

Linear Regression : close to local minimum

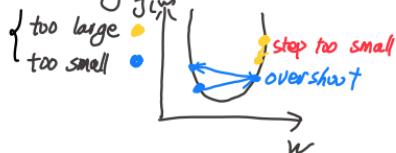
m : number of data.

$$\text{cost function: } J(w, b) = \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

$$\text{model: } f_{w,b}(x) = wx + b$$

$$\text{gradient descent: } \begin{cases} w = w - \alpha \frac{\partial}{\partial w} J(w, b) = \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x^{(i)} \\ b = b - \alpha \frac{\partial}{\partial b} J(w, b) = \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) \end{cases}$$

α is the learning rate



$$\text{with multiple feature: } f_{w,b}(\vec{x}) = \vec{w} \cdot \vec{x} + b = w_1 \cdot x_1 + w_2 \cdot x_2 + \dots + b$$

$$\text{Gradient descent (in multiple)} \quad \begin{cases} \vec{w} = \vec{w} - \alpha \frac{\partial}{\partial \vec{w}} J(\vec{w}, b), \quad \frac{\partial}{\partial w_i} J(w, b) = \frac{1}{m} \sum_{i=1}^m (f_{w,b}(\vec{x}^{(i)}) - y^{(i)}) x_i^{(i)} \\ b = b - \alpha \frac{\partial}{\partial b} J(w, b) \end{cases}$$

repeat } update each feature.

$$w_1 = w_1 - \alpha \frac{\partial}{\partial w_1} J(w, b)$$

$$\vdots$$

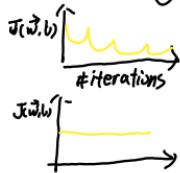
$$w_n = w_n - \alpha \frac{\partial}{\partial w_n} J(w, b)$$

Feature scaling: Scale the data \rightarrow

have some transform to the numbers

$$\text{mean normalization} \Rightarrow x_1 = \frac{x_1 - \mu_1}{\max - \min} \quad \text{Z-score normalization} \Rightarrow x_1 = \frac{x_1 - \mu_1}{\sigma_1}$$

Choose learning rate



Feature engineering. transform or combine parameter

$$f_{w,b}(\vec{x}) = w_1 x_1 + w_2 x_2 + \dots + b \rightarrow f_{w,b}(\vec{x}) = w_1 x_1 + w_2 x_2 + w_3 x_3 + b$$

$$X_3 = x_1 x_2$$

Polynomial regression use other function

$$f_{\vec{w}, b}(x) = w_1 x + w_2 x^2 + w_3 x^3 + b$$

change the function and do the gradient descent.

Classification: Logistic Regression.

(分类)

Sigmoid function
(logistic)

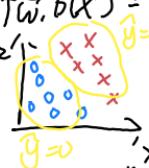


$$\boxed{z = \vec{w} \cdot \vec{x} + b} \rightarrow z \rightarrow \boxed{\begin{matrix} g(z) = \text{logistic}(z) \\ \text{scale to } (0,1) \end{matrix}} \Rightarrow f_{\vec{w}, b}(\vec{x}) = g(\vec{w} \cdot \vec{x} + b)$$

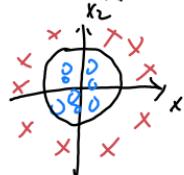
"logistic regression"

result: + \vec{w} , $b(\vec{x}) = 0.7$: 70% chance y is 1 (True)

Decision boundary: $f(\vec{w}, b)(\vec{x}_j) = P(x_j | \vec{w}, b)$ { when $f(\vec{x}_j | \vec{w}, b) > 0.5$
 $\hat{y} = 1$



Non-linear boundaries



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

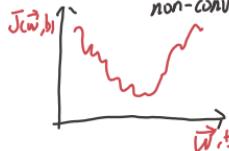
if $\geq x_1^2 + x_2^2 - 1 = 0$

$x_1^2 + x_2^2 = 1$ (boundary)

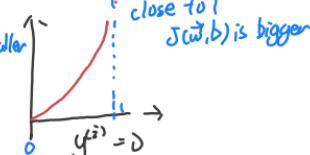
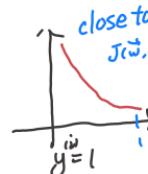
Cost function: $J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m \text{Loss}_i$

(or logistic regression)

non-convex



$$\text{Loss of logistic regression.}$$



Simplified Cost function

new Loss: $L(\vec{w}, b(x^{(i)}), y^{(i)}) = \underbrace{-y^{(i)} \log(\vec{w}, b(x^{(i)}))}_{y, close\ to\ 0, this\ close\ to\ 0} - \underbrace{(1-y^{(i)}) \log(1-\vec{w}, b(x^{(i)}))}_{y\ close\ to\ 1, this\ close\ to\ 1}$

Cost function: $J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m L(\vec{w}, b(x^{(i)}), y^{(i)})$

Use gradient descent (training)

derivative of $J(\vec{w}, b)$

$$\frac{\partial}{\partial w_j} J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m (\vec{w}, b(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

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

Process: repeat {

$$w_j = w_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (\vec{w}, b(x^{(i)}) - y^{(i)}) x_j^{(i)} \right]$$

$$b = b - \alpha \left[\frac{1}{m} \sum_{i=1}^m (\vec{w}, b(x^{(i)}) - y^{(i)}) \right]$$

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

Problem of Overfitting

use different model to fit the problem



when it overfitted



Solve problem:

(1) More data

(2) use fewer features. (feature selection)

(3) Regularization (L1/L2)

setting parameter smaller rather than 0
keep the feature, decrease the impact

Apply Regularization.

make parameter smaller. ($w_j \approx 0$)
(the penalize one)

given that parameters are w_1, w_2, \dots, w_i, b

$$J(\vec{w}, b) = \frac{1}{2m} \sum_{i=1}^m (\hat{t}_{\vec{w}, b}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2 \quad (\lambda > 0)$$

λ means Regularization parameter

when λ is bigger.
it has more impact
of w_j . it will shrink
it smaller

Sample 1: (Linear Regression)

$$\text{model : } \hat{t}_{\vec{w}, b}(\vec{x}) = w_1 x_1 + w_2 x_2^2 + w_3 x_3^3 + w_4 x_4^4 + b$$

$$\text{cost : } J(\vec{w}, b) = \frac{1}{2m} \sum_{i=1}^m (\hat{t}_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$

$$\text{Process: } w_j = w_j - \alpha \underbrace{\frac{\partial}{\partial w_j} J(\vec{w}, b)}_{\text{derivative}} \rightarrow \frac{1}{m} \sum_{i=1}^m (\hat{t}_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_j + \frac{\lambda}{m} w_j$$

$$b = b - \alpha \underbrace{\frac{\partial}{\partial b} J(\vec{w}, b)}_{\text{derivative}} \rightarrow \frac{1}{m} \sum_{i=1}^m (\hat{t}_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) \quad \begin{array}{l} \text{no necessary} \\ \text{for regularize} \\ b \end{array}$$

Sample 2 (Logistic Regression)

$$\text{model: } \hat{t}_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-c \vec{w} \cdot \vec{x} + b}}$$

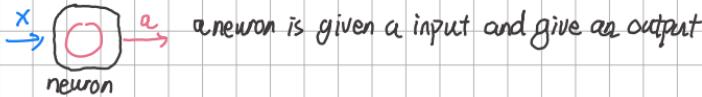
$$\text{cost : } J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m \underbrace{[y^{(i)} \log(\hat{t}_{\vec{w}, b}(\vec{x}^{(i)})) + (1 - y^{(i)}) \log(1 - \hat{t}_{\vec{w}, b}(\vec{x}^{(i)}))]}_{\text{loss of this Regression}} + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$

$$\text{Process: } w_j = w_j - \alpha \underbrace{\frac{\partial}{\partial w_j} J(\vec{w}, b)}_{\text{derivative}} \rightarrow \frac{1}{m} \sum_{i=1}^m (\hat{t}_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_j + \frac{\lambda}{m} w_j$$

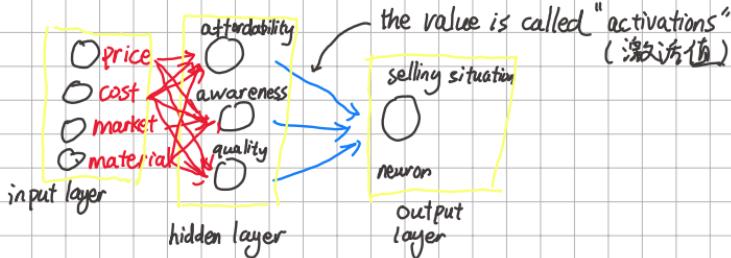
$$b = b - \alpha \underbrace{\frac{\partial}{\partial b} J(\vec{w}, b)}_{\text{derivative}} \rightarrow \frac{1}{m} \sum_{i=1}^m (\hat{t}_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) \quad \begin{array}{l} \text{no necessary} \\ \text{for regularize} \\ b \end{array}$$

$\hat{t}(\vec{w}, b)$ is Logistic function

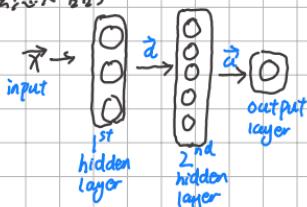
Neuron Network:



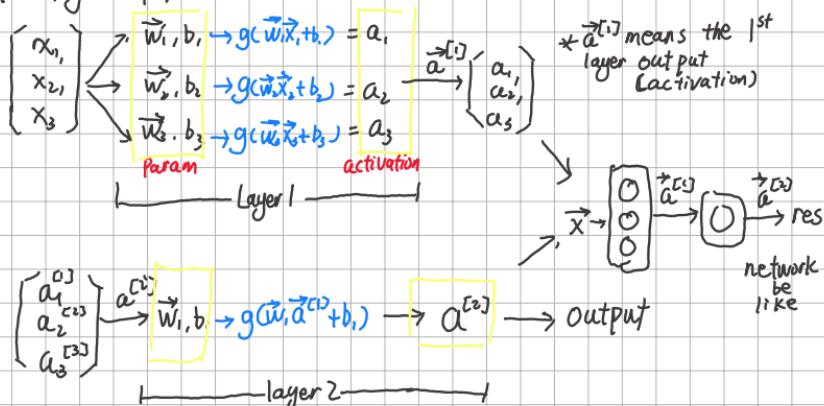
Demand Prediction:



Multiple perceptron: multiple layers
(多層感知器)



Build a layer of Network:



For the l th layer: $a_j^{[l]} = g(\vec{w}_j \cdot \vec{a}^{[l-1]} + b_j)$

forward propagation

parameter of layer l and unit j

Matrix multiply

轉置矩陣

$$\vec{a} = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \rightarrow \vec{a}^T \begin{bmatrix} 1 & 2 \end{bmatrix}$$

column
rows

$$1) \begin{bmatrix} 1 & 2 \end{bmatrix} \cdot \begin{bmatrix} 3 & 5 \\ 4 & 6 \end{bmatrix} = \begin{bmatrix} 1 \times 3 + 2 \times 4 & 1 \times 5 + 2 \times 6 \end{bmatrix}$$

$$[\leftarrow \vec{a}^T \rightarrow] \begin{bmatrix} \uparrow & \uparrow \\ \vec{w}_1 & \vec{w}_2 \\ \downarrow & \downarrow \end{bmatrix} = [\vec{a}^T \vec{w}_1 \quad \vec{a}^T \vec{w}_2]$$

$$2) \begin{bmatrix} 1 & 2 \\ -1 & -2 \end{bmatrix} \cdot \begin{bmatrix} 3 & 5 \\ 4 & 6 \end{bmatrix}$$

$$\begin{bmatrix} \leftarrow \vec{a}_1 \rightarrow & \uparrow & \uparrow \\ \leftarrow \vec{a}_2 \rightarrow & \vec{w}_1 & \vec{w}_2 \\ \downarrow & \downarrow & \downarrow \end{bmatrix}$$

$$= \begin{bmatrix} \vec{a}_1 \vec{w}_1 & \vec{a}_1 \vec{w}_2 \\ \vec{a}_2 \vec{w}_1 & \vec{a}_2 \vec{w}_2 \end{bmatrix}$$

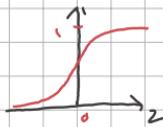
Training: ① define model $f_{\vec{w}, b}(\vec{x})$

$$② \text{loss \& cost: } J(\vec{w}, b) = \underbrace{\frac{1}{m} \sum_{i=1}^m l(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)})}_{\text{cost}} \quad \text{Loss function}$$

③ train data to minimize $J(\vec{w}, b)$

Other Activation function: $a_+^{(i)} = g(\vec{w} \cdot \vec{x} + b)$

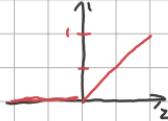
① Sigmoid



$$g(z) = \frac{1}{1+e^{-z}} = p(y=1|\vec{x})$$

good to solve
binary problem

② ReLU



$$g(z) = \max(0, z) \quad \text{on } 0, \text{ or } \pm$$

③ Softmax (multi-classification) sample: $y=1, 2, 3, 4$

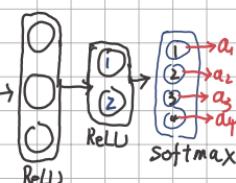
$$z_1 = \vec{w}_1 \cdot \vec{x} + b_1 \Rightarrow \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}} = a_1 = p(y=1|\vec{x})$$

$a_1 + a_2 + a_3 + a_4 = 1$

$$z_2 = \vec{w}_2 \cdot \vec{x} + b_2 \Rightarrow \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}} = a_2 = p(y=2|\vec{x})$$

$$z_3 = \vec{w}_3 \cdot \vec{x} + b_3 \Rightarrow \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}} = a_3 = p(y=3|\vec{x})$$

$$z_4 = \vec{w}_4 \cdot \vec{x} + b_4 \Rightarrow \frac{e^{z_4}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}} = a_4 = p(y=4|\vec{x})$$



↓ to sum up

$$\left\{ \begin{array}{l} z_j = \vec{w}_j \cdot \vec{x} + b_j, j=1, \dots, N \\ a_j = \frac{e^{z_j}}{\sum_{k=1}^N e^{z_k}} = P(y=j | \vec{x}) \end{array} \right.$$

Cost of softmax Regression:

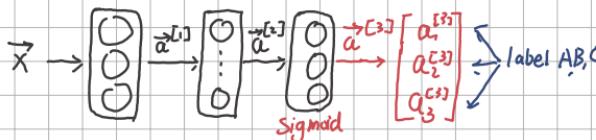
Loss: $\text{loss}(a_1 + \dots + a_N, y) = \begin{cases} -\log a_1 & \text{if } y=1 \\ -\log a_2 & \text{if } y=2 \\ \vdots & \vdots \\ -\log a_N & \text{if } y=N \end{cases}$

$$\text{Cost} : J(\vec{w}) = \frac{1}{m} \sum_{n=1}^m \text{loss}_{(n)}$$



Multiple Classification

1) multi-label



2) multi-class



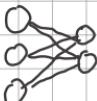
when using multiClass classification,
you have to build more network.

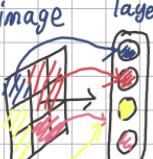


Other learning Algorithm:

Adam automatically adjust learning rate α

Other Layer Type

Dense:  each get all neuron's activation

Convolutional Layer: 

each activation focus of different part
benefit: (1) less overshooting
(2) need less data.
(3) Faster computation.

Evaluate model

1) Dataset → training set 70% → $J_{\text{train}}(\vec{w}, b)$
 → test set 30% → $J_{\text{test}}(\vec{w}, b)$

2) Dataset → Training set 60%
 → cross-validation set 20% → test different model
 → test set 20%

Model selection

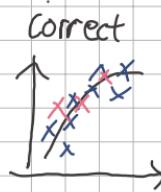
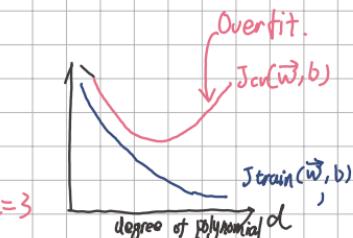
$$1. \vec{w}, b(\vec{x}) = w_1 x_1 + b \quad d=1$$

$$2. \vec{w}, b(\vec{x}) = w_1 x_1 + w_2 x^2 + b \quad d=2$$

$$\vdots \quad 10. \vec{w}, b(\vec{x}) = w_1 x_1 + w_2 x^2 + \dots + w_{10} x^{10} + b \quad d=3$$

Cross-validation: $J_{cv}(\vec{w}^{(4)}, b^{(4)})$ use to choose which model

test set use to evaluate the final model



$J_{\text{train}}(\vec{w}, b)$ is low
 $J_{cv}(\vec{w}, b)$ is low

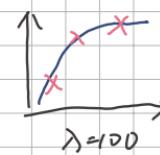
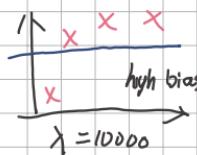


$J_{\text{train}}(\vec{w}, b)$ is low
 $J_{cv}(\vec{w}, b)$ is high

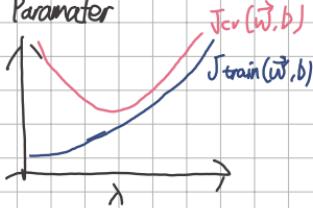
Choose Regularization · model: $\vec{w}, b(x) = w_1 x_1 + w_2 x^2 + w_3 x^3 + \dots + b$

$$J(\vec{w}, b) = \frac{1}{2m} \sum_{i=1}^m (\vec{w}, b(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n w_j^2$$

Regularization



For Regularization Parameter

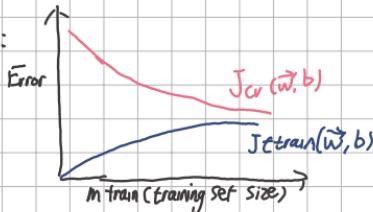


Judge the performance

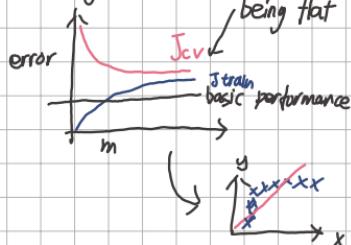
- 1) Baseline Performance (高) if $① < ② \Rightarrow$ high bias
- 2) Training Error if $① < ③ \Rightarrow$ overfit
- 3) Cross Validation Error

Learning Curve

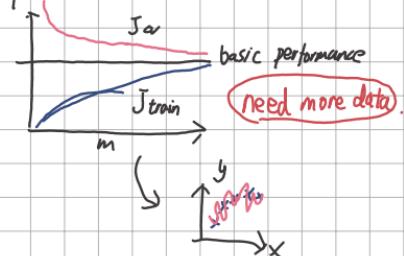
For a model:



For high bias model



For overfit model



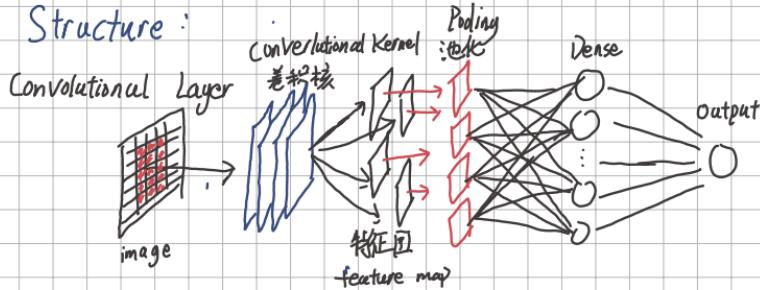
Debug a model

- 1) get more data \rightarrow high variance
- 2) smaller features \rightarrow high variance
- 3) additional features \rightarrow high bias

- 4) add polynomial features (x_1^2, x_1x_2, \dots , etc.)
↳ high bias
- 5) decrease $\lambda \rightarrow$ high bias
- 6) increase $\lambda \rightarrow$ high variance

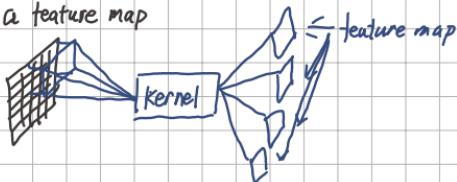
Convolutional Neuron Network 卷积神经网络.

Structure :

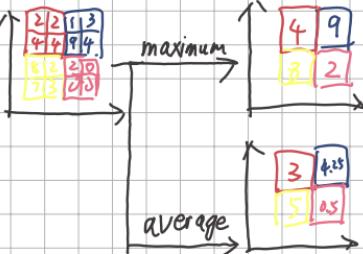


Convolutional Kernel:

match each area of the input, compute each part into a feature map



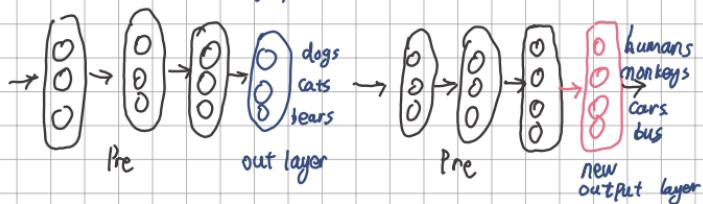
Pooling:



make features easier to calculate

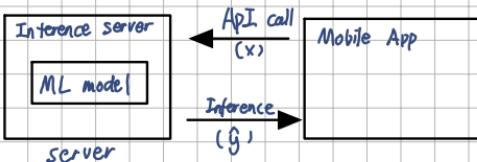
Transfer learning (迁移学习)

use pre-trained model to train another model
(just change the output layer)



Change layer, and train the new model (further train)
 { train new part
 { train whole model

Deploy a model



Precision 精度 (P) $\frac{\text{True Pos}}{\text{True Pos} + \text{False Pos}}$ (True among pos)

Recall 召回 (R) $\frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}}$ (True among True Positive)

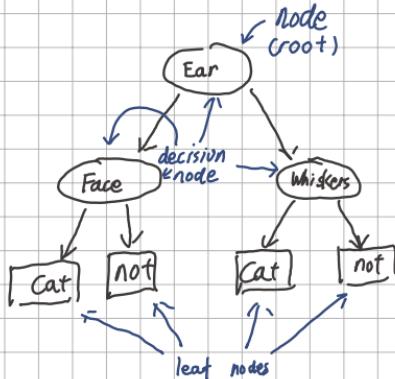
$$F1 \text{ Score} : \frac{1}{\frac{1}{P} + \frac{1}{R}} = 2 \cdot \frac{RP}{P+R}$$

Decision tree (决策树)

1. How to split?

2. How to stop split?

- node is 100% pure
- exceeding maximum depth
- purity score.

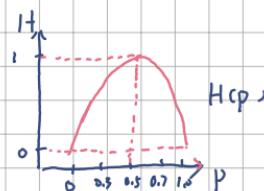


Entropy (H): measure of impurity

p_i : a class of the set

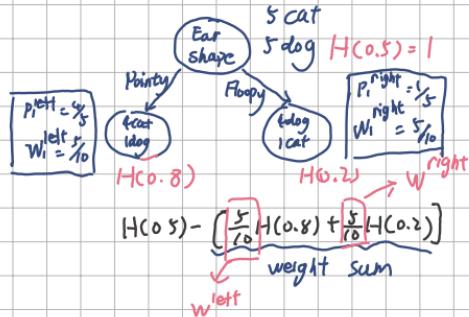
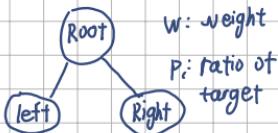
$$P_0 = 1 - P_1$$

$$H(p_i) = -P_1 \log_2(P_1) - P_0 \log_2(P_0)$$



Information gain 信息增益

$$\text{define: } H_C(p_i)^{\text{root}} = (w_{\text{left}} H(p_i^{\text{left}}) + w_{\text{right}} H(p_i^{\text{right}}))$$



Whole Process:



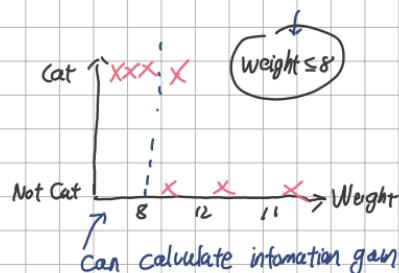
One-hot encoding (Process multiple features)

in k features, 0 means not, 1 means has the feature

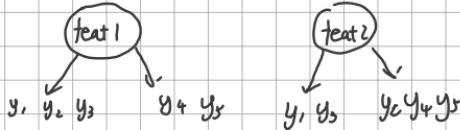
	Pointy ear	Floppy ear	Oval ear
1	1	0	0
2	0	1	0
3	1	0	0
4	0	0	1

Continuous values features

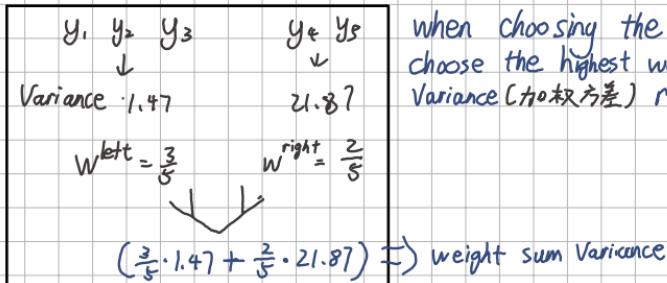
split at a threshold which can get highest information gain



Regression tree

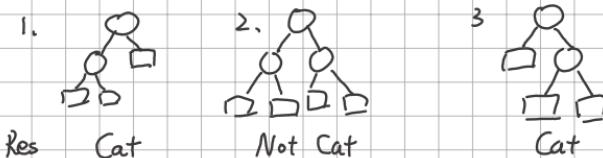


count Variance as purity and get information gain
(方差)



Tree ensemble.

when you have multiple trees, use both tree to judge and choose the best result



let the trees vote the result, and choose the one who has highest vote.

How to come up with those trees?

Sampling with replacement (带放回抽样)



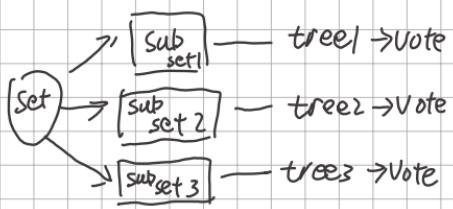
each time get one and put it back, finally will get different random training set. to use that way you can build a new training set.

1	red	red	blue
2	green	blue	yellow
3	yellow	green	red

Random forest:

given a dataset of size m .

loop } For $b=1$ to $B(100, 64, 128 \dots)$
 use sampling with replacement. get a new training set at size of m and get a new tree
 when predicting, each tree will give a vote



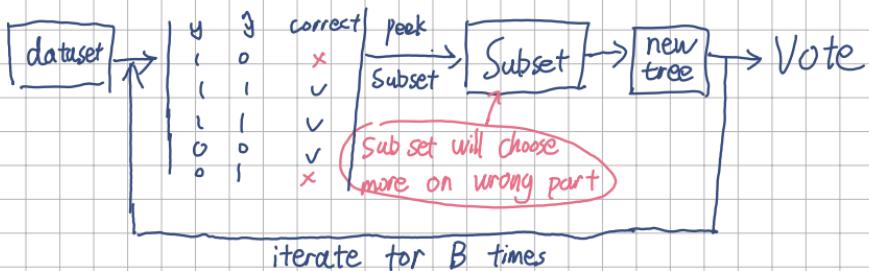
Randomizing the feature choice

from n features, only choose a subset of $K < n$ features to allow algorithm to choose.

usually: $K = \sqrt{n}$

XG Boost:

For the loop in Random forest



XG Boost (Extreme Gradient Boosting) provide a way to choose subset and fit the dataset

When to use decision trees:

tree:

- structured data (data table)
- fast training

Neuron Networks

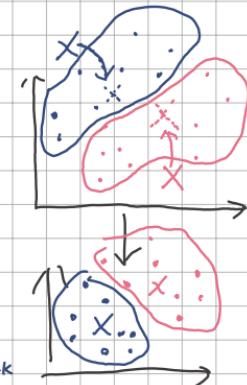
- all types data (images, audio)
- mixed data.
- slower
- work with transfer learning

Unsupervised Learning

Clustering (聚类)

k means:

- 1) choose random cluster centroid (簇质心)
- 2) Assign points depends on the distance to the cluster centroid
- 3) move the cluster centroid to the average location of the point of itself
- 4) loop



in program:

Choose initialize K cluster centroids $\mu_1, \mu_2, \dots, \mu_K$

Repeat {

for $i=1$ to m

$c^{(i)} = \text{Index (from 1 to } K\text{) of cluster centroid closest to } x^{(i)}$
 In math: $\min_k \|x^{(i)} - \mu_k\|_2^2$
 $L_2 \text{ norm (L2范数)}$

for $k=1$ to K

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

$$\mu_k = \frac{1}{4} [x^{(1)} + x^{(2)} + x^{(3)} + x^{(4)}]$$

Cost function of k-means:

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

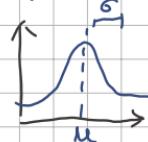
Distortion (失真)

Randomly choose 100 different initialization and choose the lowest J value

Anomaly Detection (异常检测): find unnormal parts

Gaussian Distribution:

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



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

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

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

Process:

1. Choose n features x_i as anomalous examples
2. Fit parameters: $\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2$

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

3. give a new example x , compute $p(x)$:

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

- 4: Anomaly if $p(x) < \epsilon$

Algorithm evaluate

If . 10000 good sample
20 anomalous sample.

Training set: 6000 good samples, $C_y = 0$

Cross Validation: 4000 good samples ($y=0$)
20 anomalous ($y=1$)

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

evaluation metrics.

- True pos, false pos, false neg, true neg
- Precision/Recall
- F_1 -Score

Recommend System (Predict rating in unlabeled)

For users:	Alice(1)	Bob(2)	Carol(3)	x_1 (romance)	x_2 (action)
Movie1	5	3	4	0.9	0
Movie2	0	?	?	1.0	0.01
Movie3	0	0	?	0.1	1.0

Notation:

$r(c_i, j)$ if user j has rated movie i (0 or 1)
 $y^{(i,j)}$ = rating given by user j on movie i

i : movie
 j : user

$w^{(i)}, b^{(i)}$ parameters for user j

$x^{(i)} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$ feature vector for movie i

Predict rating: $w^{(i)} \cdot x^{(i)} + b^{(i)}$

$m^{(j)}$ = number of movies rated by user j

Cost:

$$\min_{w^{(i)}, b^{(i)}} J(w^{(i)}, b^{(i)}) = \frac{1}{2m^{(j)}} \sum_{i: r(c_i, j)=1} (w^{(i)} \cdot x^{(i)} + b^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2m^{(j)}} \sum_{k=1}^n (w_k^{(i)})^2$$

Total Cost: (fitting the user parameter)

$$(1) J(w^{(1)}, \dots, w^{(n)}) = \frac{1}{2} \sum_{j=1}^{n_u} \sum_{i: r(c_i, j)=1} (w^{(i)} \cdot x^{(i)} + b^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (w_k^{(i)})^2$$

Feature Cost (get the proper feature $x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$ of movies)

$$(2) J(x^{(1)}, x^{(2)}, \dots, x^{(n)}) = \frac{1}{2} \sum_{i=1}^{n_m} \sum_{j: r(c_i, j)=1} (w^{(i)} \cdot x^{(i)} + b^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2$$

Collaborative filtering: (协同过滤)

Put together:

$$\min_{\substack{w^{(1)}, \dots, w^{(n_m)} \\ b^{(1)}, \dots, b^{(n_u)} \\ x^{(1)}, \dots, x^{(n)}}} J(w, b, x) = \frac{1}{2} \sum_{(i,j): r(c_i, j)=1} (w^{(i)} \cdot x^{(i)} + b^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (w_k^{(i)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2$$

this formula made w, b, x both as parameters while learning from chart. It's available to count

$$+ \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2$$

$\rightarrow i$ Cost from the movie that has rated

	Alice	Bob	Carol	
1	5	3	?	rating
2	?	2	3	

Train the model

Linear regression

Repeat {

$$w_i^{(j)} = w_i^{(j)} - \alpha \frac{\partial}{\partial w_i^{(j)}} J(w, b, x)$$

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

$$x_k^{(j)} = x_k^{(j)} - \alpha \frac{\partial}{\partial x_k^{(j)}} J(w, b, x)$$

}

while calculating, w, b, x are real numbers, and $J(w, b, x)$ comes from the dataset (people's rate) to fit.

For binary label

like: did the user clicked? (%)

did it attract user? (%)

Model:

$$y^{(i,j)}: f_{(w,b,x)}(x) = g(w^{(i,j)} \cdot x^{(i,j)} + b^{(j)})$$

Loss:

$$L(f_{(w,b,x)}(x), y^{(i,j)}) = -y^{(i,j)} \log(f_{(w,b,x)}(x)) - (1-y^{(i,j)}) \log(1-f_{(w,b,x)}(x))$$

二元交叉熵

Loss function:

$$J(w, b, x) = \sum_{(i,j): r(i,j)=1} L(f_{(w,b,x)}(x), y^{(i,j)})$$

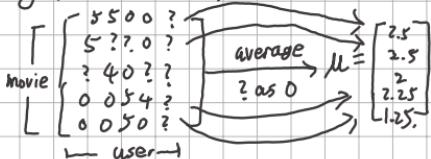
Add new user

	Alice	Bob	Eve
1	5	?	?
2	0	0	?
3	?	3	?

(均值) $\bar{x} = \frac{1}{n} \sum x_i$

use Mean normalization:

get the data from chart



$$\Rightarrow \text{Count: } x - \mu \quad \begin{bmatrix} 2.5 & 2.5 & 2.5 & 2.5 \\ 2.5 & ? & ? & 2.5 \\ 2.25 & 2 & -2 & ? \\ 2.25 & -2.25 & 2.75 & 1.25 \\ 1.25 & -1.25 & 2.75 & 1.25 \end{bmatrix}$$

For user j on movie i , we predict: $w_i^{(j)} \cdot x_i^{(j)} + b^{(j)} + M_i$

Find Related items

$$\text{feature } \vec{x}^{(i)} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$\text{Similarity } \parallel \vec{x}^{(k)} - \vec{x}^{(i)} \parallel^2$$

Content-Based filtering

$r(i, j)$ user j rated item i

$y(i, j)$ rating given by user j on item i (if defined)

User features

- Age
- Gender
- Country
-

Movie feature

- Year
- Genre
- Reviews
-

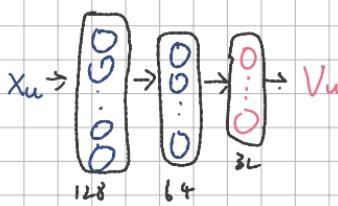
$$\left. \begin{array}{l} \text{User features} \\ \text{Movie feature} \end{array} \right\} \quad \left. \begin{array}{l} \vec{x}_u^{(j)} \text{ for user } j \\ \vec{x}_m^{(i)} \text{ for movie } i \end{array} \right\} = \left[\begin{array}{c} 0.3 \\ 0.42 \\ \vdots \\ 2.7 \end{array} \right] \Rightarrow \vec{v}_u^{(j)} \quad \left[\begin{array}{c} 0.2 \\ 0.51 \\ \vdots \\ 7.3 \end{array} \right] \Rightarrow \vec{v}_m^{(i)}$$

no need to
be same dimension

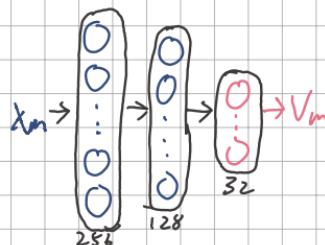
Some dimension
↓
count $v_u^{(j)}, v_m^{(i)}$
as predicted
ratings

How to get $\vec{v}_u^{(j)}, \vec{v}_m^{(i)}$?

$x_u \rightarrow v_u$ User Network



$x_m \rightarrow v_m$ Movie network



Prediction: $\vec{v}_u^{(j)} \cdot \vec{v}_m^{(i)}$

$g(\vec{v}_u^{(j)} \cdot \vec{v}_m^{(i)})$ to predict the probability of $y^{(i,j)}$ is 1

Cost function. $J = \sum_{(i,j) \in r(u, m)} (V_u^{(j)} \cdot V_m^{(i)} - y^{(i,j)})^2 + NN \text{ regularization term.}$

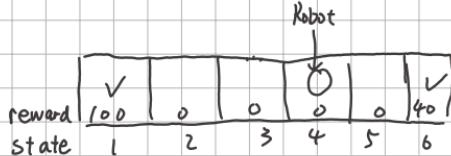
Reinforcement Learning

强化学习

Return

Assuming a robot reach a goal can have reward and each move have a Discount

$$\text{the Reward: } \text{go left} \cdot 0 + \underbrace{(0.9) \cdot 0 + (0.9)^2 \cdot 0 + (0.9)^3 \cdot 100}_{\text{Discount factor}} = 72.9$$



The discount Factor represents the discount of each move
The Formula:

$$\text{Return} = R_1 + rR_2 + r^2R_3 + \dots \quad (\text{until terminal state})$$

R: reward r: discount factor

Policy (策略)

The goal of Reinforcement Learning is to find a policy that can make states reflect into actions.

$$\text{state } s \xrightarrow{\pi} \text{action } a$$

State action value function

$Q(s, a)$

- start from state s
- take action a (once)
- then behave optimally after that
best behave

$Q(s, a)$ means the reward you get when you do a action and do the best behave after you do that

- 1) The best possible return from state s is $\max_a Q(s, a)$
- 2) The best possible action in state s is action a that gives $\max_a Q(s, a)$

Bellman equation

R_{CS} = Reward of current state

S : current state

a : current action

s' : state you get to after taking action a

a' : action that you take in state s .

formula:

$$Q(S, a) = R_S + \gamma \max_{a'} Q(S', a')$$

$Q(S, a)$

Reward you
get right
away.

Return from behave
optimally starting from s'

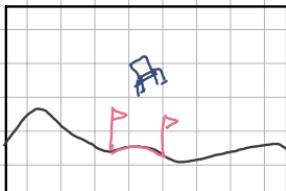
$$\text{example: } R_1 + \gamma R_2 + \gamma^2 R_3 + \gamma^3 R_4 + \dots$$

$$= R_1 + \gamma [R_2 + \gamma R_3 + \gamma^2 R_4 + \dots]$$

$\underbrace{\qquad\qquad\qquad}_{Q(S', a')}$

It's a recursion function

MOON Lander



moon lander, can fire engine (push up) or fire side thruster (push angle), algorithm
need to control the lander to the pad

State:

$$S = \begin{bmatrix} x \\ y \\ \dot{x} \\ \dot{y} \\ \theta \\ l \\ r \end{bmatrix} \quad \begin{array}{l} x \text{ axis} \\ y \text{ axis} \\ x \text{ speed (change rate)} \\ y \text{ speed} \\ \text{angle} \\ \text{left leg landed;} \\ \text{right leg landed;} \end{array}$$

Reward function

- Getting to land pad: 100 - 1/e0
- Additional reward for moving toward the pad
- Crashed: -100
- Soft landing: +100
- Leg grounded: +10
- Fire main engine: -0.3
- Fire side thruster: -0.03

Task: learn a policy π that given S picks action $a = \pi(S)$

so as to maximize the return ($R = 0.985$)

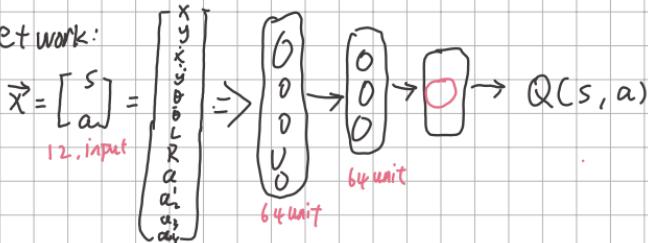
Pick state function

In state s , use NN to compute

$Q(s, \text{main})$, $Q(s, \text{left})$, $Q(s, \text{right})$, $Q(s, \text{nothing})$

And Ultimately maximize $Q(s, a)$ (one-hot)

Network:



Learn of Model

$$\text{Bellman Equation: } Q(s, a) = R(s) + \gamma \max_a Q(s', a')$$

Dataset:

$(s^{(1)}, a^{(1)}, R(s^{(1)}), s'^{(1)})$	$x^{(1)}$	$y^{(1)}$
$(s^{(2)}, a^{(2)}, R(s^{(2)}), s'^{(2)})$	$x^{(2)}$	$y^{(2)}$
$(s^{(3)}, a^{(3)}, R(s^{(3)}), s'^{(3)})$	$x^{(3)}$	$y^{(3)}$

$$y^{(i)} = R(s^{(i)}) + \gamma \max_{a'} Q(s^{(i)}, a') \quad \text{comes from learning}$$

Process:

Initialize neural Network randomly as guess of $Q(s, a)$
Repeat

Get $(s, a, R(s), s')$

Store recent data tuples \leftarrow Replay Buffer

Train:

Create set of example of $x = (s, a)$, $y = R(s) + \gamma \max_{a'} Q(s', a')$ training

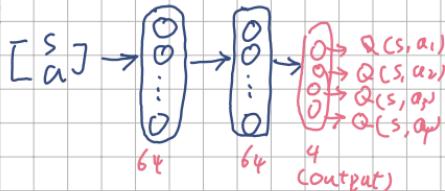
Train Q_{new} such that $Q_{\text{new}}(s, a) \approx y$

Set $Q = Q_{\text{new}}$

↓

Improve Net work

prevent from calculate Q too much times. Network can have multiple value



Choose $Q(s, a_i)$

option 1:

Pick action a that maximize $Q(s, a)$

Option 2 (ϵ -greedy)

With 0.95 prob. pick the max $Q(s, a)$ "Exploitation"

With 0.05 prob. pick a randomly "Exploration"

With ϵ -greedy policy, algorithm can discover different way of finding $\max Q(s, a)$ better (more discovery)

Start ϵ high
 $1.0 \rightarrow 0.01$
Gradually decrease

End.