}$

$x^{(i)} \in \mathbb{R}^n$ drop ($x_0 = 1$) convention

Randomly initialize K cluster centroids $\mu_1, \mu_2, \dots, \mu_K \in \mathbb{R}^n$

Repeat {

clusters assigned

for $i = 1$ to m

holding

μ_{const}

min w.r.t

c

$c^{(i)} :=$ index (from 1 to K) of cluster centroid
closest to $x^{(i)}$ $\min_k \|x^{(i)} - \mu_k\|$

move centroid

holding

μ_{const}

min w.r.t

μ

for $k = 1$ to K

$\mu_k :=$ average (mean) of points assigned to k^{th} cluster

optimization objective

$$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c(i)}\|^2$$

distortion function

min

$c^{(1)}, \dots, c^{(m)}$

μ_1, \dots, μ_K

$$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$$

Random initialization

Should have $K \leq m$

Randomly pick K training examples

Set μ_1, \dots, μ_K equal to these K examples

To avoid local optima

for $i = 1$ to 100 {

randomly initialize K-means

Run K-means set $c^{(1)}, \dots, c^{(m)}$, μ_1, \dots, μ_K

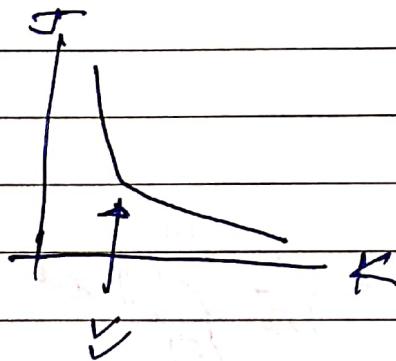
Compute cost function (distortion)

}

Pick clustering that gave lowest $J(C, \mu)$

Choosing number of clusters

① elbow method.



② evaluate K-means based on a metric for how well it performs for that later purpose.

Dimensionality Reduction & PCA.

Principal Component Analysis (PCA)

Data pre-processing

Training set $x^{(1)}, x^{(2)}, \dots, x^{(m)}$

$$x_j^{(i)} := \frac{x_j^{(i)} - \mu_j}{s_j} \quad \begin{array}{l} \text{mean normalization} \\ \text{feature scaling} \end{array}$$

PCA algorithm

Compute covariance matrix

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)}) (x^{(i)})^T \quad \cdot \frac{1}{m} * X' * X$$

Compute eigenvectors of Σ

$$[U, S, V] = \text{SVD}(\Sigma \text{ (sigma)});$$

singular value decomposition.

$$U_{\text{reduce}} = U(:, 1:k_r);$$

$$Z = U_{\text{reduce}}' * x_j;$$

reconstruction

$$x_{\text{approx}}^{(i)} = U_{\text{reduce}} * Z^{(i)}$$

choosing K number of principal components

Algorithm

Try PCA with $K=1$

Compute $Z^{(1)}, Z^{(2)}, \dots, Z^{(m)}$ $x_{\text{app}}^{(1)}, x_{\text{app}}^{(2)}, \dots, x_{\text{app}}^{(m)}$

$$\text{check if } \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{app}}^{(i)}\|^2 = 1 - \frac{\sum_{i=1}^K s_{ii}}{\sum_{i=1}^m s_{ii}} \leq 0.01$$

$$\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2$$

99% variance retained.

Supervised learning speedup

$(x^{(1)}, y^{(1)}) \quad (x^{(2)}, y^{(2)}) \dots$
PCA on $x^{(i)}$ to $z^{(i)}$

$(z^{(1)}, y^{(1)}) \quad (z^{(2)}, y^{(2)}) \dots$

Note mapping $x^{(i)} \rightarrow z^{(i)}$ should be defined by running PCA only on the training set -
this mapping can be applied as well to $x_{\text{cv}}^{(i)}$ $x_{\text{test}}^{(i)}$.

try running algo w/o PCA.

Anomaly detection

Dataset $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$
is x_{test} anomalous?

model $p(x)$

$p(x_{\text{test}}) < \epsilon$ flag anomaly
 $p(x_{\text{test}}) \geq \epsilon$ OK

Gaussian distribution

$$x \sim N(\mu, \sigma^2)$$

σ standard deviation

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

If each feature $\in \mathbb{R}^n$

$$p(x) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2)$$

algorithm

- choose features that you think might be indicative of anomalous examples.

- fit parameters $\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2$

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

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$\sigma_j^2 = \frac{1}{m-1} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$$

given new example compute $p(x)$

$$p(x) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2)$$

Anomaly if $p(x) < \epsilon$

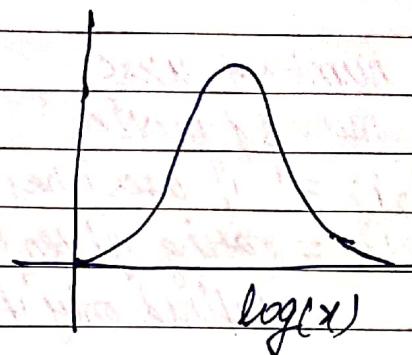
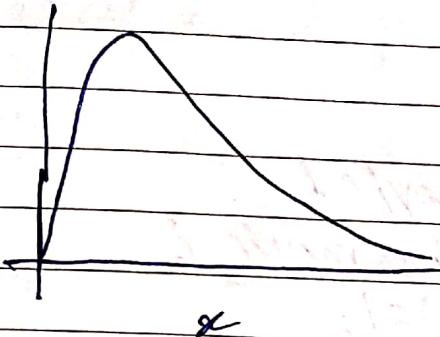
Algorithm evaluation

- 1 Fit model ~~on~~ \hat{y} on training set $\{x^{(1)}, \dots, x^{(m)}\}$
- 2 On a cross validation/test example x predict y
use precision and recall and F₁ score to evaluate
Can also use cross validation test to choose E .

Anomaly detection vs Supervised learning

① Very small number of +ve examples ($y=1$) ($0-20$ is common)	Large number of +ve and -ve examples,
② Large number of -ve examples	
③ Many different types of anomalies. Hard for any algorithm to learn from the examples what the anomaly looks like. future anomalies look very nothing like the ones we've seen so far.	Enough +ve examples for algorithm to get a sense of what the eg are like. future +ve egs are likely to be similar to ones in training set.
eg Fraud detection, manufacturing, monitoring machines in data center	email spam classification, weather prediction, genres classification.

choosing what features to use.



$\log(x+c)$

x/c

choose features that might take on unusually large or small values in the event of an anomaly

Multivariate gaussian distribution

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

$$\Sigma = \frac{1}{m} \sum (x^{(i)} - \mu)(x^{(i)} - \mu)^T$$

original model

manually create features to capture anomalies where x_1, x_2 take unusual combinations of values.

computationally cheaper scales better to large n

multivariate gaussian

automatically captures correlations b/w features

computationally more expensive $m > n$

Recommender systems

n_u number of users

n_m number of movies

$r(i,j) = 1$ if user j has rated movie i .

$y^{(i,j)}$ = rating given by user j to movie i .
defined only if $r(i,j) = 1$.

Content based recommenders systems

For each user j learn a parameter $\theta^{(j)} \in \mathbb{R}^{n_m}$ predict user j
as rating movie i with $(\theta^{(j)})^T x^{(i)}$ stars.

Problem formulation

$\theta^{(j)}$ parameter vector for user j .

$x^{(i)}$ feature vector for movie i

For user j , movie i predicted rating $(\theta^{(j)})^T x^{(i)}$

To learn $\theta^{(j)}$

The optimization objective:

To learn $\theta^{(j)}$ (parameters for user j):

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

To learn $\theta^{(j)}$ $\forall j = 1, 2, \dots, n_u$

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

gradient descent update

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

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \left(\sum ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} + \lambda \theta_k^{(j)} \right) \quad k \neq 0$$

collaborative filtering

Given $\theta^{(1)}, \dots, \theta^{(n_u)}$, to learn $x^{(i)}$:

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

given $\theta^{(1)}, \dots, \theta^{(n_u)}$, to learn $x^{(1)}, \dots, x^{(n_m)}$:

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

given $x^{(1)}, \dots, x^{(n_m)}$ (and movie ratings)
can estimate $\theta^{(1)}, \dots, \theta^{(n_u)}$

given $\theta^{(1)}, \dots, \theta^{(n_u)}$
can estimate $x^{(1)}, \dots, x^{(n_m)}$

Collaborative filtering optimization objective

Minimizing $x^{(1)} \dots x^{(n_m)}$ and $\theta^{(1)} \dots \theta^{(n_u)}$ simultaneously

$$J(x^{(1)} \dots x^{(n_m)}, \theta^{(1)} \dots \theta^{(n_u)}) = \frac{1}{2} \sum_{ij: r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2$$

$$+ \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^{n_x} (x_k^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^{n_x} (\theta_k^{(j)})^2$$

Collaborative filtering algorithm

- 1) Initialize $x^{(1)} \dots x^{(n_m)}$, $\theta^{(1)} \dots \theta^{(n_u)}$ to small random values
- 2) Minimise $J(x^{(1)} \dots x^{(n_m)}, \theta^{(1)} \dots \theta^{(n_u)})$ by gradient descent

$$x_k^{(i)} := x_k^{(i)} - \alpha \sum_{j: r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) \theta_j^{(i)} + \lambda x_k^{(i)}$$

$$\theta_k^{(i)} := \theta_k^{(i)} - \alpha \sum_{j: r(i,j)=1} ((\theta^{(i)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} + \lambda \theta_k^{(i)}$$

- 3) for a user with parameters θ and movie with (learned) features x predict a star rating $\theta^T x$

$$X = \begin{bmatrix} (x^{(1)})^T \\ \vdots \\ (x^{(n,m)})^T \end{bmatrix}$$

$$\Theta = \begin{bmatrix} (\Theta^{(1)})^T \\ \vdots \\ (\Theta^{(k)})^T \end{bmatrix}$$

$X \Theta^T$ predicted ratings

Mean normalization

$$\theta^{(j)} + x^{(i)} + \mu_i$$

Large scale machine learning.

Stochastic gradient descent

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

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

① Randomly shuffle data set

② Repeat {

for $i = 1, \dots, m$ }

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

(for $j = 0, \dots, n$) simultaneous update

}

Mini-batch gradient descent

Batch gradient descent: Use all m examples in each iteration

Stochastic gradient descent: Use 1 example in each iteration

Mini-Batch gradient descent: Use 5 examples in each iter.

Say $t = 10$, $m = 1000$

Repeat {

for $i = 1, 11, 21, \dots, 991$ }

$$\theta_j := \theta_j - \alpha \frac{1}{10} \sum_{k=i}^{i+9} (h_{\theta}(x^{(k)}) - y^{(k)}) x_j^{(k)}$$

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

}

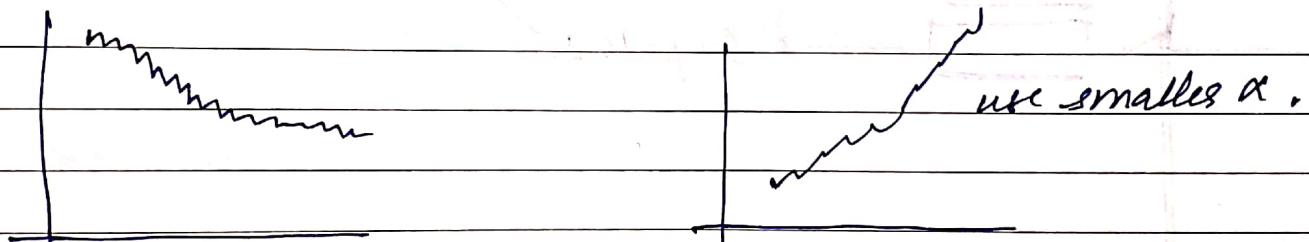
Checking for convergence

(1) Batch gradient descent :

- Plot $J_{\text{train}}(\theta)$ as a func of number of iterations of gradient descent.
- $J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$

(2) Stochastic gradient descent

- cost($\theta, (x^{(i)}, y^{(i)})$) = $\frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$
- During learning, compute cost($\theta, (x^{(i)}, y^{(i)})$) before updating θ using $(x^{(i)}, y^{(i)})$
- Every 1000 iterations (say), plot cost($\theta, (x^{(i)}, y^{(i)})$) averaged over the last 1000 examples processed by the algo.



average over
more examples.

- We can slowly decrease α over time if we want θ to converge

$$\alpha = \frac{\text{cost}}{\text{IterationNumber} + \text{const2}}$$

Online learning

Repeat forever {

set (x, y)

update θ using (x, y) :

$$\theta_j := \theta_j - \alpha (\phi(x) - y) x_j \quad j = 0 \dots n$$

}

can adapt to changing user preferences.

Map reduce

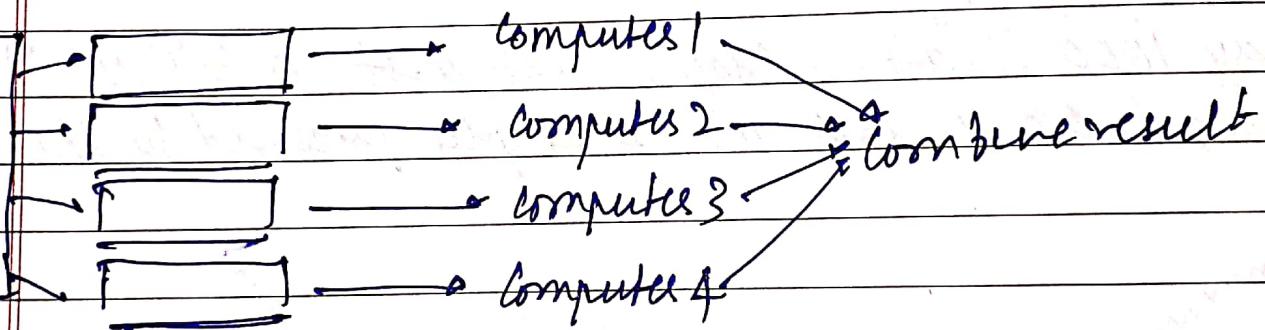
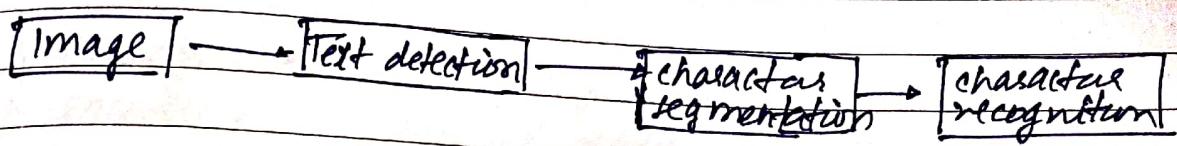


Photo OCR problem



On getting more data

- 1) Make sure you have a low bias classifier before expending the effort.
(Plot learning curves) Keep ^{increasing} number of features / number of hidden units in a neural network until you have a low bias classifier.
- 2) How much work would it be to get 10x as much data as we currently have?
 - Artificial data synthesis
 - Collect label data yourself
 - Crowd source' Amazon Mechanical task.

celling analysis.

Component	Accuracy
Overall system	72%
Text detection	89% ↓ 17%
Character segmentation	90% ↓ 1%
Character recognition	100% ↓ 10%