

application of ML \Rightarrow self driving car, computer vision, etc

What is ML \Rightarrow field of study that gives computer

Arthur Samuel \Rightarrow the ability to learn without being explicitly programmed

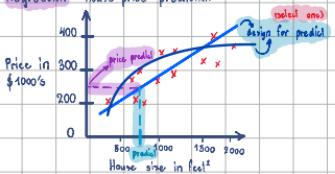
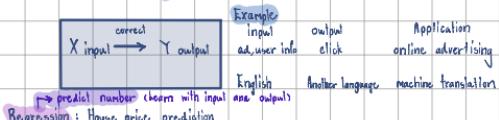
ML has 4 types 1) Supervised learning (Used most in real world)

2) Unsupervised learning

3) Recommender systems

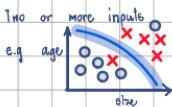
4) Reinforcement learning

1) Supervised learning \Rightarrow learn with example with right answer



Classification \Rightarrow predict categories

output \Rightarrow class, category (small number of possible outputs)



2. Unsupervised learning \Rightarrow make algorithm figure by yourself
 $|X \text{ input} \neq \text{no has output } Y|$

clustering \Rightarrow group of data that has same or similar topics

Anomaly detection \Rightarrow find unusual data points

Dimensionality reduction \Rightarrow Compress data using fewer numbers



Terminology \Rightarrow data used to train the model

x = "input" variable or feature

y = "output" variable or "target" variable

(x_i, y_i) = single training example

$(x^{(i)}, y^{(i)})$ = i^{th} training example

m = number of training examples

flowchart: training set \rightarrow features \rightarrow targets

learning algorithm

or hypothesis \rightarrow function \rightarrow \hat{y} ("y-hat" estimated y)

feature \rightarrow function \rightarrow \hat{y} ("y-hat" estimated y)

How to represent f ? $f_{w,b}(x) = wx + b$ or weights, parameters, coefficients

model generates best linear function but not just only linear

you can use curve, line for more accuracy

univariate linear regression

call it squared error cost function

Cost function \Rightarrow tell how well the model is doing
 find w, b so $\hat{y}^{(i)}$ is close to $y^{(i)}$ for all $(x^{(i)}, y^{(i)})$

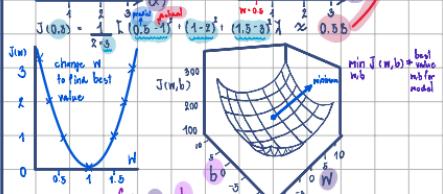
$\hat{y} - y$ call error
 $J(w, b) = \frac{1}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$

from $\hat{y}^{(i)} = f_{w,b}(x^{(i)})$ $J(w, b) = \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$

$\frac{1}{m}$ = calculate average value

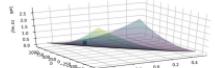
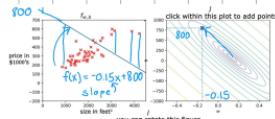
1 \Rightarrow make the derivative cleaner (clean)

Idea or concept behind the formula



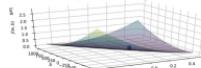
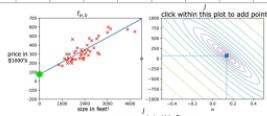
contour plots \Rightarrow a convenient way to visualize 3-D cost function $J(w)$

but if plotted in 2D



Stanford ONLINE @deeplearning.AI

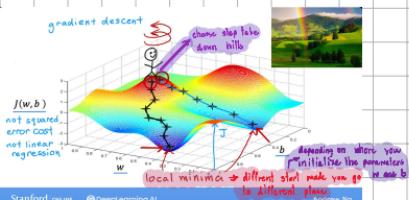
Andrew Ng



Stanford ONLINE @deeplearning.AI

Andrew Ng

Gradient Descent → use in every machine learning, using by keep changing w, b to reduce $J(w, b)$ until settle at or near a minimum.



Stanford ONLINE @deeplearning.AI

Andrew Ng

Implementing Gradient Algorithm

Learning rate α controls size your take down hills

$$w = w - \alpha \frac{\partial}{\partial w} J(w, b)$$

Derivative → what direction you want to take your step

$$b = b - \alpha \frac{\partial}{\partial b} J(w, b)$$

ab

Simultaneously update w and b → update 2 parameter real time

$$\text{temp_}w = w - \alpha \frac{\partial}{\partial w} J(w, b)$$

$$\text{temp_}b = b - \alpha \frac{\partial}{\partial b} J(w, b)$$

$$w = \text{temp_}w \quad b = \text{temp_}b$$

We want to update w before using it
calculate temp_b

} not fit

Gradient Descent Intuition

$$w = w - \alpha \frac{\partial}{\partial w} J(w)$$

$w \leftarrow w - \alpha \cdot (\text{positive num})$
make you close to local minimum

$$w = w - \alpha \cdot (\text{negative num})$$

ratio of negative direction

Learning Rate

(long time) too small & small step using

long time to achieve local minima (over shoot)
 $w = w - \alpha \cdot c \cdot d$

too large & skip or pass local minima, never reach minimum

} fit

Gradient Descent for Linear Regression

Linear regression model

$$f_{w,b}(x) = wx + b \quad J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

Gradient descent algorithm

repeat until convergence {

$$w = w - \alpha \frac{\partial}{\partial w} J(w, b) \rightarrow \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})x^{(i)}$$

$$b = b - \alpha \frac{\partial}{\partial b} J(w, b) \rightarrow \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})$$

(Optional)

$$\frac{\partial}{\partial w} J(w, b) = \frac{1}{m} \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2 = \frac{1}{m} \frac{1}{2m} \sum_{i=1}^m (wx^{(i)} + b - y^{(i)})^2$$

$$= \frac{1}{m} \sum_{i=1}^m (wx^{(i)} + b - y^{(i)})^2 \times \overset{(i)}{x^{(i)}} = \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})x^{(i)}$$

$$\frac{\partial}{\partial b} J(w, b) = \frac{1}{m} \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2 = \frac{1}{m} \frac{1}{2m} \sum_{i=1}^m (wx^{(i)} + b - y^{(i)})^2$$

$$= \frac{1}{m} \sum_{i=1}^m (wx^{(i)} + b - y^{(i)}) \cancel{x^{(i)}} = \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})$$

corner function → bowl shape (It has global minimum)



Running Gradient Descent = "Batch" gradient descent

Batch → each step of gradient descent uses all the training examples

Made linear regression more powerful

Multiple Features (variables)

$$\text{original} \quad y = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b$$

but with multiple features using by this

x_j = j^{th} feature

n = number of features

\vec{x} = features of i^{th} training example

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

e.g. $x_1 \ x_2 \ x_3 \ x_4 \ \dots \ x_n \ \text{feature } (n=4)$

	x_1	x_2	x_3	x_4
1000	10	2	5	
1500	30	1	6	
2000	30	3	1	
$\vec{x}^{(3)}$	1	3	0	1

$$\text{model} \quad f_{\vec{w}, b}(\vec{x}) = \vec{w}_1 x_1 + \vec{w}_2 x_2 + \dots + \vec{w}_n x_n + b$$

$\vec{w} = [w_1, w_2, \dots, w_n]^T$ linear algebra
count from 1

b is a number

$\vec{x} = [x_1, x_2, \dots, x_n]^T$ dot product

while again $f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$

Vectorization → made code shorter and enable gpu hardware

using for manage vector allow code to run more efficiently

$$f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b \quad f = \text{np.dot}(\vec{w}, \vec{x}) + b$$

Why it runs so fast?

Without vectorization

```
for j in range(0, 16):
    z = z + w[0] * x[0]
    ...
    z = z + w[15] * x[15]
```

Vectorization

```
z = np.dot(w, x)
```

$w[0]$	$w[1]$	\dots	$w[15]$
$x[0]$	$x[1]$	\dots	$x[15]$

in parallel

$$z = w[0]*x[0] + w[1]*x[1] + \dots + w[15]*x[15]$$

efficient → scale to large datasets

Gradient descent

```
w = (w1, w2, ..., w16)  $\rightarrow$  parameters
d = (d1, d2, ..., d16)
w = np.array([0.5, 1.3, ..., 3.4])
d = np.array([0.3, 0.2, ..., 0.4])
```

compute $w_j = w_j - 0.1d_j$ for $j = 1 \dots 16$

Without vectorization

```
w1 = w1 - 0.1d1
w2 = w2 - 0.1d2
...
w16 = w16 - 0.1d16
```

for j in range(0, 16):
 w[j] = w[j] - 0.1 * d[j]

With vectorization

$w = \vec{w} - 0.1 \vec{d}$

$w = w - 0.1 \cdot d$

Gradient Descent for Multiple Regression (features > 2)

$$\begin{aligned} w_0 &= w_0 - \alpha \frac{1}{M} \sum_{i=1}^M (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_0^{(i)} \\ &\vdots \\ w_n &= w_n - \alpha \frac{1}{M} \sum_{i=1}^M (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_n^{(i)} \\ b &= b - \alpha \frac{1}{M} \sum_{i=1}^M (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) \end{aligned}$$

(An alternative to linear regression)

- Normal Equation 1) Only for linear regression

2) Solve for w, b without iterations

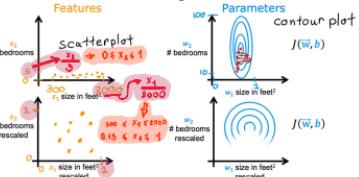
Disadvantages 1) Doesn't generalize with other learning algorithms

2) Slow when number of feature is large ($n > 10,000$)

Gradient descent is the recommended method to find w, b

Feature Scaling ⇒ make gradient descent run much faster

Feature size and gradient descent



divide maximum value is not only way to feature scaling

another way call Mean normalization

$$x_i' = \frac{x_i - \mu_i}{\sigma_i} \quad \mu_i = \text{average value of all data in } i^{\text{th}} \text{ feature}$$



And another method is Z-score normalization

if you want for data points to have standard deviation = 1 (calculate in each features)

$$x_i = \frac{x_i - \mu_i}{\sigma_i}$$

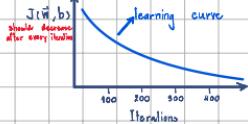
$$\sigma = \sqrt{\frac{1}{N-1} \sum_{i=1}^N (x_i - \mu)^2}$$



aim for about $-1 \leq x_i \leq 1$ for each features (acceptable ranges).

Checking Gradient Descent for convergence

+ make sure gradient descent is working correctly



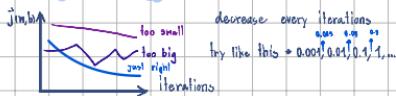
Automatic convergence test

$$\epsilon = "epsilon" \approx 10^{-3}$$

If $J(w, b)$ decreases by ϵ in one iteration,

stop model training

Choosing the learning Rate \Rightarrow select from $J(\vec{w}, b)$ will



Feature Engineering

$$\text{e.g. } f_{\vec{w}, b}(\vec{x}) = \vec{w}_1 x_1 + \vec{w}_2 x_2 + b$$

transforming or combining original features
new feature

multiple linear regression feature engineering

$$f_{\vec{w}, b}(x) = \vec{w}_1 x_1 + \vec{w}_2 x_2 + b$$

or $f_{\vec{w}, b}(\vec{x}) = w_1 x_1 + w_2 x_2 + b$

Choice of features

$$f_{\vec{w}, b}(x) = w_1 x_1 + w_2 x_2 + b$$

slope

Never forget to study the last 3 optional tabs in Module 2

Classification with logistic regression

1) binary classification $\Rightarrow 2$ possible class or category

e.g. is this email spam? yes/no

Threshold \hat{y} : decision line

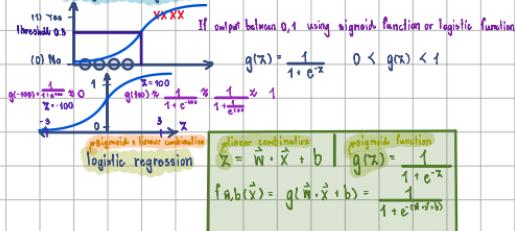
False 0: misclassified

True 1: correct

if $f_{\vec{w}, b}(\vec{x}) < 0.5 \Rightarrow \hat{y} = 0$

if $f_{\vec{w}, b}(\vec{x}) \geq 0.5 \Rightarrow \hat{y} = 1$

2) logistic Regression (use most in real world)



Decision Boundary

$$f_{\vec{w}, b}(\vec{x}) - g(\vec{x}) = g(\vec{w}^T \vec{x} + b) - 1$$

$$= \vec{w}_1 x_1 + \vec{w}_2 x_2 + b - 1 = 0$$

$$= \vec{w}_1 x_1 + \vec{w}_2 x_2 - 1 = 0$$

$$= \vec{w}_1 x_1 + \vec{w}_2 x_2 = 1$$

$$= \vec{w}_1 x_1 + \vec{w}_2 x_2 = 1$$

$$\text{Yes: } \hat{y} = 1$$

$$f_{\vec{w}, b}(\vec{x}) \geq 0.5$$

$$g(\vec{x}) \geq 0.5$$

$$\vec{w}^T \vec{x} + b \leq 0$$

$$\vec{w}^T \vec{x} + b \geq 0$$

$$\hat{y} = 0$$

Non-linear decision boundaries

$$f_{\vec{w}, b}(\vec{x}) = g(\vec{x}) = \frac{1}{1 + e^{-(\vec{w}_1 x_1^2 + \vec{w}_2 x_2^2 - 1)}}$$

$$\vec{x} = x_1^2 + x_2^2 - 1 = 0$$

$$x_1^2 + x_2^2 = 1$$

Cost Function for logistic Regression \Rightarrow How to choose \vec{w} and b

Squared error cost $J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2$

linear regression $f_{\vec{w}, b}(\vec{x}) = \vec{w}^T \vec{x} + b$

logistic regression $f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w}^T \vec{x} + b)}}$

loss is lowest when $f_{\vec{w}, b}(\vec{x}^{(i)})$ predicts close to true label $y^{(i)}$

$h(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = \begin{cases} -\log(f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 1 \\ -\log(1 - f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 0 \end{cases}$

check how close the predicted value is to the true label $y^{(i)}$

example: model predicts 0.01 $\rightarrow -\log(0.01) \approx 4.6$ (High loss)

model predicts 0.99 $\rightarrow -\log(1 - 0.99) \approx 0.01$ (Low loss)

loss is lowest when $f_{\vec{w}, b}(\vec{x}^{(i)})$ predicts close to true label $y^{(i)}$

$h(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = \begin{cases} -\log(f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 1 \\ -\log(1 - f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 0 \end{cases}$

check predicted value closer to 1 than 0

check predicted value closer to 0 than 1

Simplified Cost Function

$$\text{original } J_{\hat{w}, b}(\hat{x}^{(i)}, y^{(i)}) = -y^{(i)} \log(f_{\hat{w}, b}(\hat{x}^{(i)})) - (1 - y^{(i)}) \log(1 - f_{\hat{w}, b}(\hat{x}^{(i)}))$$

$$\text{if } y^{(i)} = 1: J_{\hat{w}, b}(\hat{x}^{(i)}, y^{(i)}) = -\log(f_{\hat{w}, b}(\hat{x}^{(i)}))$$

$$\text{if } y^{(i)} = 0: J_{\hat{w}, b}(\hat{x}^{(i)}, y^{(i)}) = -\log(1 - f_{\hat{w}, b}(\hat{x}^{(i)}))$$

from cost function: $J(\hat{w}, b) = \frac{1}{M} \sum_{i=1}^M [J_{\hat{w}, b}(\hat{x}^{(i)}, y^{(i)})]$

$$= \frac{1}{M} \sum_{i=1}^M [y^{(i)} \log(f_{\hat{w}, b}(\hat{x}^{(i)})) + (1 - y^{(i)}) \log(1 - f_{\hat{w}, b}(\hat{x}^{(i)}))]$$

Gradient Descent Implementation

linear regression same like logistic regression

but here this is different

- 1) linear regression
- 2) logistic regression

$$f_{\hat{w}, b}(\hat{x}) = \hat{w} \cdot \hat{x} + b \quad f_{\hat{w}, b}(\hat{x}) = \frac{1}{1 + e^{-(\hat{w} \cdot \hat{x} + b)}}$$

Same concepts:

- 1) Monitor gradient descent (learning curve)

- 2) Vectorized implementation

- 3) Feature scaling

The problem of overfitting

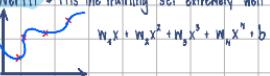
\rightarrow high bias
underfit \rightarrow does not fit the training set well



\rightarrow generalization \rightarrow fits training set pretty well



\rightarrow overfit \rightarrow fits the training set extremely well



Addressing Overfitting

How to fix overfit + 1) collect more training examples

- 2) select useful features

wrong or useful features could be lost

- 3) Regularization (reduce the size of parameters W_j)

overfit \rightarrow regularization

Cost Function with Regularization

note: W_j , b mostly small (≈ 0)

$M \times X \times M \times X \times M \times X \approx 0$

choose λ too small \rightarrow still overfit

$\frac{\partial J}{\partial w_j} = \frac{\partial}{\partial w_j} \left[\frac{1}{M} \sum_{i=1}^M (f_{\hat{w}, b}(\hat{x}^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^n W_j^2 \right]$

$\frac{\partial J}{\partial b} = \frac{\partial}{\partial b} \left[\frac{1}{M} \sum_{i=1}^M (f_{\hat{w}, b}(\hat{x}^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^n W_j^2 \right]$

choose λ too large \rightarrow prediction can't be straight line

\rightarrow $y = a$ ($a \neq b$)

if we have n features

$$J(\hat{w}, b) = \frac{1}{M} \sum_{i=1}^M (f_{\hat{w}, b}(\hat{x}^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^n W_j^2 + \frac{\lambda}{2} b^2$$

- 1) balance bias/grad
- 2) fit data
- 3) decrease W_j and regularization term

Regularized Linear Regression + Gradient Descent

from Gradient descent

$$W_j = W_j - \alpha \frac{1}{M} \sum_{i=1}^M (f_{\hat{w}, b}(\hat{x}^{(i)}) - y^{(i)}) X_j^{(i)} + \frac{\lambda}{M} W_j$$

don't have to regularize b

$$W_j = 1 W_j - \alpha \left(\frac{1}{M} \sum_{i=1}^M (f_{\hat{w}, b}(\hat{x}^{(i)}) - y^{(i)}) X_j^{(i)} \right) + \frac{\lambda}{M} W_j$$

$$W_j = (1 - \frac{\lambda}{M}) W_j$$

$$W_j = (1 - 0.001) W_j$$

$$W_j = 0.999 W_j$$

usual update

Regularized logistic regression

$$f_{\hat{w}, b}(\hat{x}) = \frac{1}{1 + e^{-\hat{w} \cdot \hat{x}}}$$

$$J_{\hat{w}, b} = -\frac{1}{M} \sum_{i=1}^M y^{(i)} \log(f_{\hat{w}, b}(\hat{x}^{(i)})) + (1 - y^{(i)}) \log(1 - f_{\hat{w}, b}(\hat{x}^{(i)})) + \frac{\lambda}{2M} \sum_{j=1}^n W_j^2$$

$n \rightarrow$ number of features
 $n \leq W_j$