

Machine learning

— Coursera

Andrew Ng



Start off from Multivariate linear regression

$x_j^{(i)}$ = value of feature j in i th training example

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

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

Gradient descent: Repeat

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

Therefore cost function will be:

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=0}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)} \quad (\text{m examples})$$

Simultaneously update θ_j ($j=0, 1, \dots, n$)

We can speed up gradient descent by (it will take less iterations)
feature scaling and mean normalization.

$$x_i := \frac{x_i - \bar{x}_i}{s_i}$$

if α is sufficient small, $J(\theta)$ will decrease every iteration.

Polynomial regression:

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 \quad (\text{let } x_2 = x_1^2)$$

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

if choose feature this way, feature scaling becomes really important.

Normal equation (without using Gradient Descent)

Gradient descent	Normal equation
choose α	No need to choose α
need many iterations	No need to iterate
$O(kn^2)$	$O(n^3)$
works well if n is large	slow if very large.

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

$$\begin{array}{c|c} x_0 & y \\ x_1 & \\ x_2 & \\ x_3 & \\ \vdots & \vdots \\ \text{Some number} & \text{some number} \end{array}$$

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

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

$$= x^T \theta - \bar{y}^T \bar{y}$$

or can understand as:

$$x\theta = \bar{y}$$

$$\underline{x^T(x\theta - \bar{y}) = 0}$$

order make difference

What if $X^T X$ is invertible:

- Redundant feature
- Delete some feature, or use regularization.

Octave

x or(1,0) ans=1

% to comment

disp() (= print)

sprintf('2 decimal: %0.2f', a)

format long

A = [12; 34; 56]

V = [1; 2; 3]

V = 1:0.1:2

V = 1:6

ones(2,3)

zeros(3,3)

rand(3,3)

sqrt(a)

hist()

eye(4)

help eye.

size(A)

size(A,2) the second number in size

length(a) ⇒ vector

pwd

cd

ls

load file

load('file')

who
whos more detail command

clear <variable>

<variable>(1:10) index using ()

save hello.mat V;

save hello.txt V -ascii

A([13], :)

A = [A, [Column]]

A(:) % put all element together

C = [A B] == C = [A, B]

C = [A ; B]

A * C dot

A.*C element wise

A.^2 square

1./V every element

log(e)

V(e)

abs(a)

V+1

A' A transpose

val = max(a)

[val, ind] = max(a) val and index of a.

a < 3. ⇒ [1, 0, 1, 1] if for True or False

find(a>>) all the index of value

A = magic(3)

[r, c] = find(A >= 7)

A(2,3) = 7

sum(a)

ceil(a)

floor(a)

prod(a)

max(rand(3)), rand(3) ⇒ two signs

max(A, [], 1) first dimension max

A =
magic(3)

max(max(A))

sum(A, 1) [.....]

sum(A, >> (:))

A.*eye(9)

flipud(eye(9))

pinv(A) % inverse of A.

plotting

plot(x, y)

hold on;

plot(x, y2)

xlabel('...')

legend('sin', 'cos')

title('my plot')

save file:

print -dpng 'myplot.png'

close

figure(1); plot(t, y1);

subplot(1, 2, 1);

subplot(1, 2, 2);

axis([xmin xmax ymin ymax])

imagesc(A)

imagesc(A), colorbar, colormap gray;



for i=1:10,
 v(i) = 2^i;

end;

break

continue,

i=5

while i <= 5,

 v(i)=100;

 i = i+1

end;

if, elseif, else end;

function y = squareThisNumber(x)

 y = x^2; a file with .m

function [y1, y2] = y2(x)

::::

function J = costFunctionJ(X, y, theta)

 m = size(X, 1)

 predictions = X * theta

 sqrErrors = (predictions - y). ^ 2

 J = 1 / (2 * m) sum(sqrErrors)

Vectorization

Classification:

Do not usually use linear regression for classification problem.

will be influenced by extreme cases, we need to design a new function

Logistic regression:

$$\begin{cases} h_{\theta}(x) = g(\theta^T x) \\ g(z) = \frac{1}{1+e^{-z}} \end{cases}$$

Sigmoid function
Logistic function

$h_{\theta}(x) = P(y=1 | x; \theta)$ probability $y=1$, given x parameterized by θ

$y=1$ if $h_{\theta}(x) \geq 0.5 \Rightarrow \theta^T x \geq 0$
 $y=0$ if $h_{\theta}(x) < 0.5 \Rightarrow \theta^T x < 0$

Decision boundary: separate two cases, $h_{\theta}(x) = 0.5$
 it is a property of parameter instead of a property of a dataset

Non-linear decision boundary:

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

Cost function

Convex function

if $y=1$, but predict $y=0$, we give it big penalty.

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

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

(principle of maximum likelihood optimization to get cost function)

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

$$= -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log\left(\frac{1}{1+e^{-\theta^T x}}\right) + (1-y^{(i)}) \log\left(1-\frac{1}{1+e^{-\theta^T x}}\right)$$

update of θ_j :

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} (J(\theta)) \Rightarrow \text{compute by chain rule.}$$

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

Vectorized implementation:

$$\theta := \theta - \frac{\alpha}{m} x^T (g(x\theta) - \vec{y})$$

other optimization algorithm:

- Gradient descent
- Conjugate descent
- BFGS
- L-BFGS

No need α
 faster than gradient descent

function [jval, gradient] = costFunction(theta)

$$\vec{jval} = (\theta_0 - 5)^2 + (\theta_1 - 5)^2$$

$$\text{gradient} = \text{zeros}(1, 2)$$

$$\text{gradient}(1) = 2 * (\theta(1) - 5)$$

$$\text{gradient}(2) = 2 * (\theta(2) - 5)$$

options = optimset('GradObj', 'on', 'MaxIter', '100')

initialTheta = zeros(2, 1)

[optTheta, functionVal, exitFlag] = fminunc(@costFunction, initialTheta, options)

multi-class classification: one vs all method

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

$\max_i h_{\theta}^{(i)}(x)$ to be the correct category

Overfitting:

assumption causing bias

underfit high bias

overfitting low bias

Address overfitting:

1. Reduce number of features

- Manually select which feature to keep.

- Model selection algo

2. regularization

- keep all features reduce value/magnitude of θ_j

- works well if we have a lot of features each contribute a bit to y .

small values θ_0, θ_1

- 'simpler' hypothesis

- less prone to over fitting

$$J(\theta) = \frac{1}{m} \left[\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right] (R_2)$$

conventionally, usually start at 1

Gradient descent

Repeat {

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

$$\theta_j := \theta_j \underbrace{\left(1 - \frac{\alpha \lambda}{m}\right)}_{\downarrow 0.99 \text{ shrinking } \theta_j \text{ a little bit.}} - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Normal equation regularization

$$(X^T X + \lambda [I_{n+1}])^{-1} X^T y \quad \frac{\partial}{\partial \theta_j} J(\theta) = 0$$

optional if $\lambda > 0$, \downarrow is invertible

Regularized logistic regression:

Repeat {

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

$$\theta_j := \theta_j \underbrace{\left(1 - \frac{\alpha \lambda}{m}\right)}_{\text{(simultaneously update all } \theta_j, j=1, \dots, n\text{)}} + \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

(same as above)

Cost function:

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

function [J val, gradient] = costFunction(theta)

J val = [Code to compute $J(\theta)$]

gradient(1) = [Code to compute $\frac{\partial}{\partial \theta_0} J(\theta)$]

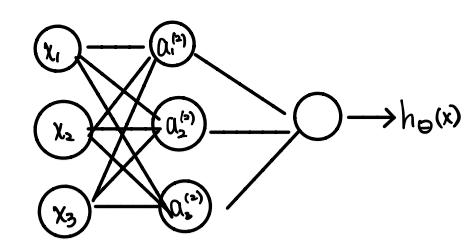
$$\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

gradient(2) = [Code to compute $\frac{\partial}{\partial \theta_1} J(\theta)$]

⋮

gradient($n+1$) = [Code to compute $\frac{\partial}{\partial \theta_n} J(\theta)$]

$$\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_n^{(i)} + \frac{\lambda}{m} \theta_{n+1}$$



$a_i^{(j)}$ = action of unit i in layer j
 $\Theta^{(j)}$ = matrix of weights controlling function mapping from layer j to layer $j+1$
 if network has s_j units in layer j , s_{j+1} units in layer $j+1$.

then $\Theta^{(j)}$ will be dimension $s_{j+1} \times (s_j + 1)$

Forward propagation:

$$z^{(j)} = \Theta^{(j)} a^{(j)} \text{ (add } a_0^{(j)} \text{ to)}$$

$$a^{(j)} = g(z^{(j)}) \text{ (g: logistic)}$$

$\Theta_{ij}^{(j)}$ from j to i
 Layer n Layer $n+1$

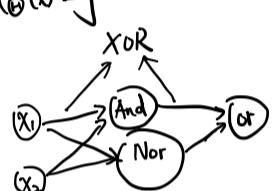
$$a_i^{(j)} = g(\Theta_{10}^{(j)} a_0^{(j)} + \Theta_{11}^{(j)} a_1^{(j)} + \Theta_{12}^{(j)} a_2^{(j)} + \Theta_{13}^{(j)} a_3^{(j)})$$

for logical function:

$$h_{\text{And}}(x) = g(-3x_0 + 2x_1 + 2x_2) \text{ AND}$$

$$h_{\text{Or}}(x) = g(-1x_0 + 2x_1 + 2x_2) \text{ OR}$$

$$h_{\text{Not}}(x) = g(10 - 2x_0) \text{ NOT}$$



Neural network (classification)

L = total layers S_l = no. units (not counting bias unit) in layer l
 k = number of output units classes $h_{\Theta}(x)_k$: hypothesis result in k th output.

$$J(\Theta) = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K [y_k^{(i)} \log(h_{\Theta}(x^{(i)})_k) + (1-y_k^{(i)}) \log(1-h_{\Theta}(x^{(i)})_k)] \\ + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{S_l} \sum_{j=1}^{S_{l+1}} (\Theta_{j,i}^{(l)})^2. \quad (\text{by convention, do not add in } x_0)$$

for multiclass classification.

back propagation:

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

$$\delta_j^{(4)} = a_j^{(4)} - y_j \quad \delta^{(4)} = a^{(4)} - y$$

$$\delta^{(3)} = (\Theta^{(3)})^T \delta^{(4)} \cdot g'(z^{(3)}) \quad a^{(3)} \cdot (1-a^{(3)})$$

prove ignored:

If ignore regularization term:

$$\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = a_j^{(l)} \delta_i^{(l+1)}$$

Back propagation algorithm:

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

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

For $i = 1 \dots m$

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

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

Using $y^{(i)}$, Compute $\delta^{(l)} = a^{(l)} - y^{(i)}$

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

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

$$\Delta_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)} \text{ if } j \neq 0$$

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

$$\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = \Delta_{ij}^{(l)}$$

(FP using $x^{(i)}$ followed by BP using $y^{(i)}$, FP using $x^{(i)}$ followed by BP using $y^{(i)}$)

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

The capital delta matrix Δ is used as an "accumulator" to add up our values as we go along and eventually compute our partial derivative. Thus we get $\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = \Delta_{ij}^{(l)}$

Forward propagation:

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

$$z^{(2)} = \Theta^{(1)} a^{(1)}$$

$$a^{(2)} = g(z^{(2)}) \text{ add } a_0^{(2)}$$

$$z^{(3)} = \Theta^{(2)} a^{(2)}$$

$$a^{(3)} = g(z^{(3)}) = h_{\Theta}(x)$$

$$\delta^{(l)} = (\Theta^{(l)})^T \delta^{(l+1)} \cdot g'(z^{(l)}) \quad \text{using chain rule.}$$

$$= ((\Theta^{(l)})^T \delta^{(l+1)}) \cdot a^{(l)} \cdot (1 - a^{(l)})$$

Neural network in practise (In Octave)

Example

$$\Theta^{(0)} \in \mathbb{R}^{10 \times 11}, \Theta^{(1)} \in \mathbb{R}^{10 \times 11}, \Theta^{(2)} \in \mathbb{R}^{10 \times 11}$$

$$D^{(1)} \in \mathbb{R}^{10 \times 11}, D^{(2)} \in \mathbb{R}^{10 \times 11}, D^{(3)} \in \mathbb{R}^{10 \times 11}$$

$$\thetaetaVec = [\Thetaeta1(:, :); \Thetaeta2(:, :); \Thetaeta3(:, :)]$$

$$DVec = [D1(:, :); D2(:, :); D3(:, :)]$$

unroll:

$$\Thetaeta1 = reshape(thetaVec(1:110), 10, 11);$$

$$\Thetaeta2 = reshape(thetaVec(111:220), 10, 11);$$

$$\Thetaeta3 = reshape(thetaVec(221:231), 10, 11);$$

Learning Algo:

initial parameters $\Theta^{(0)}, \Theta^{(1)}, \Theta^{(2)}$

Unroll to get initialTheta to pass to

fminunc(@costFunction, initialTheta, options)

Function [Jval, gradientVec] = costFunction(thetaVec)

From thetaVec, get $\Theta^{(0)}, \Theta^{(1)}, \Theta^{(2)}$

Use forward, backward prop to compute $D^{(1)}, D^{(2)}, D^{(3)}$ and $J(\Theta)$.

Unroll $D^{(1)}, D^{(2)}, D^{(3)}$ to get gradientVec.

Numeric estimation of gradients:

$$\text{gradApprox} = \frac{(J(\Theta + \epsilon) - J(\Theta - \epsilon))}{2\epsilon}$$

Remember to do gradient checking.

epsilon = 10^{-4} .

for $i = 1:n$

thetaPlus = theta;

thetaPlus(i) += epsilon;

thetaMinus = theta;

thetaMinus(i) -= epsilon;

$$\text{gradApprox}(i) = (J(\theta + \epsilon) - J(\theta - \epsilon)) / (2 * \epsilon)$$

end;

Neural network training:

Random initialization:

PROBLEM:

zero initialization:

$$a_1^{(2)} = 0_2^{(2)}, \text{ Also } \delta_1^{(2)} = \delta_2^{(2)}$$

$$\frac{\partial}{\partial \theta_{j1}^{(1)}} J(\theta) = \frac{\partial}{\partial \theta_{j1}^{(2)}} J(\theta)$$

After update, parameters are still the same.

→ therefore, pick theta value $\in (-\varepsilon, \varepsilon)$.

(e.g. $\text{rand}(0, b) * (2 * \text{init_epsilon}) - \text{init_epsilon}$;

Training a neural network

1. randomly initialize weights
2. implements forward propagation to get $h_\theta(x^{(i)})$ for any $x^{(i)}$
3. implement code to compute $J(\theta)$
4. implement backprop to compute partial derivatives $\frac{\partial}{\partial \theta_{j1}^{(1)}} J(\theta)$

→ for $i = 1:m$

Perform forward and backward prop using $(x^{(i)}, y^{(i)})$

(Back prop algorithm)

5. Use gradient checking to compare $\frac{\partial}{\partial \theta_{jk}^{(1)}} (J(\theta))$ computed

using backprop vs numeric gradient of $J(\theta)$

Then disable checking code.

6. Use gradient descent or advanced optimization method
with backprop to try to minimize $J(\theta)$. (not convex)

Week 6.

Evaluating an algorithm:

training set test set
70% 30%

$m_{\text{test}} = \text{No. of test set cases}$

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

$$J_{\text{test}}(\theta) = \frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} y_{\text{test}}^{(i)} \log(h_{\theta}(x^{(i)})) + (1-y_{\text{test}}^{(i)}) \log(1-h_{\theta}(x^{(i)})) \quad \text{classification}$$

Misclassification error

$$\text{err}(h_{\theta}(x), y) = \begin{cases} 1 & \text{if } h_{\theta}(x) \geq 0.5 \quad y=0 \quad \text{OR} \quad h_{\theta}(x) \leq 0.5 \quad y=1 \\ 0 & \end{cases}$$

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

choose model selection on valid set.

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

others is same as others.

For this example:

1. optimize parameter θ using training set.
2. find polynomial degree d with least error using cross validation set.
3. estimate generalization error using test set with $J_{\text{test}}(\theta^d)$

Bias (underfit):

$$\left. \begin{array}{l} J_{\text{train}}(\theta) \text{ will be high} \\ J_{\text{cv}}(\theta) \approx J_{\text{train}}(\theta) \end{array} \right\}$$

Variance (overfit)

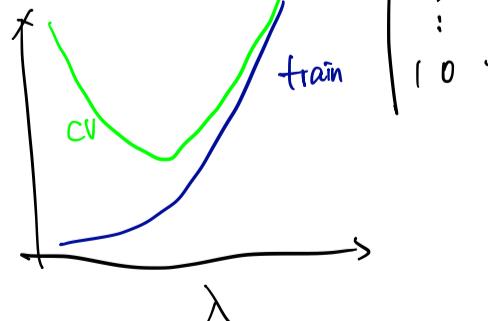
$J_{\text{train}}(\theta)$ will be low

$$J_{\text{cv}}(\theta) \gg J_{\text{train}}(\theta)$$

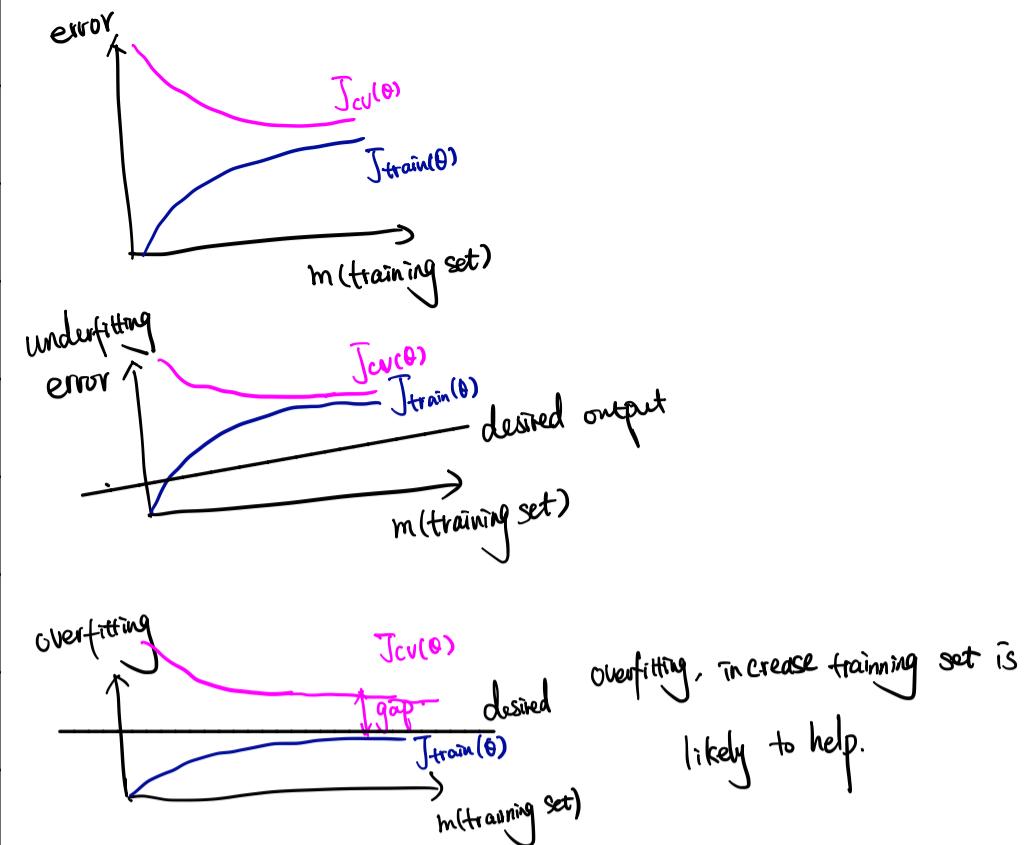
$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^m \theta_j^2$$

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

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



$$\lambda = \begin{cases} 0 \\ 0.01 \\ 0.02 \\ 0.04 \\ \vdots \\ 10 \end{cases}$$



Get more training examples : Fix high variance

try small sets of features : Fix high variance

try extra features : Fix high bias

adding polynomial features Fix high bias

Decreasing λ Fix high bias

Increasing λ Fix high variance

Building Spam classifier.

prioritizing what to work on:

- collect lots of data
- Develop sophisticated features (using email header data in spam emails)
- Develop algorithms to process your input in different ways (recognizing misspellings)

Recommended approach:

- start with simple algorithm you can implement quickly.
Implement it and test it on your cross validation data.
- Plot learning curves to decide if more data, more features, etc are likely to help.
- Error analysis: Manually examine examples (in cross validation set) that your algorithms made errors on. See if you spot any systematic trend in what type of examples it is making errors on.

Skewed dataset problem:

accuracy: $\frac{\text{correct cases}}{\text{all cases}}$

$$\text{precision} = P = \frac{TP}{TP+FP}$$

$$\text{recall} = R = \frac{TP}{TP+FN}$$

$$F_1 = \frac{2}{\frac{1}{R} + \frac{1}{P}} = \frac{2PR}{P+R} = \frac{2TP}{2TP+FP+FN}$$

$$\text{F measure} = (\beta^2+1) \frac{P+R}{\beta^2 P + R}$$

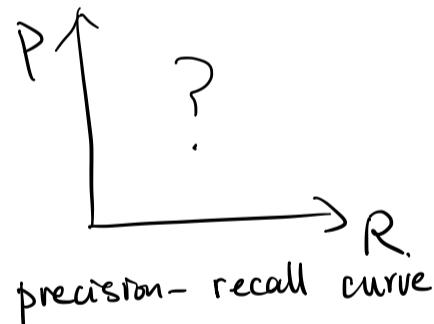
	Relevant	Non-relevant
Retrieved	TP	FP
Not Retrieved	FN	TN

Last letter stands for Predicted
Combined letter stands for actual

Trading off precision and recall.

predict 1 if $h_\theta(x) \geq 0.9$ higher precision, lower recall.
predict 0 otherwise

Vice versa.



threshold affects precision & recall.

Large data rationale:

Useful test: Given input x , can human expert confidently predict y . If can, a large dataset can help.

1. Use low bias algorithms.

2. use very large training set \rightarrow low variance

A large training set is unlikely to help when:

- ① humans cannot predict accurately.
- ② too large number of features (low bias)
- ③ we're not using regularization.

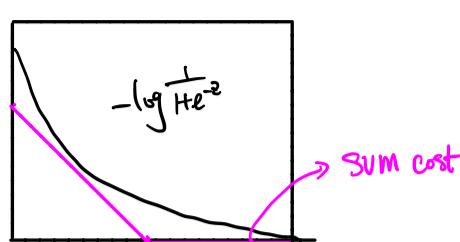
1. human test.

2. large low bias algorithm + large training set.

can ensure high accuracy.

SVM.

Logistic regression cost:



Logistic regression cost:

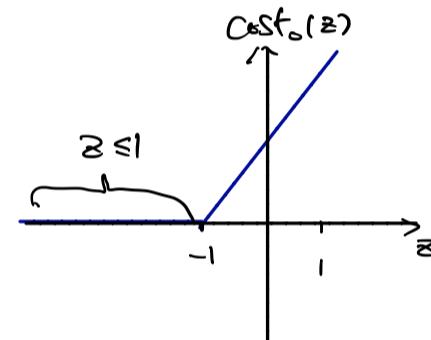
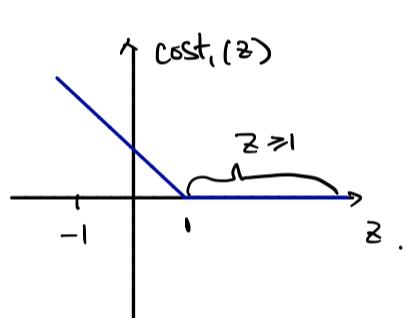
$$J(\theta) = \frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log(h_\theta(x^{(i)})) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

Sum: (Large margin classifier)

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

C is $\frac{1}{\lambda}$ also controls regularization.

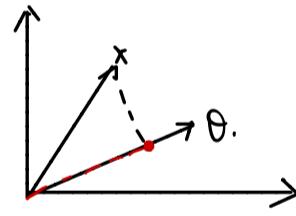
$h_\theta(x) \begin{cases} 1 & \text{if } \theta^T x \geq 0 \\ 0 & \text{otherwise.} \end{cases}$ For SUMs.



If $y=1$, we want $\theta^T x \geq 1$ (not just ≥ 0)

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

s.t. $\theta^T x^{(i)} \geq 1$ if $y^{(i)}=1$
 $\theta^T x^{(i)} \leq -1$ if $y^{(i)}=0$



Decision boundary is always perpendicular to θ .

Intuition:

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2 = \frac{1}{2} \|\theta\|^2 \quad p \text{ is projection } x \text{ onto } \theta.$$

s.t. $p^{(i)} \cdot \|\theta\| \geq 1$, if $y^{(i)}=1$
 $p^{(i)} \cdot \|\theta\| \leq -1$, if $y^{(i)}=0$

Kernel for sums:

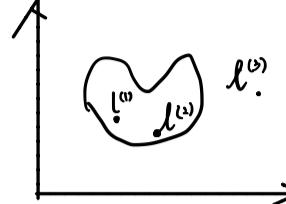
high polynomial feature function:

choose lemma $l^{(1)}, l^{(2)}, \dots, l^{(n)}$

Gaussian kernel:

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

bigger σ^2 , the extend get more spread out.



(logistic regression does not work well with kernel)

choose lemma l as training set.

we get feature vector

$$f_i = \text{similarity}(x^{(i)}, l^{(i)}) \quad f = \begin{bmatrix} f_0 \\ f_1 \\ \vdots \\ f_m \end{bmatrix} \quad f_0 = 1$$

predict " $y=1$ " if $\theta^T f \geq 0$

Training:

$$\min_{\theta} C \sum_{i=1}^m y^{(i)} \text{cost}_1(\theta^T f^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T f^{(i)}) + \frac{1}{2} \sum_{j=1}^n \theta_j^2 \quad (n=m)$$

$$\sum_j \theta_j^2 = \theta^T \theta \Rightarrow \theta^T M \theta \quad (\text{ignore } \theta_0)$$

This is implemented to allow fast calculation.

C ($= \frac{1}{\lambda}$): Large C : low bias, high variance.
 Small C : high bias, low variance.

σ^2 : Large σ^2 : Feature vary more smoothly. higher bias, lower variance.
 Smaller σ^2 : higher variance, low bias

Need to specify: ① choice of parameter C
 ② choice of kernel.)

Gaussian kernel.: need to choose σ^2

Note: perform feature scaling before using the Gaussian kernel.

(need to satisfy technical condition called "Mercer's Theorem". to converge.

Kernel available:

- Polynomial kernel. $k(x, l) = (x^T l + \text{const})^{\text{degree}}$
- More esoteric: String kernel, chi square kernel, histogram intersection kernel...

For multiclass, rest vs all.

$n = \# \text{ of features}$. $m = \# \text{ of training examples}$

If n is large. (compare to m) \Rightarrow logistic regression, linear kernel.

If n is small. m is intermediate. Use sum with Gaussian Kernel.

If n is small, m is large, add on more features. Use logistic R or linear kernel.

Unsupervised learning:

Clustering algorithm:

K-means clustering algorithm:

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

k clusters

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

K-means algorithm

Randomly initializes k cluster centroids $\mu_1, \dots, \mu_k \in \mathbb{R}^n$.

Repeat f

for $i = 1$ to m

$c^{(i)} := \text{index}(\text{from } 1 \text{ to } k) \text{ of cluster centroid}$

closest to $x^{(i)}$ ($\min_{k=1}^K \|x^{(i)} - \mu_k\|^2$)

for $k = 1$ to k

$\curvearrowright c^{(i)}$

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

} if one cluster have 0 points, then it should be eliminated
or reinitialize it.

K-means algorithm objective.

$c^{(i)} = \text{index of cluster which } x^{(i)} \text{ is assigned.}$

$\mu_k = \text{cluster of centroid } k (\mu_k \in \mathbb{R}^n)$

$\mu_{c^{(i)}} = \text{cluster centroid of cluster to which example } x^{(i)} \text{ assigned}$

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

if J increase, there is something wrong in algorithm.

K-means random initialization.

local optima.

using training example $x^{(i)}$ as μ_k .

more random initialization. and choose one have

least cost.

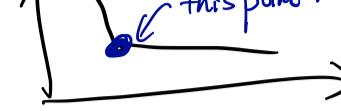
For $k \leq 10$, multiple random initialization may help

a lot but for $k \geq 10$, it may not help so much.

choose number of K .

usually manually picking.

Elbow method



or depend on downstream purpose / human insight.
later.



PCA (Principal Component analysis)

Find k vectors $u^{(1)}, u^{(2)}, \dots, u^{(k)}$ onto which to project data, so as to minimize error.
span of

$$\text{Cost function: } J = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2.$$

Principle Component Analysis Algorithm.

Reduce data from n dimensions to k -dimensions

Compute "Covariance matrix":

$$\Sigma \rightarrow \Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)})(x^{(i)})^T \Rightarrow (n \times n) \text{ matrix}$$

Compute "eigenvectors" of matrix Σ :

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

$$U = \begin{bmatrix} | & | & | \\ u^{(1)} & u^{(2)} & \dots & u^{(n)} \\ | & | & | \end{bmatrix} \in \mathbb{R}^{n \times n} \Rightarrow \text{we get first } k \text{ columns of data.}$$

$$Z = \begin{bmatrix} | & | & | \\ u^{(1)} & u^{(2)} & \dots & u^{(n)} \\ | & | & | \end{bmatrix}^T X = \begin{bmatrix} | & | & | \\ u^{(1)} & u^{(2)} & \dots & u^{(n)} \\ | & | & | \end{bmatrix} X^{(i)} \Rightarrow (k \times 1) \text{ matrix.}$$

Note:

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

Octave implementation:

$$\Sigma = \frac{1}{m} X' * X;$$

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

$$U_{\text{reduce}} = U(:, 1:k); \quad \text{ignoring } x_0=1.$$

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

Reconstruction from compressed data:

$$\underbrace{x_{\text{approx}}}_{n \times 1} = \underbrace{U_{\text{reduce}}}_{n \times k} \cdot \underbrace{Z}_{k \times 1}.$$

how to choose k :

choose k to be smallest:

$$\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq 0.01$$

I would say 99% percent of variance is retained.

$$[U, S, V] = \text{svd}(\Sigma)$$

$$S = \begin{bmatrix} S_{11} & & 0 \\ & \ddots & \\ 0 & & S_{nn} \end{bmatrix}$$

$$\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} = 1 - \frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}} \leq 0.01.$$

advice for applying PCA:

Supervised learning speedup.

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

$$\downarrow$$

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

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

Application of PCA:

— Compression:

- Reduce memory
- speed up learning algorithm.

— Visualization:

It is not a good way to avoid overfitting.

before implementing PCA, try running with raw data first.
if cannot, can use PCA.

Anomaly detection:

problem: x_{test} is anomalous?

$p(x_{\text{test}}) < \varepsilon \rightarrow \text{flag anomaly}$

$p(x_{\text{test}}) \geq \varepsilon \rightarrow \text{OK.}$

Fraud Detection ; Manufacturing ; Monitoring computer

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

Moder $P(x)$

Judge / identify.

Gaussian distribution:

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

Parameter estimation: $N(\mu, \sigma^2)$

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

($\frac{1}{m-1}$ in statistics)

Anomaly detection algorithm,

1. Choose features that you think might be example of anomaly.

2. Fit $\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2$.

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

3. Given new example, Compute $p(x)$:

$$p(x) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi}\sigma_j} \exp\left\{-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right\}$$

Normal evaluation

training data: unlabelled 6000 good engines (can contain some anomaly)
CV. labelled dataset 2000 good 20 bad.

$$y = \begin{cases} 1 & \text{if } p(x) < \varepsilon \text{ (anomaly)} \\ 0 & \text{if } p(x) \geq \varepsilon \text{ (OK)} \end{cases}$$

possible evaluation metrics: Precision, recall, F₁-score.

choose ε on valid set.

Anomaly detection vs. Supervised learning:

Anomaly detection | supervised learning

<ol style="list-style-type: none"> ① very few anomalies ② large number normal ③ many types of anomalies, future ones might be never met before learn from normal ones 	<ol style="list-style-type: none"> ① large number of positive & negative ② enough normal to learn what normal looks like.
--	---

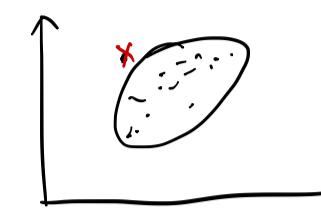
What data to use.

① make feature more gaussian, plot it in hist.

$$\log(x) \cdot x^{0.5}$$

② more feature distinguish anomaly from normal.

→ create new features from existing features.



to address this problem.

use Multivariate Gaussian (Normal) distribution

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

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

$$\Sigma = \begin{bmatrix} \sigma_1^2 & & & 0 \\ 0 & \sigma_2^2 & & \dots \\ & & \ddots & \sigma_n^2 \end{bmatrix} \text{ is original model.}$$

Original model = assume no correlation between different models.
Computationally cheap.

For multivariable gaussian model: must $m > n$, or Σ is non-invertible
If Σ is non-invertible:

① $m < n$

② linearly dependent features.

Recommendation systems:

Movie recommendation:

$n_u = \text{no. users}$ $n_m = \text{no. movies}$ $r(i,j) = 1$ if user j has rated movie i .

$y^{(i,j)}$ = rating given by user j to movie i .

Problem formulation:

$r(i,j) = 1$ if user j has rated movie i (0 otherwise)

$y^{(i,j)}$ = rating by user j on movie i (if defined)

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

$\chi^{(i)}$ = parameter vector for movie i

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

$m^{(j)} = \text{no. of movies rated by user } j$

+ to learn $\theta^{(j)}$:

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

Total optimization learning algorithm:

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

all the movie to one user

Gradient descent update:

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

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \sum_{i:r(i,j)=1} ((\theta^{(j)})^T \chi^{(i)} - y^{(i,j)}) \chi_k^{(i)} + \lambda \theta_k^{(j)} \quad (k \neq 0)$$

pretty much as logistic regression.

Collaborative Filtering:

if user $\theta^{(1)} \dots \theta^{(n_u)}$, to learn $\chi^{(i)}$:

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

to learn $\chi^{(1)} \dots \chi^{(n_m)}$:

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

↓ This is j . now.

Guess $\theta \rightarrow \chi \rightarrow \theta \rightarrow \chi \rightarrow \theta \dots$ simultaneously learn θ and χ .

Minimizing $\chi^{(1)}, \dots, \chi^{(n_m)}$, and $\theta^{(1)}, \dots, \theta^{(n_u)}$ simultaneously:

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

1. Initialize $\chi^{(1)}, \dots, \chi^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)}$ to random small value.

$$\begin{aligned} \chi_k^{(i)} &:= \chi_k^{(i)} - \alpha \left[\sum_{j:r(i,j)=1} ((\theta^{(j)})^T \chi^{(i)} - y^{(i,j)}) \theta_k^{(j)} + \lambda \chi_k^{(i)} \right] \\ \theta_k^{(j)} &:= \theta_k^{(j)} - \alpha \left[\sum_{i:r(i,j)=1} ((\theta^{(j)})^T \chi^{(i)} - y^{(i,j)}) \chi_k^{(i)} + \lambda \theta_k^{(j)} \right] \end{aligned}$$

Note: There is no $\chi_0 = 1$ here. so we do not write like that.

3. predict. star rating $\theta^T \chi$.

$$\chi = \begin{bmatrix} -(\chi^{(1)})^T \\ -(\chi^{(2)})^T \\ \vdots \\ -(\chi^{(n_m)})^T \end{bmatrix} \quad \textcircled{H} = \begin{bmatrix} -(\theta^{(1)})^T \\ -(\theta^{(2)})^T \\ \vdots \\ -(\theta^{(n_u)})^T \end{bmatrix}$$

Predicted ratings: $\chi \textcircled{H}^T$ Low rank matrix factorization.

Find related movies:

small $\|\chi^{(i)} - \chi^{(j)}\| \rightarrow$ similar movies.

mean normalization:

minus the mean value in that column,

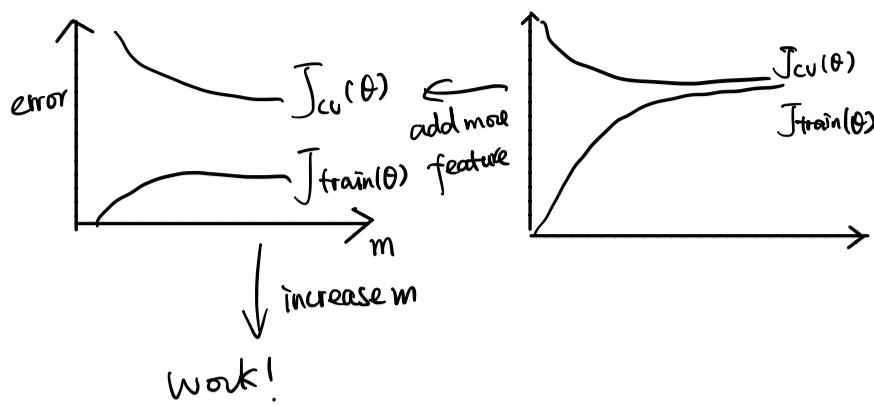
For user j , the rating score,

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

because ratings are already comparable, so
do not need to do range normalization.

Machine learning with large dataset.

with small dataset first.



Stochastic gradient descent: (use 1 example each iteration)

(Batch gradient descent. whole batch) whole dataset.

1. random shuffle (reorder) training examples

2. repeat { (for some iterations)

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

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

(for all j)

}

Mini batch gradient descent:

let $b=10, m=1000$.

Repeat {

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

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

{

↑
one advantage parallel computation.

Checking for convergence.

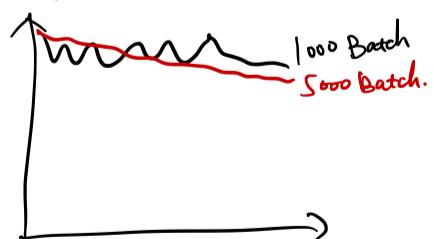
SGD:

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

During learning compute $\text{Cost}(\theta, (x^{(i)}, y^{(i)}))$ before updating θ .

Every 1000 iterations, plot average of cost.

Checking for convergence for $J(\theta)$.



Online learning:

Product search

User search for "Android phone 1080 camera", give 10 results, 10 training pair.
learn $p(y=1/x; \theta)$, update parameters.

① change update with time, of user need.

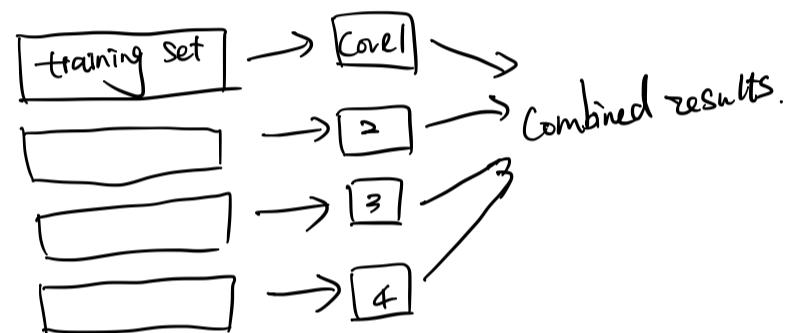
② update 1 at a time, and disregard that training set.

Map reduce:

(Jeffrey Dean) legendary person
google.
split work evenly to
different machines and send to one central machine.

Map reduce needs learning algorithm to perform summation over training set.
Many learning algorithm can be expressed as computing sum of functions
over training set.

Multi-core machine:



OCR
Pipeline
Sliding windows

Artificial data:

- ① Create Data (OCR font)
- ② Data amplification. (Audio, background knowledge)

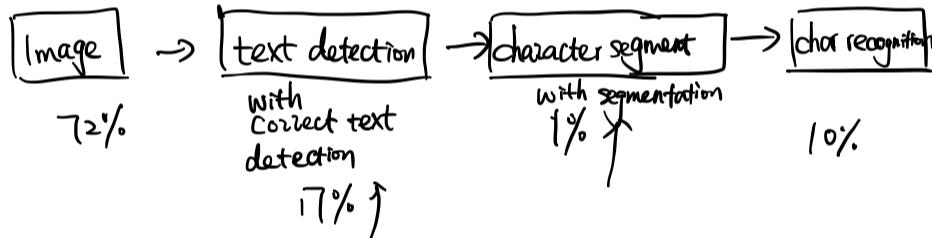
Data:

1. make sure you have a low bias classifier before expanding effort. (plot learning curves)
E.g. keep increasing # of features / hidden units in 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
- Label / collect yourself.
- "crowd source" (E.g. Amazon Mechanical Turk)

Ceiling analysis:

give correct labels:



pipeline is very important.

Summary : Main topics

Supervised learning:

Linear regression, logistic regression, neural networks, SVMs.

Unsupervised learning:

k-means, PCA, Anomaly detection.

Special Applications / special topics

— Recommender system, large scale machine learning.

Advice on building a machine learning system

- Bias/Variance, regularization, deciding what to do next, evaluation of learning algorithms, learning curves, error analysis, ceiling analysis.