

$$L = -y \log \alpha - (1-y) \log (1-\alpha)$$

$$\begin{aligned}\frac{dL}{d\alpha} &= -\frac{1}{\alpha}(1) - (1-y)\frac{-1}{1-\alpha} \\ &= \frac{-y}{\alpha} + \frac{1-y}{1-\alpha} = \frac{-y(1-\alpha) + \alpha(1-y)}{\alpha(1-\alpha)} \\ &= \frac{\cancel{y}-y-\cancel{\alpha}+\alpha}{\alpha(1-\alpha)} \\ &= \frac{\alpha-y}{\alpha(1-\alpha)}\end{aligned}$$

$$\begin{aligned}\frac{d\alpha}{dz} &= \frac{d}{dz} \pi(z) \\ &= \pi(z)(1-\pi(z))\end{aligned}$$

$$\begin{aligned}\frac{dL}{dz} &= \frac{dL}{d\alpha} \frac{d\alpha}{dz} \\ &= \frac{(a-y) \cdot \cancel{\alpha(1-\alpha)}}{\cancel{\alpha(1-\alpha)}} \\ &= a-y\end{aligned}$$

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

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

$$\begin{aligned}\frac{d}{dw} J(w, b) &= \frac{d}{dw} \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2 \\ &= \frac{d}{dw} \frac{1}{2m} \sum_{i=1}^m (w x^{(i)} + b - y^{(i)})^2\end{aligned}$$

$$\begin{aligned}&= \frac{1}{2m} \sum_{i=1}^m (w x^{(i)} + b - y^{(i)}) \cancel{2x^{(i)}} \\ &= \frac{1}{m} \sum_{i=1}^m (w x^{(i)} + b - y^{(i)}) x^{(i)}\end{aligned}$$

$$\begin{aligned}
\frac{\partial}{\partial b} J(w, b) &= \frac{1}{2b} \frac{1}{2m} \sum_{t=1}^m (f_{w,b}(x^{(t)}) - y^{(t)})^2 \\
&= \frac{1}{2b} \frac{1}{2m} \sum_{t=1}^m (w^T x^{(t)} + b - y^{(t)})^2 \\
&= \frac{1}{2m} \sum_{t=1}^m (w^T x^{(t)} + b - y^{(t)})^2 \\
&= \frac{1}{m} \sum_{t=1}^m (w^T x^{(t)} + b - y^{(t)}) \\
&= \frac{1}{m} \sum_{t=1}^m (f_{w,b}(x^{(t)}) - y^{(t)})
\end{aligned}$$

$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(w_1, w_2, \dots, w_n, b)$$

$$b = b - \alpha \frac{\partial}{\partial b} J(w_1, w_2, \dots, w_n, b)$$

$$\rightarrow w_j = w_j - \alpha \frac{d}{dw_j} J(\vec{w}, b)$$

$$b = b - \alpha \frac{d}{db} J(\vec{w}, b)$$

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

$$\rightarrow w - \alpha \left(\frac{1}{2m} \sum_{t=1}^m (f_{w,b}(x^{(t)}) - y^{(t)})^2 \right)$$

$$\rightarrow w - \alpha \frac{1}{2m} \sum_{t=1}^m (w^T x^{(t)} - y^{(t)}) R(x^{(t)})$$

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

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

$$= b - \alpha \frac{1}{2m} \sum_{t=1}^m f_{w,b}(x^{(t)}) - y^{(t)}$$

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

$$= w - \alpha \left(\frac{1}{2m} \sum_{t=1}^m (f_{w,b}(x^{(t)}) - y^{(t)})^2 \right)$$

.

$$\begin{aligned}
 &= N - \alpha \frac{1}{m} \sum_{i=1}^m (w x^{(i)} - y^{(i)}) \\
 &\leftarrow w - \alpha \frac{1}{m} \times L \times \sum_{i=1}^m (w x^{(i)} - y^{(i)}) x^{(i)} \\
 b &= b - \alpha \frac{1}{m} \sum_{i=1}^m (w x^{(i)} - y^{(i)})
 \end{aligned}$$

② Mean Normalization

$$\frac{x - \mu}{\max - \min}$$

③ Z score Normalization (with standard deviation)

$$\frac{x - \mu}{\sigma}$$

④ Feature Scaling

$$\begin{aligned}
 -1 \leq x_i \leq 1 \\
 -3 \leq x_i \leq 3 \\
 -0.3 \leq x_i \leq 0.3
 \end{aligned}
 \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} \text{all acceptable}$$

if different values of features ranging too large.
we have to do feature scaling

$0 \leq x_1 \leq 3$	okay, no rescaling	} to prevent local minimum for gradient descent
$-2 \leq x_2 \leq 0.5$	okay, no rescaling	
$-100 \leq x_3 \leq 100$	too large \rightarrow rescale	
$-0.001 \leq x_4 \leq 0.001$	too small \rightarrow rescale	
$98.6 \leq x_5 \leq 105$	too large \rightarrow rescale	

⑤ Gradient Descent

$$\begin{aligned}
 w_j &= w_j - \eta \frac{\partial L}{\partial w_j}, L = J(\vec{w}, b) \\
 b &= b - \eta \frac{\partial L}{\partial b}
 \end{aligned}$$

10b

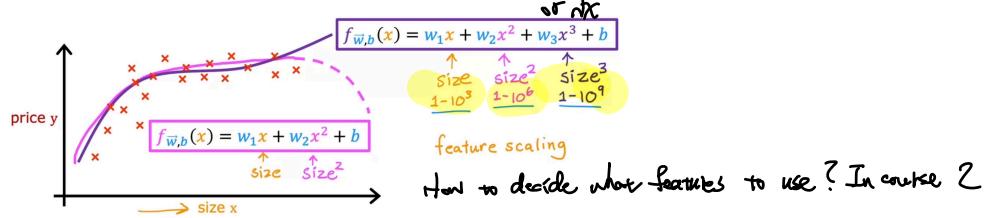
① Feature Engineering

- Utilizing known features to create new features
By combining or transforming

② Polynomial Regression

→ Multiple Linear Regression + Feature Engineering
goal: fit curves

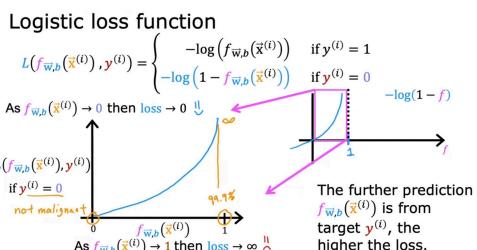
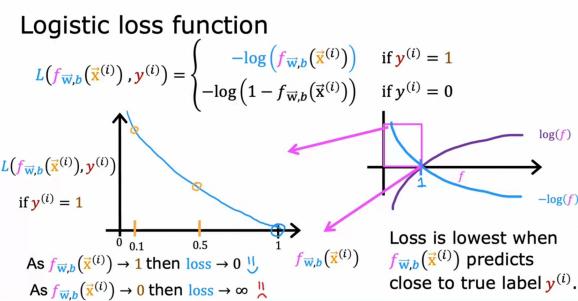
→ feature scaling is needed



③ Logistic loss function

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

$\in [0, 1]$



Cst

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

$$\begin{aligned} &= \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(x^{(i)}) \cdot y^{(i)}) \\ &\text{mean all the loss to obtain the overall cost} \\ &= \begin{cases} -\log(f_{\vec{w}, b}(x^{(i)})) & \text{if } y^{(i)} = 1, \text{ plot of the cost func is "convex", so GD is applied to find optimal} \\ -\log(1 - f_{\vec{w}, b}(x^{(i)})) & \text{if } y^{(i)} = 0 \end{cases} \end{aligned}$$

* if using the MSE for logistic regression, the loss is non-convex, which cause local minima

$$\rightarrow \text{loss}(f_{\vec{w}, b}(x^{(i)}), y^{(i)}) = (-y^{(i)} \log(f_{\vec{w}, b}(x^{(i)})) - (1-y^{(i)}) \log(1 - f_{\vec{w}, b}(x^{(i)})))$$

② Address Overfitting

1. More data collection to train model
2. Feature selection (throwing some features)
but might lose useful features
3. Regularization on w_1, w_2, \dots, w_n (Reduce the sizes of params)

③ Regularization on linear regression

1. Penalize all the w_j params (because we don't know which should be included/excluded)

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

$\lambda > 0$
regularization parameter

→ minimize the original cost

$$J(\vec{w}, b) = \min \left[\frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(x^{(i)}) - y^{(i)}) + \frac{\lambda}{2m} \sum_{j=1}^m w_j^2 \right]$$

[
 bottom: minimization to run
 middle: cause error
 top: linear regression to run]

\leftarrow more steps

regularization with

keep W_j small

if λ is very big, to reduce original error $\rightarrow W_j$ would be small

Gradient Descent

$$W_j = W_j - \alpha \frac{\partial}{\partial W_j} J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} W_j^{(i)}$$

$$b = b - \alpha \frac{\partial}{\partial b} J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) \quad \text{don't have to regularize "b"}$$

for GD:

$$\begin{aligned} \text{loop } \left\{ \begin{array}{l} W_j = W_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} W_j^{(i)} \right] \\ b = b - \alpha \left(\frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) \right) \end{array} \right\} \end{aligned}$$

$$W_j = \underbrace{(W_j - \alpha \frac{\lambda}{m} W_j)}_{W_j(1 - \alpha \frac{\lambda}{m})} - \underbrace{\alpha \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x_j^{(i)}}_{\text{usual update}}$$

slightly less than 1

$$0.1 \frac{1}{50} = 0.002 \rightarrow 1 - \alpha \frac{\lambda}{m} = 1 - 0.002 = 0.998$$

② Regularization on logistic regression

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

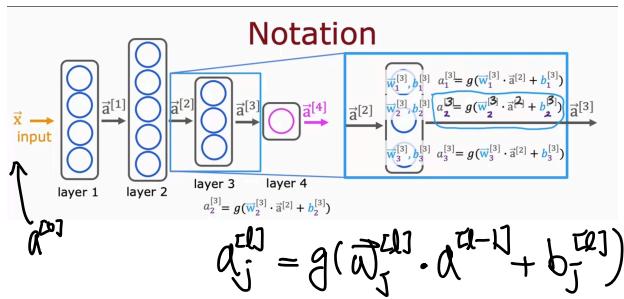
$$J(\vec{w}, b) = -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \log(f_{w,b}(x^{(i)})) + (1-y^{(i)}) \log(1-f_{w,b}(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n W_j^2$$

Gradient Descent

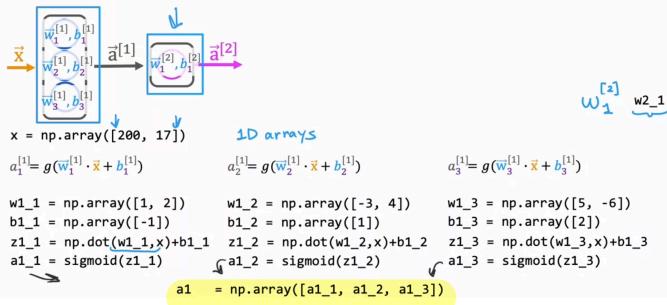
$$W_j = W_j - \alpha \frac{\partial}{\partial W_j} J(\vec{w}, b) \rightarrow = \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} W_j^{(i)}$$

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

• X C₂



forward prop (coffee roasting model)



• X Numerical Roundoff Errors

Ex: Logistic Regression

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

Logistic loss

$$\text{loss} = -y \log(a) - (1-y) \log(1-a)$$

More accurately:

$$\text{loss} = -y \log\left(\frac{1}{1+e^{-z}}\right) - (1-y) \log\left(-\frac{1}{1+e^{-z}}\right)$$

logits

```

More numerically accurate implementation of softmax
Softmax regression
    (a1,...,a10) = g(z1,...,z10)
    Loss = L(d,y) = { -log(a1) if y = 1
                      { -log(a1) if y = 10
                        { -log(a1) + ... + log(ai) if y = i
More Accurate
    L(d,y) = { -log(a1) + ... + log(ai) if y = 1
                { -log(a1) + ... + log(ai) if y = 10
                { -log(a1) + ... + log(ai) if y = i
model.compile(loss='SparseCategoricalCrossentropy', from_logits=True)

```

Andrew Ng
although unfortunately, it is a

○ Multi-labeled Classification

1. Treat as separate classification problems

2. A NN to do multiple classification

$$\begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix} \begin{array}{l} \text{car} \\ \text{bus} \\ \text{pedestrian} \end{array}$$

⑤ Regularization and bias/variance

Model: $f_{\vec{w}, b}(x) = w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4 + b$

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

as λ bigger, model would try
to keep w_j smaller.

but if $\lambda=0$, there is no regularization \rightarrow overfitting

J_{train} is small, but

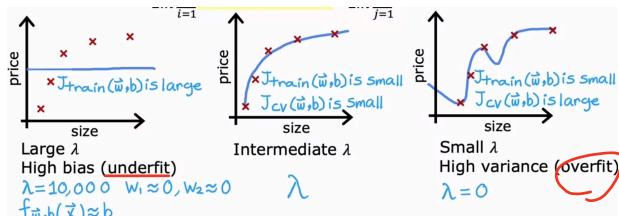
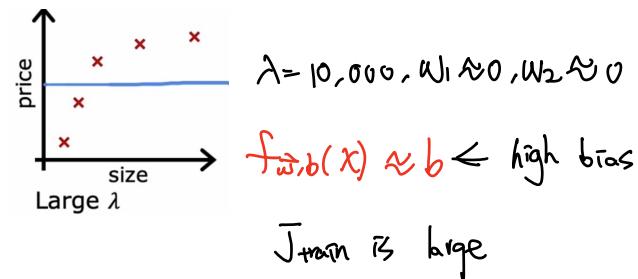
J_{cv} is large

Overfitting: J_{train} is small

J_{cv} is large

Underfitting: J_{train} is large,

J_{cv} is more larger

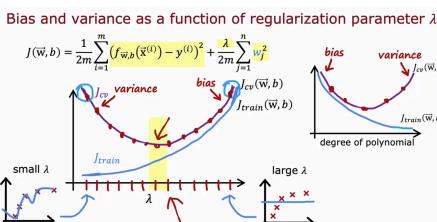


choose the λ value by examining the J_{cv}

* if we have λ with multiple values \rightarrow we train model with multiple \vec{w}, b pairs
then we choose the pair with lowest J_{cv} . apply the pair \vec{w}, b to test error

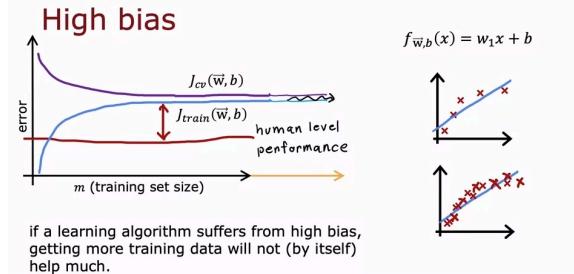
$\lambda \uparrow, J_{\text{cv}} \uparrow \rightarrow J_{\text{cv}}(̄w, b) \approx b$ (high bias)

$\lambda \downarrow, J_{\text{cv}} \downarrow \rightarrow$ overfitted



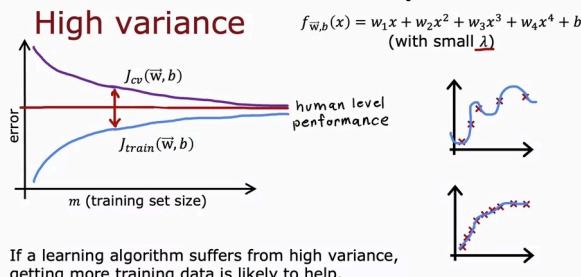
• ✘- high bias / high variance

1. high bias → underfitting (该 algo 需要改进, 虽然 data 量不会改善 bias)



2. high variance → overfitting

lower $\lambda \rightarrow$ higher w , far away from high bias situation



- ✘ Debugging the learning algorithm

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

Debugging a learning algorithm

You've implemented regularized linear regression on housing prices

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

But it makes unacceptably large errors in predictions. What do you try next?

- Get more training examples
- Try smaller sets of features x, x^2, x^3, \dots
- Try getting additional features
- Try adding polynomial features $(x_1^2, x_2^2, x_1 x_2, \dots)$
- Try decreasing λ
- Try increasing λ

fixes high variance
fixes high variance
fixes high bias
fixes high bias
fixes high bias
fixes high variance

没有合适的 training samples fit option

high bias → 主要是 model 不是很好

① 希望 low bias → $\lambda \downarrow$, $\lambda \downarrow$

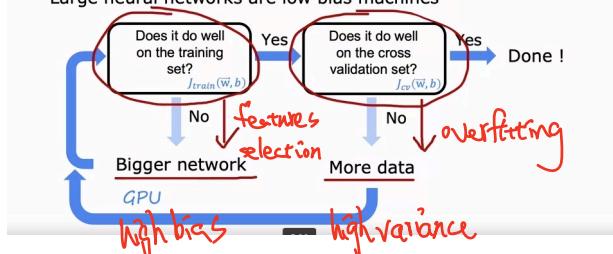
② 少 feature → model underfitting

③ polynomial fit features too

→ 模型拟合非线性

Neural networks and bias variance

Large neural networks are low bias machines



※ NN and bias/variance

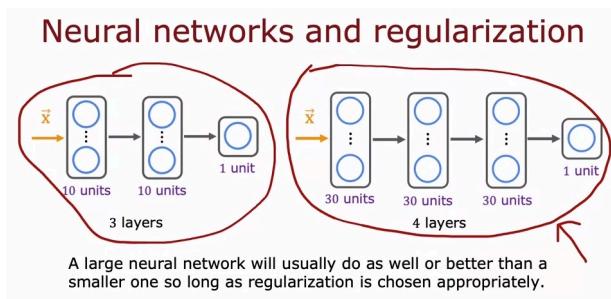
1. 通常大NN都能在小dataset中有小bias,不然就是NN要慢慢
放大(e.g. 多一层layers) 所以通常需要 solve 的是 high variance
大NN problem, 不会是大bias problem

2. performance: ① 达到 human performance baseline, 并看 J_{train} & J_{cv}

3. KNN 不一定会有小bias, 因为 造成小bias 需要 ① 交叉 feature 使 model fit 好一点
② $N \uparrow, \downarrow$

4. 加大NN不一定好, 在 feature 不夠構林時,
計算太 expensive ③ 複雜-交叉的 feature (non-linear)

5. 小改KNN來降低 bias 會引起 high variance (overfitting) → 对大NN用 regularization



6. check underfitting (high bias): 看 $J_{train} \approx J_{cv}$ 差距

check overfitting (high variance): 看 $J_{train} \gg J_{cv}$ 差距

※ NL process

< Error Analysis:

① Group the mis-classified examples

將分錯的做 grouping 並找出共通點 (eg. miss-spelling)

② 分錯的 "少數 group", 可 collect more data 來 model
去學 #examples 太少

③ spam email classification 也可透過 suspicious URLs 之間的
links 做為 feature 來 train URLs-based spam emails
classifier

④ the core idea from Bias/Variance

1. collect more data or not?

2. extract / minimize features?

3. do more/less on L2 Regularization λ ?

4. Think more on model / data

⑤ Error Analysis only works on what we are able to "observe" human

- ✗ Adding data on improving overfitting (high variance)

1. Error Analysis: After finding "weak" collect more data
to improve it for λ

2. Data Augmentation: Modify existing examples to better train
model (e.g. flip/distortion/rotate... on image)

⑥ The "distortion" should also be in the testing set!

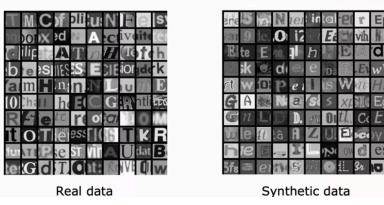
→ help us know whether the distortion for training process is useful

3. Data synthesis: creating brand new examples

1. Generate the realistic data similar to real data

2.

Artificial data synthesis for photo OCR



Focus more on 'data' could
be more efficient way to
collaborate with existed
advanced algorithms.

✗ Transfer learning

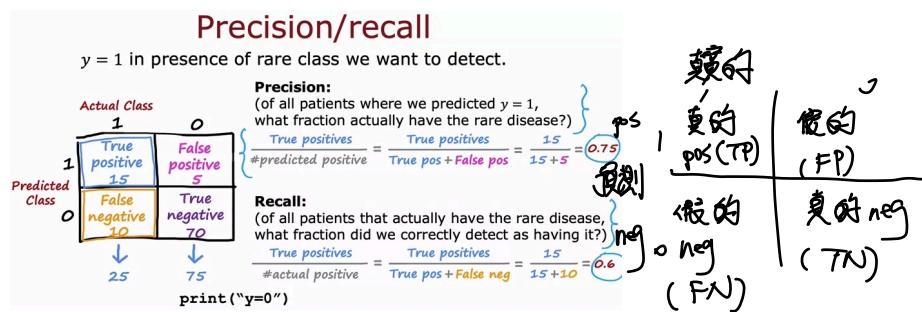
1. Use the pre-trained model, but only train the final output layer params

2. The dimension of our data and the pre-trained data should be the same.

3. Use the pre-trained model, then train it entirely on our dataset

-X Skewed datasets

1. Problem: if train a model, it would be good at classifying the type of majority in the pre-trained data
2. Error metric: Precision / Recall



衡量 algo 在 skewed dataset 上的 precision, recall

$$\frac{TP}{TP+FP} \quad \frac{TP}{TP+FN}$$

recall 的出現是為了防止 $TP+FP=0$

→ precision 和 recall 的 trade-off 依情況而定

High precision: classifier 真正識別出有病者
中, 真正病的機率高

High recall = $\frac{TP}{TP+FN}$ recall ↑, 但 FN ↓,
實際上無病 → classifier 將有病者
但認為沒病 認錯的機率最低

-X Decision Trees

1. How to choose features at each node?

→ Maximum purity

2. To estimate purity: "Entropy" *keep the tree smaller*

3. When to stop splitting? → avoid overfitting *pre-defined*

① 100% pure on classes

② When exceeding the depth of tree

- ③ When the improvements in purity score are below a threshold
 - ④ When the number of example in one node is below or threshold
- threshold
- ① max depth
 - ② improvement of purity
 - ③ #examples in nodes

4. Choose what features to split on?

→ main goal: maximum purity!

Key: When to stop splitting!

Define a decision node!

5. Entropy as a measure of impurity

$$H(p_1) = -p_1 \log(p_1) - p_2 \log(p_2)$$

$$= -p_1 \log(p_1) - (1-p_1) \log(1-p_1)$$

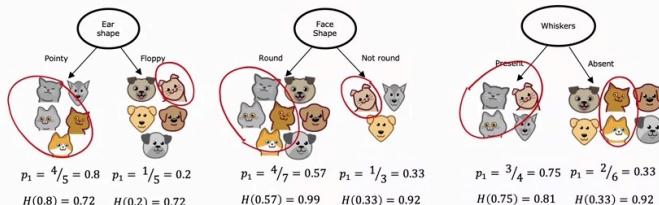
⑥ Decision Nodes: Information Gain

1. The choice of decision nodes hope to reduce the entropy

2. Information Gain = Reduction of entropy

3.

Choosing a split



$$+ H(0.8) = -0.1 \cdot \log(0.8) - 0.2 \cdot \log(0.2)$$

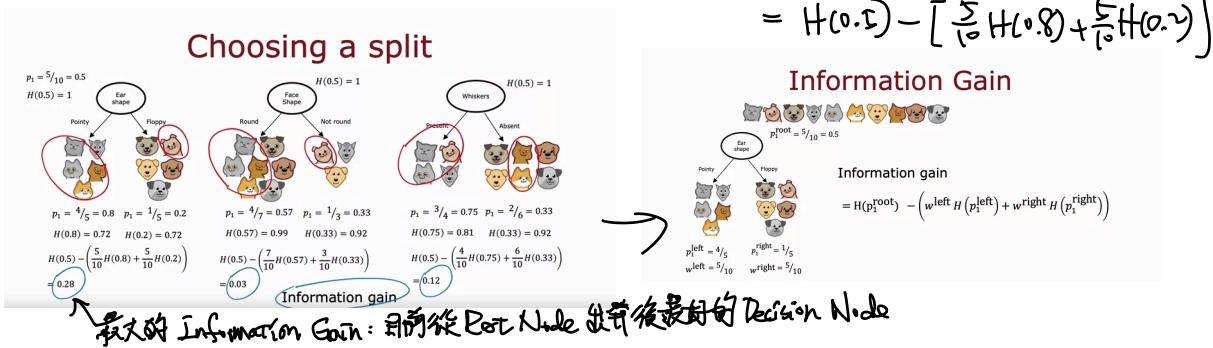
$$= H(0.2) = 0.72$$

上面 - weighted sum, weights by the # samples if it's true

$$\frac{4}{5} \cdot 0.72 + \frac{1}{5} \cdot 0.72 = 0.72$$

解釋: cat, dog 各半: $H(0.5) = -0.5 \log(0.5) - 0.5 \log(0.5)$

→ Information Gain = Reduction of entropy



Decision Tree Learning

- Start with all examples at the root node : We get the initial distribution
- Calculate information gain for all possible features, and pick the one with the highest information gain → ~~but it's empty~~
- Split dataset according to selected feature, and create left and right branches of the tree
- Keep repeating splitting process until stopping criteria is met:
 - When a node is 100% one class
 - When splitting a node will result in the tree exceeding a maximum depth
 - Information gain from additional splits is less than threshold
 - When number of examples in a node is below a threshold

止住

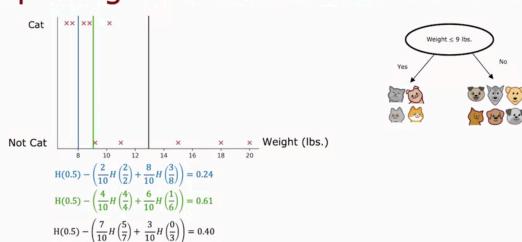
止住
1. 100%
2. max depth
3. class if sample < threshold
4. Information Gain < threshold

① Categorical Feature Encoding

- one-hot encoding: can be fed into NN, decision trees, logistic regression...
- label encoding → let input takes on more discrete value
- ordinal encoding
- interval encoding

② Encoding for continuous valued features

- Like we said, one-hot encoding takes care of discrete-number features
- Now, we need other methods to encode continuous-valued features
- Try different threshold, and calculate different information gain → then get the threshold

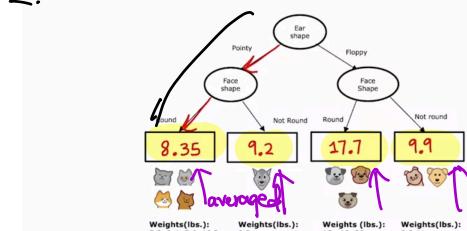


② Use decision tree to predict a number rather than a "class"

1. Called Regression Tree

2.

Regression with Decision Trees



if an animal with "Round face shape", the 8.35 is the averaged weights of training examples' weights

Regression with Decision Trees: Predicting a number

Ear shape	Face shape	Whiskers	Weight (lbs.)
Pointy	Round	Present	7.2
Floppy	Not round	Present	8.8
Floppy	Round	Absent	15
Pointy	Not round	Present	9.2
Pointy	Round	Present	8.4
Pointy	Round	Absent	7.6
Floppy	Not round	Absent	11
Floppy	Round	Absent	10.2
Floppy	Round	Present	18
Floppy	Round	Absent	20

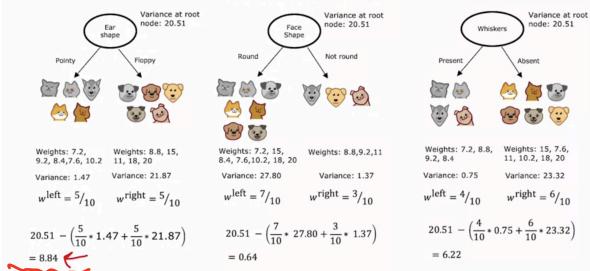
predict the weights of a given animal

3. 若在決定 decision node 時，分下來因爲只有對象的 continuous valued features

計算 information gain \rightarrow 看 continuous values of variance

The values of variance affect information gain, 依分數取量做 weighted sum.

Choosing a split



we get the "reduction in variance", then choose the subtree returns the "most reduction in variance"

即 information gain 一様，要選角原的 entropy / variance 成後 "最大" 者
為當下 decision node. 因而 改良愈多 (reduction in entropy / variance)
才是收斂的目標

③ Using multiple decision trees: Random Forest

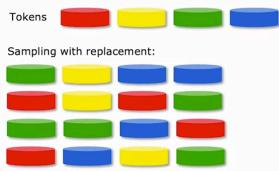
1. Why? Using single decision tree is sensitive to small changes in data

\rightarrow to be more robust, we need more decision tree to do ensemble

2. Build lots of decision tree to do tree ensemble

3. Tree Ensemble: Sampling with replacement: not get the same four tiles everytime

Sampling with replacement



And, some samples may repeat.

We use this technique to create new training for ensemble voting

Sampling with replacement

Ear shape	Face shape	Whiskers	Cat
Pointy	Round	Present	1
Floppy	Not round	Absent	0
Pointy	Round	Absent	1
Pointy	Not round	Present	0
Floppy	Not round	Absent	0
Pointy	Round	Absent	1
Floppy	Round	Present	1
Floppy	Not round	Present	1
Floppy	Round	Absent	0
Pointy	Round	Absent	1

4.

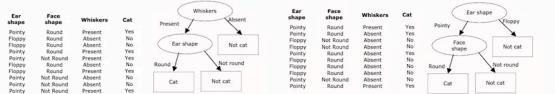
Generating a tree sample

Given training set of size m

For $b = 1$ to B

Use sampling with replacement to create a new training set of size m

Train a decision tree on the new dataset



Bagged decision tree

做 B 個 Sampling with Replacement, 每次再 train 一個 decision tree

B 的大小? $64/28\dots B$, 越大可能到某個 threshold 就不会 better 了。

↑
根據 data 這邊，沒有 default，但不會超過 100 。計算量太大且不準確。

通常可能會多次 root node 或 root node 附近的 node 在 end of Sampling with Replacement 都很像，除非更新每次 data 抽樣差異性很大

∴ 需要 randomize! 將原本可用的 features 取其 subset 不做為 decision node 的條件

Randomizing the feature choice

At each node, when choosing a feature to use to split, if n features are available, pick a random subset of $k < n$ features and allow the algorithm to only choose from that subset of features.

$$k = \sqrt{n}$$

Random forest algorithm

5. XGBoost: 這個 model focus 在分得不好的 samples 上繼續 train, 跟 L-Tree
這多個 decision tree 似 ensemble

Boosted trees intuition

Given training set of size m

For $b = 1$ to B :

Use sampling with replacement to create a new training set of size m
But instead of picking from all examples with equal $(1/m)$ probability, make it
more likely to pick misclassified examples from previously trained trees

Train a decision tree on the new dataset



1, 2, ..., b-1

分得不好的 samples 会在下輪中較高概率被抽到

XGBoost (eXtreme Gradient Boosting)

- Open source implementation of boosted trees
- Fast efficient implementation
- Good choice of default splitting criteria and criteria for when to stop splitting
- Built in regularization to prevent overfitting
- Highly competitive algorithm for machine learning competitions (e.g: Kaggle competitions)

② 不同于很多 sampling with replacement 的 dataset 来 train, XGBoost 会自行对 samples 进行采样
② 前述提到的 1. Sampling with Replacement 2. Randomizing the feature choice

都是 XGBoost 会用到的方法

① Decision Trees v.s. NN

- DT 适合处理结构化数据 (like Tabular data, spreadsheet) with
 - Categorical or Continuous-valued features
 - Predict for classification or Regression
- Interpretable: 较小的 DT 比较容易理解
- Ensemble of DT 在计算量上较大
- NN: AT 在 DT 的 structured data (spreadsheet-like) 工作得很好。
和 un-structured data & all types of data

Decision Trees vs Neural Networks

Decision Trees and Tree ensembles

- Works well on tabular (structured) data
- Not recommended for unstructured data (images, audio, text)
- Fast
- Small decision trees may be human interpretable

Neural Networks

- Works well on all types of data, including tabular (structured) and unstructured data
- May be slower than a decision tree
- Works with transfer learning
- When building a system of multiple models working together, it might be easier to string together multiple neural networks

① K-means Optimization objective

- $C^{(i)}$: index of cluster ($1, 2, 3, \dots, k$) to which example $x^{(i)}$ is currently assigned

M_k : cluster centroid k

$M_{C(i)}$: cluster centroid of cluster to which example $x^{(i)}$ has been assigned

Cost function:

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

Cost function for K-means

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

Repeat {

Assign points to cluster centroids

for $i = 1$ to m

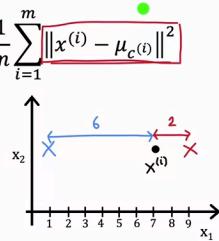
$c^{(i)}$:= index of cluster

centroid closest to $x^{(i)}$

Move cluster centroids

for $k = 1$ to K

μ_k := average of points in cluster k



}

Initializing the cluster centroids

1. Random Initialization

Random initialization

For $i = 1$ to 100 { 50-1000

Randomly initialize K-means. \leftarrow k random examples

Run K-means. Get $c^{(1)}, \dots, c^{(m)}, \mu_1, \mu_1, \dots, \mu_k \leftarrow$

Computer cost function (distortion)

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

}

Pick set of clusters that gave lowest cost J

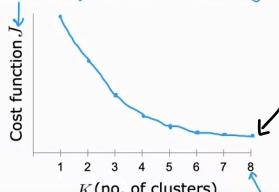
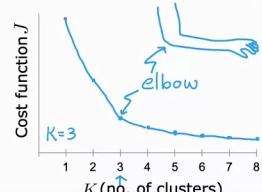
2. Choose the number of cluster centroids

Choosing the value of K

Elbow method

the right "K" is often ambiguous

Don't choose K just to minimize cost J

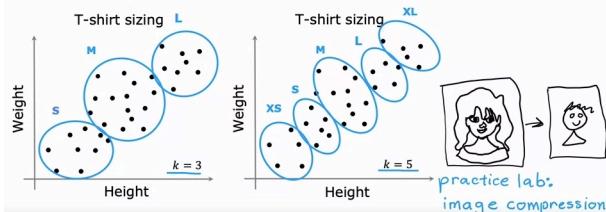


不能選擇 one - 腳步的 cluster,

這樣當然 less cost

Choosing the value of K

Often, you want to get clusters for some later (downstream) purpose.
Evaluate K-means based on how well it performs on that later purpose.



Which of these statements are true about K-means? Check all that apply.

If you are running K-means with $K = 3$ clusters, then each $c^{(i)}$ should be 1, 2, or 3.

Correct
 $c^{(i)}$ describes which centroid example(i) is assigned to. If $K = 3$, then $c^{(i)}$ would be one of 1, 2, or 3 assuming counting starts at 1.

If each example x is a vector of 5 numbers, then each cluster centroid μ_k is also going to be a vector of 5 numbers.

Correct
The dimension of μ_k matches the dimension of the examples.

The number of cluster assignment variables $c^{(i)}$ is equal to the number of training examples.

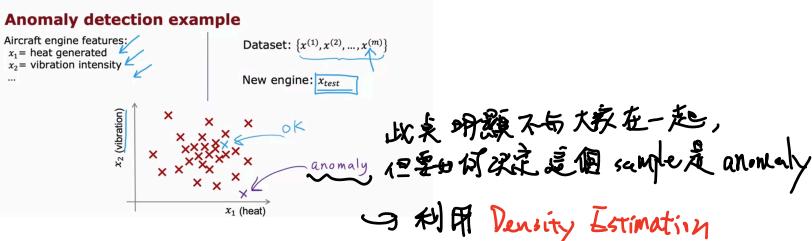
Correct
 $c^{(i)}$ describes which centroid example(i) is assigned to.

The number of cluster centroids μ_k is equal to the number of examples.

① Anomaly Detection

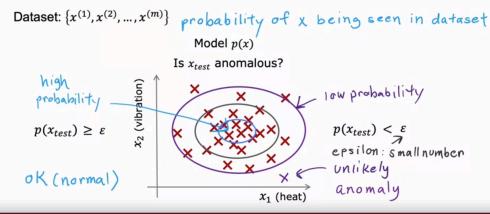
1. 異常檢測：從這些 features 做 embedding (or vectors)，判定新 sample 是不是

異常 data



2. Density Estimation

Density estimation



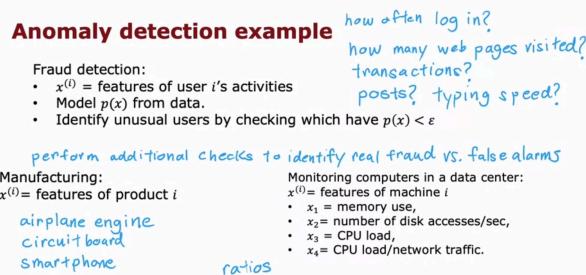
① 先以 training set + 找出常見的特徵做為 features. 建立 probability model,

(e.g. X_1, X_2 做為 features ; $X^{(1)} \sim X^{(m)}$ 都是 train set)

② 將新的 data 進來時，以建立的 probability model 計算 $p(X_{\text{test}})$ 。同時也有一個 threshold ϵ

③ 若 $p(X_{\text{test}}) > \epsilon$: Normal , else anomaly

Applications:



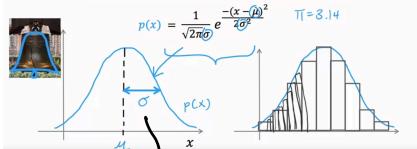
3. Methods

① Gaussian Distribution

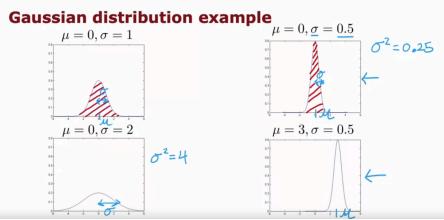
Gaussian (Normal) distribution

Say x is a number. Probability of x is determined by a Gaussian with mean μ , variance σ^2 .

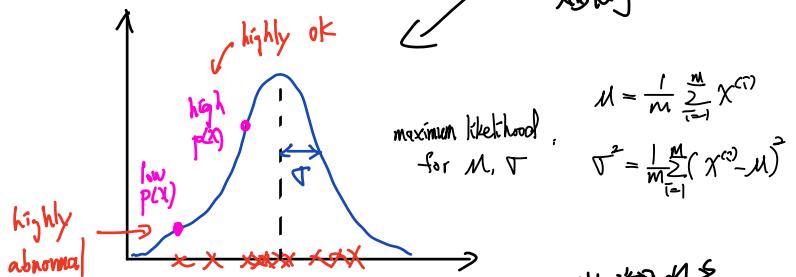
$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



↑↓, 尖山, 样本急集中



用在 Anomaly Detection:



2. Density Estimation with M, σ^2

Density estimation

Training set: $\{\vec{x}^{(1)}, \vec{x}^{(2)}, \dots, \vec{x}^{(m)}\}$
Each example $\vec{x}^{(i)}$ has n features

$$p(\vec{x}) = p(x_1; \mu_1, \sigma_1^2) * p(x_2; \mu_2, \sigma_2^2) * p(x_3; \mu_3, \sigma_3^2) * \dots * p(x_n; \mu_n, \sigma_n^2)$$

$$= \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2) \underset{\text{"add"}}{\sum} \underset{\text{"multiply"}}{\prod}$$

$p(x_1 = \text{high temp}) = 1/10$
 $p(x_2 = \text{high vibra}) = 1/20$
 $p(x_1, x_2) = p(x_1) * p(x_2)$
 $= \frac{1}{10} * \frac{1}{20} = \frac{1}{200}$

顯示出若某
temp、vibra 同時
發生的 prob 很低

針對每個 possible feature 都算 $p(x)$ ← feature of
當單獨的多

$$p(\vec{x}) = p(X_1; \mu_1, \sigma_1^2) p(X_2; \mu_2, \sigma_2^2) \dots p(X_n; \mu_n, \sigma_n^2)$$

將每個特徵視為條件獨立

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

Anomaly detection algorithm

- Choose n features x_i that you think might be indicative of anomalous examples.

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

- Given new example x , compute $p(x)$:

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

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

$$\vec{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix}$$

Anomaly if $p(x) < \epsilon$



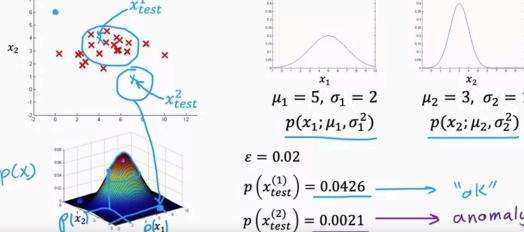
anomaly

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

當其中一個 $p(X_j)$ 很小時 → 代表著這群值，所有 $p(x)$ 相乘的結果也會超小 → $p(x) < \epsilon \rightarrow \text{anomaly}$

此時和下的問題是 它應該要少較合適

Anomaly detection example



The importance of real-number evaluation

When developing a learning algorithm (choosing features, etc.), making decisions is much easier if we have a way of evaluating our learning algorithm.

Assume we have some labeled data, of anomalous and non-anomalous examples.

Training set: $x^{(1)}, x^{(2)}, \dots, x^{(m)}$ (assume normal examples/not anomalous)
 $y=1$ for all training examples

Cross validation set: $(x_{cv}^{(1)}, y_{cv}^{(1)}), \dots, (x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$
Test set: $(x_{\text{test}}^{(1)}, y_{\text{test}}^{(1)}), \dots, (x_{\text{test}}^{(m_{\text{test}})}, y_{\text{test}}^{(m_{\text{test}})})$

根據次數的 prediction 及 labeling, 給了一張圖請問你做
自己會有一些 label off dataset

3. Evaluate the algorithm

Aircraft engines monitoring example

10000 good (normal) engines	2 to 50 flawed engines	$y=0$
2		$y=1$
Training set: 6000 good engines	train algorithm on training set	
CV: 2000 good engines ($y = 0$)	10 anomalous ($y = 1$)	tune ϵ tune x_j
Test: 2000 good engines ($y = 0$),	10 anomalous ($y = 1$)	tune ϵ tune x_j
Alternative: No test set use if very few labeled anomalous examples		
Training set: 6000 good engines	higher risk of overfitting	
CV: 4000 good engines ($y = 0$), 20 anomalous ($y = 1$)		

Algorithm evaluation

Fit model $p(x)$ on training set $x^{(1)}, x^{(2)}, \dots, x^{(m)}$
On a cross validation/test example x , predict

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

10
2000

Possible evaluation metrics:

- True positive, false positive, false negative, true negative
- Precision/Recall
- F_1 -score

Use cross validation set to choose parameter $\underline{\epsilon}$

course 2 weeks
skewed datasets

1. 以 2 个 anomaly sample 为例，分成 10 个到 cross-validation，

可以在 train 完后看 predict CV 时，三度设多少才可能都得 anomaly 拦下来。

2. 已训练完后，再到 test set 跑

3. 如果只有 2 个 anomaly，那就不是 test set 了，建议

train 和 CV，并把 anomaly 都放在 CV 让我们去 tune ϵ ，
这是唯一的方法且很可能造成 overfitting.

Anomaly Detection v.s. Supervised Learning

Anomaly detection vs. Supervised learning

Very small number of positive examples ($y = 1$). (0-20 is common).

Large number of negative ($y = 0$) examples.

$(P(x))$ Many different "types" of anomalies. Hard for any algorithm to learn from positive examples what the anomalies look like; future anomalies may look nothing like any of the anomalous examples we've seen so far.

Fraud 用于一直有

不同变形或

② Features 的分類

变动性/更新“高”

Large number of positive and negative examples.

20 positive examples

Enough positive examples for algorithm to get a sense of what positive examples are like, future positive examples likely to be similar to ones in training set.

Spam 用途已定

不变的 features
分类任务

Anomaly detection vs. Supervised learning

→ Fraud detection

→ Manufacturing - Finding new previously unseen defects in manufacturing. (e.g. aircraft engines)

→ Monitoring machines in a data center

→ Email spam classification

→ Manufacturing - Finding known, previously seen defects $y=1$ scratches

→ Weather prediction (sunny/rainy/etc.)

→ Diseases classification

features "Updating"

features "Known"

Anomaly Detection: Choose the Features

1. Supervised Learning: 一些“不重要或缺少一些不是很重要的 features 且不会对 model 性能产生大影响，因为 supervised 是有 “label” 的从 optimization

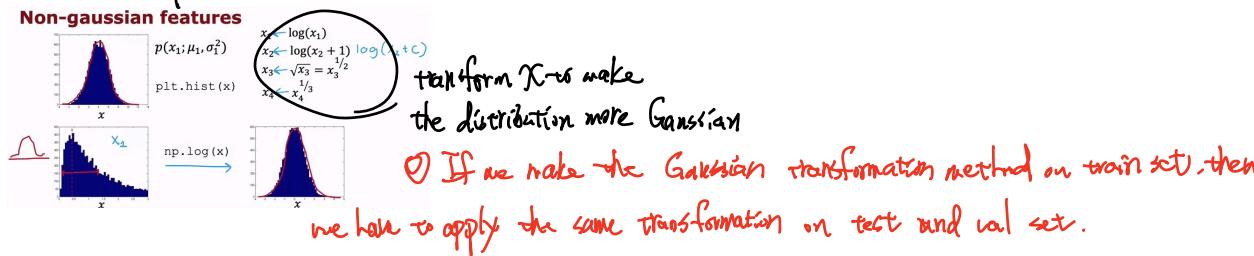
2. 但在无 label 的 condition 下，没有“对且重要的 features”才能让 algo 分得越来越好

3. So, carefully choose features & find important features is important in Anomaly Detection

4. We hope the distribution of data becomes more “Gaussian”

{ Make the non-Gaussian features more Gaussian





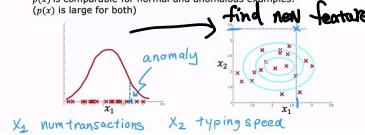
5. Error Analysis for Anomaly Detection

▷ If the trained model doesn't work well on val set, we have to apply error analysis.

② Error analysis for anomaly detection

Want $p(x) \geq \epsilon$ large for normal examples x .
 $p(x) < \epsilon$ small for anomalous examples x .

Most common problem:
 $p(x)$ is comparable for normal and anomalous examples.
 $(p(x)$ is large for both)



Monitoring computers in a data center

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

x_1 = memory use of computer
 x_2 = number of disk accesses/sec

high x_3 = CPU load
low x_4 = network traffic

$x_5 = \frac{\text{CPU load}}{\text{network traffic}}$
 $x_6 = \frac{(\text{CPU load})^2}{\text{network traffic}}$

Deciding feature choice based on $p(x)$

Large for normal examples;
Becomes small for anomaly in the cross validation set

Create new feature

X Collaborative Filtering

best-fit movie features

What if we have features of the movies?

Movie	Alice(1)	Bob(2)	Carol(3)	Dave(4)	x_1 (romance)	x_2 (action)
Love at last	5	5	0	0	0	0
Romance forever	5	?	?	0	0.9	0.01
→ Cute puppies of love	?	4	0	?	0.99	0
Nonstop car chases	0	0	5	4	0.1	1.0
Swords vs. karate	0	0	5	?	0	0.9

$$n_u = 4$$

$$n_m = 5$$

$$n = 2$$

$$x^{(1)} = [0.9]$$

$$x^{(3)} = [0.99]$$

For user 1: Predict rating for movie i as: $w^{(1)} \cdot x^{(1)} + b^{(1)}$ just linear regression

$$w^{(1)} = [5] \quad b^{(1)} = 0$$

$$x^{(3)} = [0.99]$$

$$w^{(1)} \cdot x^{(3)} + b^{(1)} = 4.95$$

→ For user j : Predict user j 's rating for movie i as $w^{(j)} \cdot x^{(j)} + b^{(j)}$

W, b 是根据不同的用户对每个特征 i 的值来调整的

Cost function

Notation:

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

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

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

For user j and movie i , predict rating: $w^{(j)} \cdot x^{(i)} + b^{(j)}$

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

To learn $w^{(j)}, b^{(j)}$

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

with per-item feature

what if we don't have advanced features?

Cost function

To learn parameters $w^{(j)}, b^{(j)}$ for user j :

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

To learn parameters $w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}, \dots, w^{(n_u)}, b^{(n_u)}$ for all users :

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

for each user

④ Collaborative filtering algorithm

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)	x_1 (romance)	x_2 (action)
Love at last	5	5	0	0	?	?
Romance forever	5	?	?	0	?	?
Cute puppies of love	?	4	0	?	?	?
Nonstop car chases	0	0	5	4	?	?
Swords vs. karate	0	0	5	?	?	?

$$\begin{aligned} w^{(1)} &= \begin{bmatrix} 5 \\ 0 \end{bmatrix}, w^{(2)} = \begin{bmatrix} 5 \\ 0 \end{bmatrix}, w^{(3)} = \begin{bmatrix} 0 \\ 5 \end{bmatrix}, w^{(4)} = \begin{bmatrix} 0 \\ 5 \end{bmatrix} \\ b^{(1)} &= 0, b^{(2)} = 0, b^{(3)} = 0, b^{(4)} = 0 \end{aligned}$$

$$\text{using } w^{(j)} \cdot x^{(i)} + b^{(j)} \\ w^{(1)} \cdot x^{(1)} \approx 5 \\ w^{(2)} \cdot x^{(1)} \approx 5 \\ w^{(3)} \cdot x^{(1)} \approx 0 \\ w^{(4)} \cdot x^{(1)} \approx 0$$

$$\rightarrow x^{(1)} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

↓

既在 $w^{(1)}$ 和 $w^{(2)}$ 那两种情况且都会是 $x^{(1)}$ 相同
才去做 linear regression & fitting

④ Collaborative Filtering

Collaborative filtering

$$\begin{aligned} \text{Cost function to learn } w^{(1)}, b^{(1)}, \dots, w^{(n_u)}, b^{(n_u)}: \\ \min_{w^{(1)}, b^{(1)}, \dots, w^{(n_u)}, b^{(n_u)}} \frac{1}{2} \sum_{j=1}^{n_u} \sum_{i:r(i,j)=1} (w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (w_k^{(j)})^2 \\ \text{Cost function to learn } x^{(1)}, \dots, x^{(n_m)}: \\ \min_{x^{(1)}, \dots, x^{(n_m)}} \frac{1}{2} \sum_{i=1}^{n_m} \sum_{j:r(i,j)=1} (w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2 \\ \text{Put them together:} \\ \min_{w^{(1)}, \dots, w^{(n_u)}, b^{(1)}, \dots, b^{(n_u)}, x^{(1)}, \dots, x^{(n_m)}} J(w, b, x) = \frac{1}{2} \sum_{(i,j):r(i,j)=1} (w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (w_k^{(j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2 \end{aligned}$$

代入部分

Cost function for binary application

Previous cost function:

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

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

$$\begin{aligned} L(f_{(w,b,x)}(x), y^{(i,j)}) &= -y^{(i,j)} \log(f_{(w,b,x)}(x)) - (1 - y^{(i,j)}) \log(1 - f_{(w,b,x)}(x)) \\ J(w, b, x) &= \sum_{(i,j):r(i,j)=1} L(f_{(w,b,x)}(x), y^{(i,j)}) \quad \text{cost for all examples} \\ &\quad g(w^{(j)} \cdot x^{(i)} + b^{(j)}) \end{aligned}$$

⑤ Mean Normalization

Users who have not rated any movies

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)	Eve (5)
Love at last	5	5	0	0	?
Romance forever	5	?	?	0	?
Cute puppies of love	?	4	0	?	?
Nonstop car chases	0	0	5	4	?
Swords vs. karate	0	0	5	?	?

$$\begin{aligned} \min_{w^{(1)}, \dots, w^{(n_u)}, b^{(1)}, \dots, b^{(n_u)}, x^{(1)}, \dots, x^{(n_m)}} & \frac{1}{2} \sum_{(i,j):r(i,j)=1} (w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (w_k^{(j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2 \\ w^{(5)} &= \begin{bmatrix} 0 \\ 0 \end{bmatrix}, b^{(5)} = 0, w^{(5)} \cdot x^{(i)} + b^{(5)} \end{aligned}$$

$$\begin{bmatrix} 5 & 5 & 0 & 0 & ? \\ 5 & ? & ? & 0 & ? \\ ? & 4 & 0 & ? & ? \\ 0 & 0 & 5 & 4 & ? \\ 0 & 0 & 5 & 0 & ? \end{bmatrix}$$

$$\begin{aligned} \text{mean } (5+5+0+?) / 4 &= 2.5 \\ (5+0)/2 &= 2.5 \\ (4+0)/2 &= 2 \\ (0+0+5+4)/4 &= 2.25 \\ (0+0+5+0)/4 &= 1.25 \end{aligned} \rightarrow \begin{bmatrix} 2.5 \\ 2.5 \\ 2 \\ 2.25 \\ 1.25 \end{bmatrix} = \mu$$

↓

$$x - \mu = \begin{bmatrix} 2.5 & 2.5 & -2.5 & -2.5 & ? \\ 2.5 & ? & ? & -2.5 & ? \\ ? & 2 & -2 & ? & ? \\ -2.5 & -2.5 & 2.25 & 1.75 & ? \\ -1.25 & -1.25 & 3.75 & -1.25 & ? \end{bmatrix}$$

For user j , on movie i predict:

$$w^{(j)} \cdot x^{(i)} + b^{(j)} + \mu$$

mean normalization: make algo run a little bit faster and gives reasonable predictions

Cost function

Given $w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}, \dots, w^{(n_u)}, b^{(n_u)}$

to learn $x^{(i)}$:

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

→ To learn $x^{(1)}, x^{(2)}, \dots, x^{(n_m)}$:

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

Gradient Descent

collaborative filtering

Linear regression (course 1)

repeat {

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

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

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

}

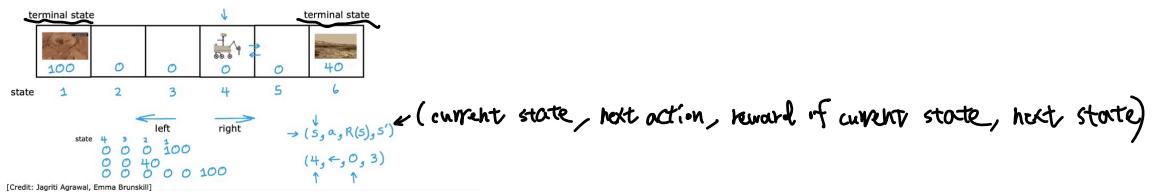
parameters w, b, x α is also a parameter

X Reinforcement Learning

1. "What to do", not "how to do it".
 algorithm has to learn itself

2.

Mars Rover Example



3. Return:

Return

	100	0	0		0	40
state	1	2	3	4	5	6
Return	$= 0 + (0.9)0 + (0.9)^2 0 + (0.9)^3 100 = 0.729 \times 100 = 72.9$					
Return	$= R_1 + \gamma R_2 + \gamma^2 R_3 + \dots$ (until terminal state)					

Discount Factor $\gamma = 0.9$ $0.99 \quad 0.999$ ← the discount factor should be close to 1
 $\gamma = 0.5$
 Return $= 0 + (0.5)0 + (0.5)^2 0 + (0.5)^3 100 = 12.5$ weighted by the discount factor

Discount Factor $\gamma = 0.999$ ← the discount factor should be close to 1
 很低的折扣因子会得到更高的回报
 (In financial set scenarios)

Example of Return

	100	50	25	12.5	6.25	40
state	1	2	3	4	5	6

← return $\gamma = 0.5$ ← reward

The return depends on the actions you take.

$4 : (0.5)^0 + (0.5)^1 0 + (0.5)^2 0 + (0.5)^3 100 = 12.5$

$3 : (0.5)^0 + (0.5)^1 0 + (0.5)^2 0 + (0.5)^3 40 = 5$

→ take the action based on what state we are at:

Example of Return

	100	50	25	12.5	6.25	40
state	1	2	3	4	5	6

The return depends on the actions you take.

$100 : (0.5)^0 + (0.5)^1 0 + (0.5)^2 0 + (0.5)^3 100 = 10$

$100 : (0.5)^0 + (0.5)^1 0 + (0.5)^2 0 + (0.5)^3 40 = 20$

Reward could be negative rewards!

discount factor

when gamma close to 1:

means the algorithm is "less impatient"

→ more patient

→ get approximate to the real reward

Policy

state	1	2	3	4	5	6
policy π						
action a	$\pi(1) = a$	$\pi(2) = a$	$\pi(3) = a$	$\pi(4) = a$	$\pi(5) = a$	$\pi(6) = a$
	100	50	25	12.5	6.25	40
	100	0	0	0	0	40

A policy is a function $\pi(s) = a$ mapping from states to actions, that tells you what action a to take in a given state s .

The goal of reinforcement learning

	100	0	0		40
state	1	2	3	4	6

$\pi(s) = a$

Find a policy π that tells you what action ($a = \pi(s)$) to take in every state (s) so as to maximize the return.

Mars rover		Helicopter		Chess	
states	6 states	position of helicopter	pieces on board	• Markov Decision Process (MDP)	
actions	$\leftarrow \rightarrow$	how to move control stick	possible move		
rewards	100, 0, 40	+1, -1000	+1, 0, -1		
discount factor γ	0.5	0.99	0.995		
return	$R_1 + \gamma R_2 + \gamma^2 R_3 + \dots$	$R_1 + \gamma R_2 + \gamma^2 R_3 + \dots$	$R_1 + \gamma R_2 + \gamma^2 R_3 + \dots$		
policy π		Find $\pi(s) = a$	Find $\pi(s) = a$		

1. future only depends on current state, not have anything prior.

choose best "a" based on $T(s)$

the "first reward" comes from "current state" without discount!

* State-action value function

State action value function (Q-function)

$Q(s, a) =$ Return if you

- start in state s ,
- take action a (once),
- then behave optimally after that.

100	50	25	12.5	20	40
100	0	0	0	0	40
$Q(2, \rightarrow)$	$Q(3, \leftarrow)$				
100 50 25 12.5 20 40	0 0 0 0 0 40				

100	50	25	12.5	20	40
100	0	0	0	0	40
$Q(2, \rightarrow)$	$Q(3, \leftarrow)$				
100 50 25 12.5 20 40	0 0 0 0 0 40				

terminal state with reward

$Q(s, a)$ returns the value by taking the action "a" with state "s"

$$\begin{aligned} &\text{--- return} \\ &\text{--- action} \\ &\text{--- reward} \\ \rightarrow & Q(2, \rightarrow) = 12.5 \\ & 0 + (0.5)0 + (0.5)^2 0 + (0.5)^3 100 \\ \checkmark & Q(2, \leftarrow) = 50 \\ & 0 + (0.5)100 \\ \downarrow & Q(4, \leftarrow) = 12.5 \\ & 0 + (0.5)0 + (0.5)^2 0 + (0.5)^3 100 \end{aligned}$$

goal: get the maximum of Q , so we have to pick the actions.

Picking actions

100	50	25	12.5	20	40
100	0	0	0	0	40

--- return
--- action
--- reward

100 100	50	12.5	25	6.25	12.5	10	6.25	20	40
100	0	0	0	0	0	0	0	0	40

$\max_a Q(s, a)$

$\rightarrow T(s) = a$

$Q(s, a) =$ Return if you

- start in state s ,
- take action a (once),
- then behave optimally after that.

The best possible return from state s is $\max_a Q(s, a)$.

The best possible action in state s is the action a that gives $\max_a Q(s, a)$. Optimal Q function

* Bellman Equation: Help us to compute the state action value function

Bellman Equation

$Q(s, a) =$ Return if you

- start in state s ,
- take action a (once),
- then behave optimally after that.

1	2	3	4	5

$R(s) =$ reward of current state

a : current action

s' : state you get to after taking action a

a' : action that you take in state s'

Bellman Equation

$$Q(s, a) = R(s) + \gamma \max_{a'} Q(s', a')$$

$$Q(s, a) = R(s) + \gamma \max_{a'} Q(s', a')$$

100	100	50	12.5	25	6.25	12.5	10	6.25	20	40
100	0	0	0	0	0	0	0	0	0	40
$Q(2, \rightarrow) = R(2) + 0.5 \max_{a'} Q(3, a')$										

$$Q(2, \rightarrow) = R(2) + 0.5 \max_{a'} Q(3, a')$$

$$= 0 + (0.5)25 = 12.5$$

$$Q(4, \leftarrow) = R(4) + 0.5 \max_{a'} Q(3, a')$$

$$= 0 + (0.5)25 = 12.5$$

Explanation of Bellman Equation

$Q(s, a) =$ Return if you

- start in state s ,
- take action a (once),
- then behave optimally after that.

\rightarrow The best possible return from state s' is $\max_{a'} Q(s', a')$

$$Q(s, a) = R(s) + \gamma \max_{a'} Q(s', a')$$

Reward you get right away

Return from behaving optimally starting from state s'

$$R_s + \gamma R_{s'} + \gamma^2 R_{s''} + \dots$$

$$Q(s, a) = R_s + \gamma [R_{s'} + \gamma [R_{s''} + \dots]]$$

Recursion

$$Q(s, a) = R(s) + \gamma \max_{a'} Q(s', a')$$

$$Q(2, \rightarrow) = R(2) + \max_{a'} Q(3, a')$$

$$Q(s, a) = R(s) + \max_{a'} Q(s', a')$$

we have to calculate the maximum with state s' and its

* Random/Stochastic Environment

possible actions a^*

1. In the natural/random environment, the real move of target is not 100% guaranteed, for example: The wheel of robot may slip to the opposite direction.
2. The goal: Choose a policy "IV" to maximum the average/expected discounted rewards.

Expected Return

100	0	0	0	0	$\frac{1}{4} \cdot 100$
1	2	3	4	5	6
0	0	0	0	100	0
0	0	0	0	0	100

↓

Expected Return = Average($R_1 + \gamma R_2 + \gamma^2 R_3 + \gamma^3 R_4 + \dots$)
 $= E[R_1 + \gamma R_2 + \gamma^2 R_3 + \gamma^3 R_4 + \dots]$

3. The modified Bellman Equation

$$Q(s, a) = R(s) + \gamma \max_{a'} Q(s', a')$$

2, ← ↓
1 or 3

4. For example:

For this problem, $\gamma = 0.25$. The diagram below shows the return and the optimal action from each state. Please compute $Q(5, \leftarrow)$.

100	25	6.25	2.5	10	40
100	0	0	0	0	40

← return
← action
← reward

$$\begin{aligned} Q(5, \leftarrow) &= R(5) + \gamma Q(4, \rightarrow) \\ &= R(5) + \gamma [R(4) + \gamma Q(5, \rightarrow)] \\ &= R(5) + \gamma R(4) + \gamma^2 R(5) + \gamma^3 R(6) \\ &= 0 + 0 + 0 + (0.25)^3 * 40 = 0.0625 + \frac{1}{4} * 40 = 0.625 \end{aligned}$$

* Continuous state spaces

Discrete vs Continuous State

Discrete State:

1	2	3	4	5	6

Continuous State:



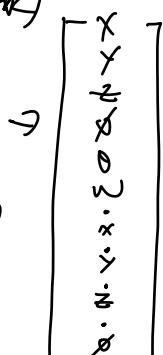
$s = \begin{bmatrix} x \\ y \\ \theta \\ \dot{x} \\ \dot{y} \\ \dot{\theta} \end{bmatrix}$

- position at (x, y)
- orientation $\in (0, 360)$
- Velocity in x, y direction
- and in orientation

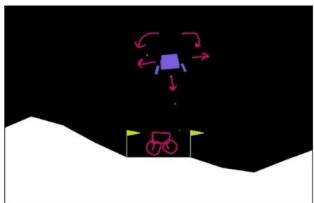
In Helicopter (直升机)

position: x, y, z

velocity: 左右滚翻 (roll)
 前.后倾斜 (pitch)
 左.右偏斜 (yaw)



Lunar Lander



actions:
do nothing
left thruster
main thruster
right thruster

$$s = \begin{bmatrix} x \\ y \\ \dot{x} \\ \dot{y} \\ \theta \\ \dot{\theta} \\ l \\ r \end{bmatrix}$$

0 or 1

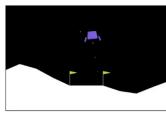
Lunar Lander Problem

Learn a policy π that, given

$$s = \begin{bmatrix} x \\ y \\ \dot{x} \\ \dot{y} \\ \theta \\ \dot{\theta} \\ l \\ r \end{bmatrix}$$

picks action $a = \pi(s)$ so as to maximize the return.

$$\gamma = 0.985$$



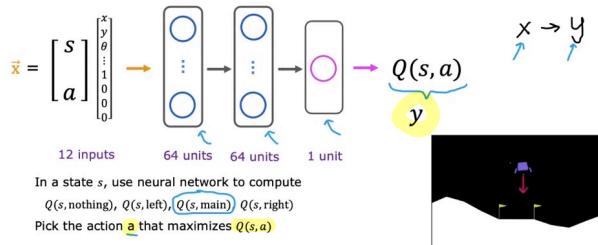
Given 's', use 'π' to pick

optimal 'a'

to maximize the return

- ① Learn the state-value function for continuous states space

Deep Reinforcement Learning



In the training data:

$$Q(s, a) = R(s) + \gamma \max_{a'} Q(s', a')$$

at first, the params of Q is random guess

$$f_{w_B}(x) \approx y$$

to construct our training data:

If we don't have the π (good policy) yet,
take actions randomly!

→ then apply Bellman equation to
improve the estimation of Q function

Bellman Equation

$$Q(x, a) = R(x) + \gamma \max_y Q(y, a)$$

random guess

$$f_{w_B}(x) \approx y$$

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

$$(s^{(1)}, a^{(1)}, R(s^{(1)}), s'^{(1)})$$

$$(s^{(2)}, a^{(2)}, R(s^{(2)}), s'^{(2)})$$

$$(s^{(3)}, a^{(3)}, R(s^{(3)}), s'^{(3)})$$

$$y^{(1)} = R(s^{(1)}) + \gamma \max_a Q(s'^{(1)}, a)$$

$$y^{(2)} = R(s^{(2)}) + \gamma \max_a Q(s'^{(2)}, a)$$

$$y^{(3)} = R(s^{(3)}) + \gamma \max_a Q(s'^{(3)}, a)$$

Learning Algorithm

Initialize neural network randomly as guess of $Q(s, a)$.

Repeat {

Take actions in the lunar lander. Get $(s, a, R(s), s')$.

Store 10,000 most recent $(s, a, R(s), s')$ tuples.



then apply Bellman equation to
optimize the params

in order to construct
more precise Q

$$\begin{array}{l} x, y \\ x', y' \\ \vdots \\ x^{10000}, y^{10000} \end{array}$$

② Algorithm refinement: E-greedy policy

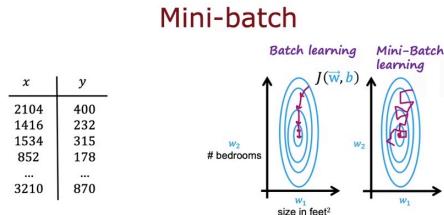
Which action should we choose?

2. 在 random initialize 的 circumstance 下，可能有些 action 的 Q 值都很大，但是这种 action 實際上是被加進來 train

→ try out something that might not be the best idea

How to choose actions while still learning?

- In some state s
- Option 1:
Pick the action a that maximizes $Q(s, a)$.
 - Option 2:
- With probability 0.95, pick the action a that maximizes $Q(s, a)$. **Greedy, "Exploitation"**
- With probability 0.05, pick an action a randomly. **"Exploration"**
- ϵ -greedy policy ($\epsilon = 0.05$)
- ~~0.95~~ → ϵ get lower = ϵ greedy
- ϵ greedy policy ($\epsilon = 0.05$)
- Start ϵ high → decrease the ϵ
 $1.0 \rightarrow 0.01$
Gradually decrease → less greedy → greedy
- ① Mini-batch & Soft Update ② gradual updates
③ convergence better
- ④ Mini-batch: ① Faster when GD ② only for large dataset because of local optima
- ⑤ Large dataset + noise造成的 loss 增大時，後續校正回來



Learning Algorithm

- Initialize neural network randomly as guess of $Q(s, a)$
- Repeat {
- Take actions in the lunar lander. Get $(s, a, R(s), s')$.
 - Store 10,000 most recent $(s, a, R(s), s')$ tuples.
- Train model:
1,000 Create training set of 10,000 examples using
 $x = (s, a)$ and $y = R(s) + \gamma \max_a Q(s', a')$
Train Q_{new} such that $Q_{new}(s, a) \approx y$.
Set $Q = Q_{new}$.
- Replay Buffer → Replay Buffer into mini-batches
- 因為 mini-batch 有時會因 samples 不夠，造成該次 loss 較大，此時若 $Q = Q_{new}$ 是不理想的，不能馬上 replace
∴ 需要 soft update → make gradual changes to Q

Soft Update

Set $Q = Q_{new}$.

w, b
 w_{new}, b_{new}

$$W = 0.01 W_{new} + 0.99 W_{old}$$

$$B = 0.01 B_{new} + 0.99 B_{old}$$

Limitations of Reinforcement Learning

- Much easier to get to work in a simulation than a real robot!
- Far fewer applications than supervised and unsupervised learning.
- But ... exciting research direction with potential for future applications.

→ update gradually, converge better