

e.g. 1st hidden ^{neuron} layer might look for a low vertical line

2nd hidden ^{neuron} layer might look for an oriented line

earliest hidden layers may look for short lines

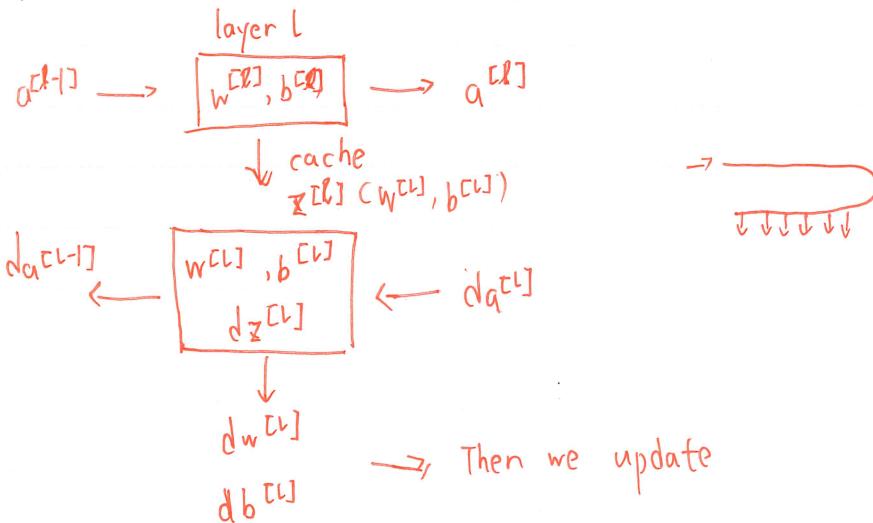
next layer...

(group of short lines) The first neuron may try to detect presence or absence of an eye

The second neuron — corner of a nose / bottom of a ear

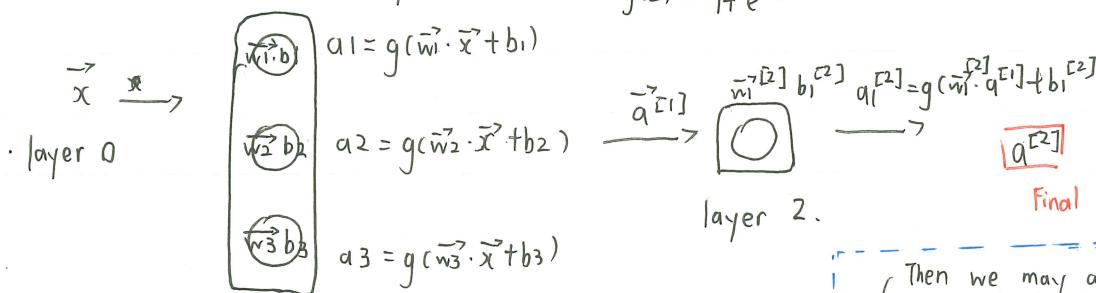
Finally — features of face shapes.

Nobody told it how to detect, all by itself!



Neural network layer

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



notation of layer numbering [1], [2]

$$\vec{w}_1 \rightarrow \vec{w}_1^{[1]} \dots \vec{w}_1^{[k]}$$

$$\vec{b}_1 \rightarrow b_1^{[1]} \dots b_1^{[k]}$$

$$\vec{a}_1 \rightarrow a_1^{[1]} \dots a_1^{[k]}$$

Then we may apply a threshold:
is $a^{[2]} \geq 0.5$?
Yes: $\hat{y} = 1$ No: $\hat{y} = 0$

Activation value of layer l, unit(neuron)j

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

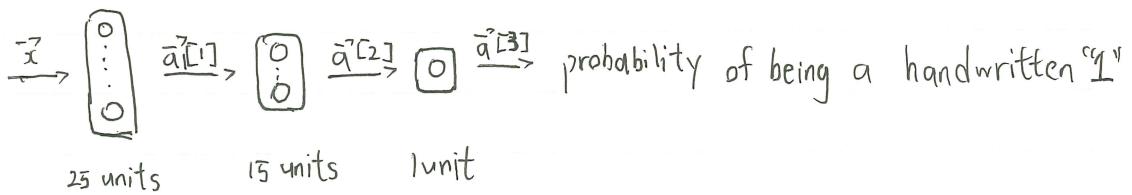
↑ ↑ ← parameters w&b of layer l, unit j
output of layer l-1

sigmoid

(activation function)

Forward propagation

Handwritten digit recognition:



forward propagation

propagating the activations of neurons.

Choosing $g(z)$ based on the demand of y

Output Layer:

Binary classification \rightarrow Sigmoid $y=0 \text{ or } 1$

Regression \rightarrow Linear activation $y = -/+$

Regression \rightarrow ReLU $y=0 \text{ or } +$ (can't below 0)

Hidden Layer:

Most common choice: — ReLU:

- Faster to compute
- flat at only one place (when $z < 0$) make gradient descent faster

- if $g(z) = z$ for all hidden layers, it's equivalent to linear regression because a linear function of a linear function is a linear function
- if g is always sigmoid, it's equivalent to logistic regression.
- The disable feature of the ReLU activation enables models to switch together linear segments to model complex non-linear functions.



Multiclass Classification

Softmax regression (N possible outputs)

$$z_j = \vec{w}_j \cdot \vec{x} + b_j \quad j=1 \dots N$$

$$a_j = \frac{e^{z_j}}{\sum_{k=1}^N e^{z_k}}$$

the exponential in the numerator of softmax magnifies small differences in value

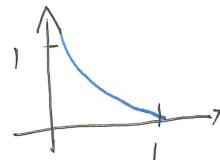
Note: $a_1 + a_2 + \dots + a_N = 1$

Logistic regression can be seen as a special case of softmax

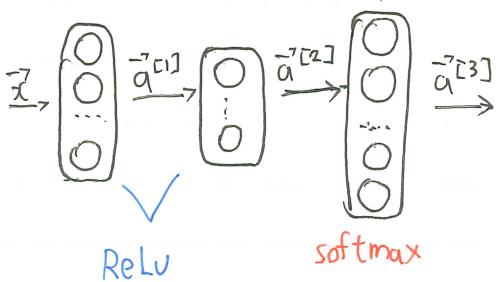
Cost:

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

loss = Sparse Categorical Crossentropy



Neural Network with Softmax output:



25 units 10 units 10 units

10 classes

$$\begin{aligned} z_1^{[3]} &= \vec{w}_1 \cdot \vec{a}^{[2]} + b_1^{[3]}, \quad a_1^{[3]} = \frac{e^{z_1^{[3]}}}{e^{z_1^{[3]}} + \dots + e^{z_{10}^{[3]}}} \\ &= P(y=1 | \vec{x}) \\ z_{10}^{[3]} &= \vec{w}_{10} \cdot \vec{a}^{[2]} + b_{10}^{[3]}, \quad a_{10}^{[3]} = \frac{e^{z_{10}^{[3]}}}{e^{z_1^{[3]}} + \dots + e^{z_{10}^{[3]}}} \\ &= P(y=10 | \vec{x}) \end{aligned}$$

logistic regression: $a_j^{[i]}$ only relates to $\vec{z}_j^{[i]}$
softmax:

$a_j^{[i]}$ relates to $a_1^{[i]}, \dots, a_N^{[i]}$

TensorFlow implementation.

`x = np.array([[200, 17]])` → [200 17] 1x2 matrix

`s = np.array([[200], [17]])` → $\begin{bmatrix} 200 \\ 17 \end{bmatrix}$ 2x1 matrix.

`x = np.array([200, 17])` → 1D vector.

Activation vector:

```
x = np.array([[200.0, 17.0]]).
```

`layer1 = Dense (units = 3, activation = 'sigmoid').`

$o1 = \text{layer_1}(x)$

tf.Tensor([0.2 0.7 0.3]), shape=(1,3) · dtype=float32)
similar to matrix 1x3 matrix 32 bits

al.numpy()

transfer from tensor to numpy

layer-2 = Dense Cunits = 25, ...

```
model = Sequential([layer_1, layer_2, ...])
```

```
model.compile(...)
```

```
x = np.array([ [0, ..., 245, ..., 17],  
              [ ..., ..., ... ] ])
```

```
y = np.array([1, 0])
```

```
model1.fit(x,y)
```

model.predict(x_new)

Neural network implementation in python

```
X = np.array([[200, 17]])  
W = np.array([[1, -3, 5], [-2, 4, -6]])  
B = np.array([-1, 1, 2]).
```

```
def dense(A_in, W, B):  
    Z = np.matmul(A_in, W) + B. / A_in @ W + B  
    A_out = g(Z) matrix multiplication  
    return A_out
```

```
def sequential(x):  
    a1 = dense(x, w1, b1).  
    .....
```

```
a4 = dense(a3, w4, b4)  
f_x = a4  
return f_x
```

```
import tensorflow as tf.  
from tensorflow.keras import Sequential.  
model = Sequential()  
[  
    tf.keras.Input(shape=(2,)), < from tensorflow.keras.layers import Dense  
    Dense(3, activation='sigmoid', name='layer1'),  
]  
① specify the model  
)  
model.summary()  
w1, b1 = model.get_layer("layer1").get_weights()  
model.compile()  
② compile the model  
loss = tf.keras.losses.BinaryCrossentropy()  
optimizer = tf.keras.optimizers.Adam(learning_rate=0.01), Using specific loss function.  
)  
model.fit(  
    xt, yt,  
    epochs=10,  
    ③ train the model  
)  
w1, b1 = ...  
predictions = model.predict(x-test)  
yhat = (predictions >= 0.5).astype(int)
```

Model Training Steps

- ① specify how to compute output given input x and parameters w, b

logic regression $f_{\tilde{w}, b}(\tilde{x}) = ?$

$$\tilde{z} = \text{np.dot}(w, x) + b$$

$$f_{\tilde{x}} = 1 / (1 + \text{np.exp}(-\tilde{z}))$$

neural
network

model = Sequential [

Dense(...)

]

- ② specify loss and cost

$$L(f_{\tilde{w}, b}(\tilde{x}), y)$$

binary cross entropy.

model.compile(

$$\text{loss} = -y * \text{np.log}(f_{\tilde{x}}) - (1-y) * \text{np.log}(1-f_{\tilde{x}}) \quad \text{loss} = \text{BinaryCrossentropy}()$$

$$J(\tilde{w}, b) = \frac{1}{m} \sum_{i=1}^m L(f_{\tilde{w}, b}(\tilde{x}^{(i)}), y^{(i)}) \quad \text{for numbers} - \text{MeanSquaredError}()$$

- ③ Train on data to minimize $J(\tilde{w}, b)$

model.fit(x, y, epochs=100)

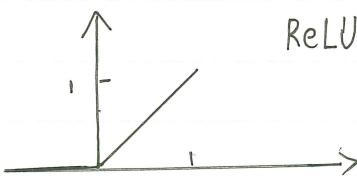
$$w = w - \alpha * \nabla_w J \quad \text{gradient descent}$$

$$b = b - \alpha * \nabla_b J$$

tensorflow compute partial derivative

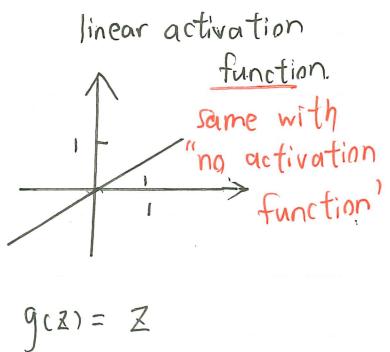
by using "back propagation"

Activation Functions:



$$g(z) = \max(0, z)$$

ReLU (Rectified Linear Unit)



$$g(z) = z$$

* Numerical Roundoff Errors

$$2.0 / 10000 \rightarrow 0.000200\ldots$$

$$1 + \frac{1}{10000} - (1 - \frac{1}{10000}) \rightarrow 0.000199\ldots$$

Due to limited memory

Original loss for logistic regression

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

More accurate loss

$$\text{loss} = -y \log(\frac{1}{1+e^{-z}}) - (1-y) \log(1 - \frac{1}{1+e^{-z}}) \quad \text{(less legible)}$$

Dense (units=1), activation = 'sigmoid' ('linear')

model.compile(loss='BinaryCrossEntropy', from_logits=True)
 $f_x = \text{tf.nn.sigmoid(model}(X))$

Then Tensorflow can rearrange terms

For softmax:

to avoid very small / big numbers

More accurate:

$$L(\vec{a}, y) = \begin{cases} -\log \frac{e^{z_1}}{e^{z_1} + \dots + e^{z_{10}}} & \text{if } y=1 \\ -\log \frac{e^{z_{10}}}{e^{z_1} + \dots + e^{z_{10}}} & \text{if } y=10 \end{cases}$$

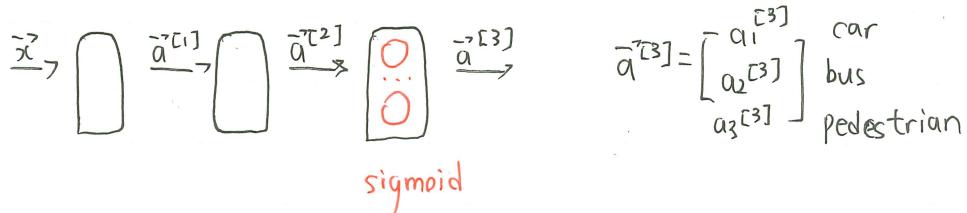
model.compile(loss='SparseCategoricalCrossEntropy', from_logits=True)

Be aware the result of X would be z_1, \dots, z_{10}

instead of a_1, \dots, a_{10}

$f_x = \text{tf.nn.softmax(model}(X))$

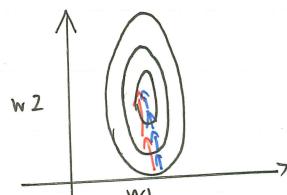
Multilabel classification different from multiclass



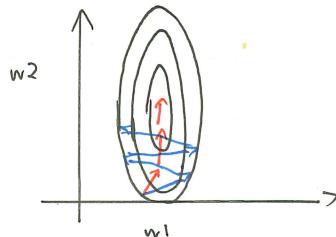
Depending on how Gradient descent is proceeding, the "Adam" algorithm in Tensorflow will adjust learning rate automatically

Adam = Adaptive Moment estimation.

$$w_j = w_j - \alpha_j \frac{\partial}{\partial w_j} J(\vec{w}, b)$$



If w_j (or b) keeps moving in same direction.
increase w_j



If w_j (or b) keeps oscillating
reduce w_j

model.compile(optimizer=tf.keras.optimizers.Adam(learning_rate=1e-3))
try to modify it!

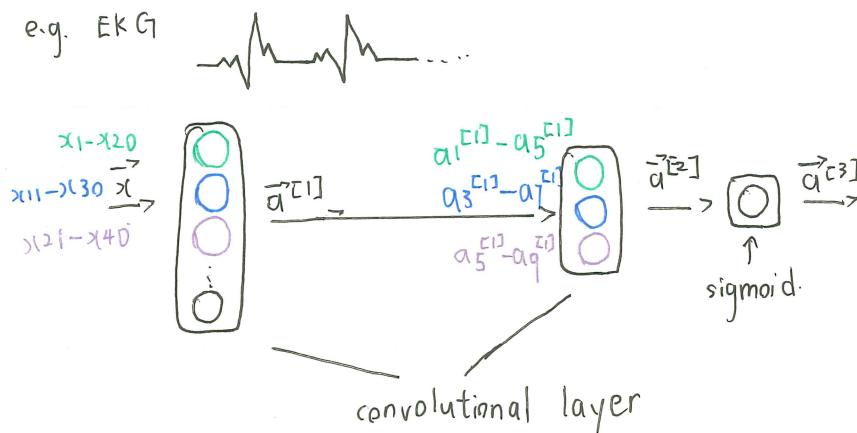
Convolutional Layer

Each Neuron only looks at part of previous layer's inputs.

- Faster computation
- less prone to overfitting

Convolutional Neural Network 卷积神经网络

e.g. EKG



Back propagation.

$$\begin{array}{c} \frac{\partial J}{\partial w} \\ \text{w} \end{array} \leftarrow \begin{array}{c} \frac{\partial J}{\partial c} \\ c \end{array} \leftarrow \begin{array}{c} \frac{\partial J}{\partial a} \\ a \end{array} \leftarrow \begin{array}{c} \frac{\partial J}{\partial d} \\ d \end{array} \leftarrow \begin{array}{c} \frac{\partial J}{\partial b} \\ b \end{array}$$

$\xrightarrow{\text{forward}}$
 $\xleftarrow{\text{back}}$

compute. $\frac{\partial J}{\partial a}$ once to compute both $\frac{\partial J}{\partial w}$ and $\frac{\partial J}{\partial b}$.

If model has N nodes and P parameters, we only comput.,
N+P step rather than N×P steps

Steps for backprop:

- ① calculate the local derivative (s) of the node
- ② use the chain rule

Debug a learning algorithm.

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

encounter unacceptably large errors:

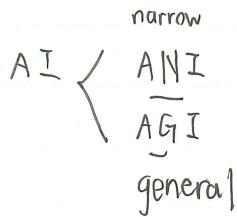
- Get more training examples fix high variance
- Try smaller sets of features fix high variance
- Try getting additional features fix high bias
- Try adding polynomial features. ($x_1^2, x_2^2, x_1x_2, \dots$) same
- Try decreasing λ / increasing λ
fix high bias high variance

Evaluate your model

Divide data set into training set and test set

Fit parameters by minimizing cost function $J(\tilde{w}, b)$

compute test error and training error



Model selection

$$d=1 \rightarrow \hat{f}_{\vec{w}, b}(\vec{x}) = w_1 x_1 + b \rightarrow J_{\text{test}}(w^{(1)}, b^{(1)})$$

$$d=10 \rightarrow \hat{f}_{\vec{w}, b}(\vec{x}) = w_{11}x_1 + w_{22}x_2^2 + \dots + w_{10}x_{10} + b \rightarrow J_{\text{test}}(w^{(10)}, b^{(10)})$$

Assume $J_{\text{test}}(w^{(5)}, b^{(5)})$ performs the best, we choose $d=5$

How well does the model perform? Directly report $J_{\text{test}}(w^{(5)}, b^{(5)})$?

No. It's likely to be an optimistic estimate.

Training, Cross validation, Test set

J_{train} J_{cv} J_{test}

$$6 \quad 2 \quad 2 \quad \rightarrow J_{\text{cv}}(w^{(1)}, b^{(1)})$$

$$\rightarrow J_{\text{cv}}(w^{(10)}, b^{(10)})$$

Pick d based on J_{cv} .

Estimate generalization error using test the set J_{test}

Not diagnostic

For model with

High bias (underfit)

"Just right"

High variance (overfit)

J_{train}

High

Low

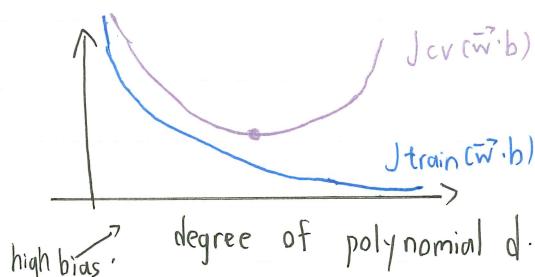
Low

J_{cv}

High

Low

High



high bias : degree of polynomial d.

High bias : $J_{\text{train}} \approx J_{\text{cv}}$ (both high)

High variance: $J_{\text{cv}} \gg J_{\text{train}}$ (J_{train} may be low)

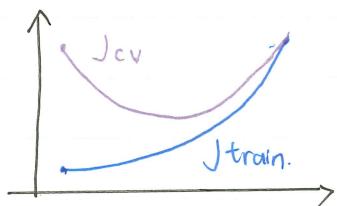
(Sometimes) High variance and high variance:

J_{train} will be high and $J_{\text{cv}} \gg J_{\text{train}}$.



Linear regression with regularization.

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



↑ high bias (reason why it looks familiar)

speech recognition example:

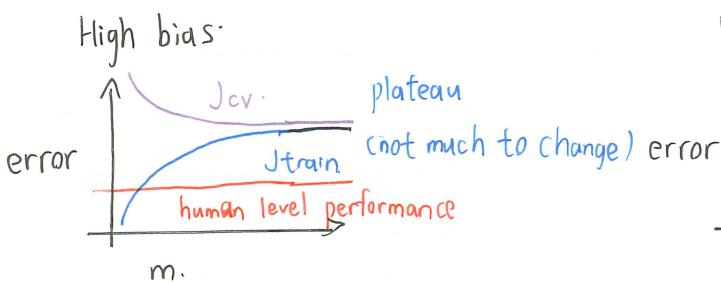
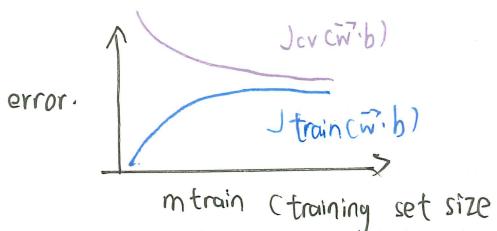
Human level performance:	10.6%	$\uparrow 0.2\%$	10.6%	10.6%
Training error:	10.8%		15.0%	15.0%
Cross validation error:	14.8%	$\downarrow 4.0\%$	15.5%	19.7%

High Variance High Bias Both High.

- Establish a baseline level of performance.

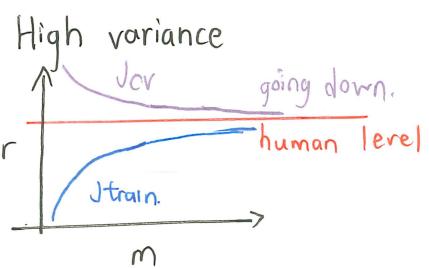
- Human level performance
- Competing model
- Guess based on experience

Learning curves



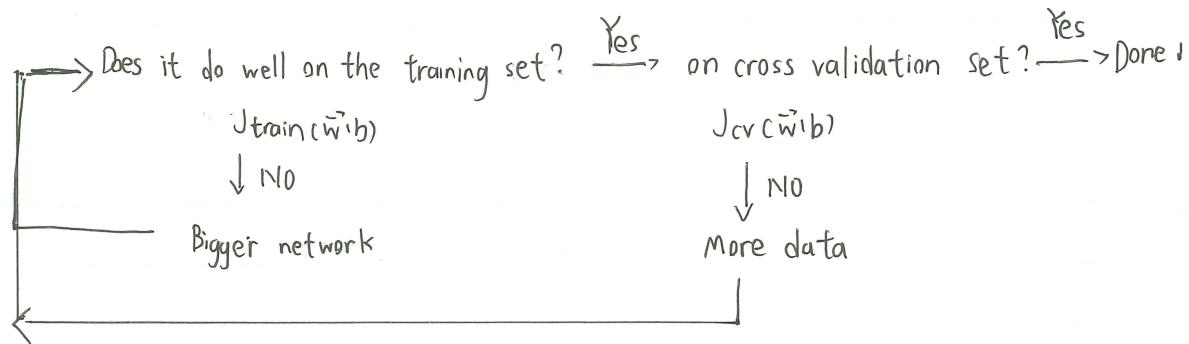
If an algorithm has high bias,
It's useless to add data

Can be used to visualize, but computing different subset of training set
is expensive.



If has high variance,
getting more training data
is likely to help

Neural networks and bias variance



By carefully choosing regularization, large model can perform well

Caveat: it's slow, but neural network is accurate with low bias

λ
 $\text{layer} = \text{Dense}(\text{units} = 25, \text{activation} = \text{"relu"}, \text{kernel_regularizer} = \text{L2}(0.01))$

Normally we don't regularize ' b ', which may cause difference.

Error Analysis

e.g. Manually analyze misclassified examples,

Categorize them based on common traits

If there's too many, fit the model with less examples..

or employ more features

come up with specific solutions.

Add data:

Add more data of types where error analysis has indicated.

Data augmentation (audios and videos)

Modifying an existing training example to create a new training example
e.g. distort / add noise must has meaning

Data synthesis

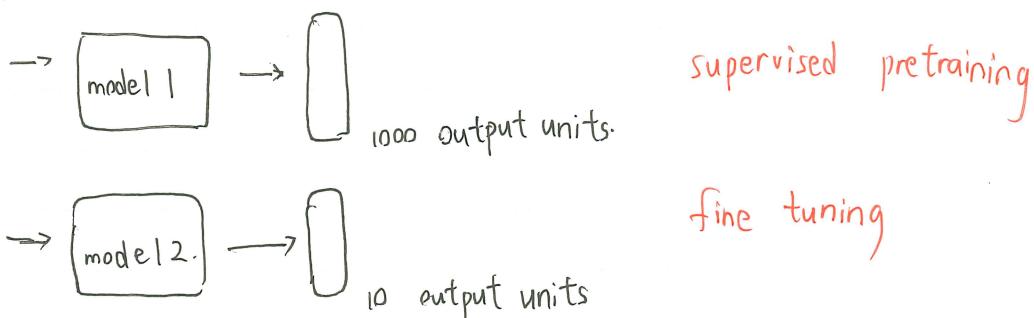
Generating data from scratch

e.g. (mostly visual) using fonts in the computer as photo

$$\text{AI} = \text{Code} + \text{Data}$$

Conventional model-centric ↑ Data-centric approach

Transfer learning: using data from a different task,



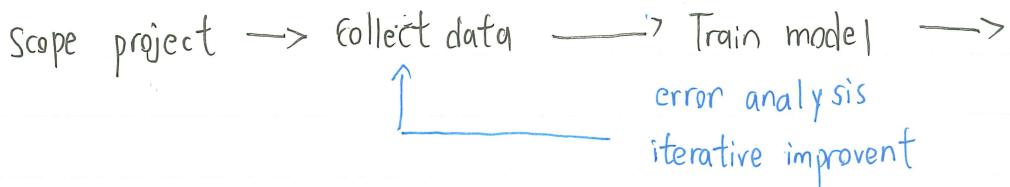
Option 1: Only train output layers parameters works better for small sets
 Option 2: train all parameters

∴ Use other's model, replace with own output layer /

do a little bit fine tune, further train the network

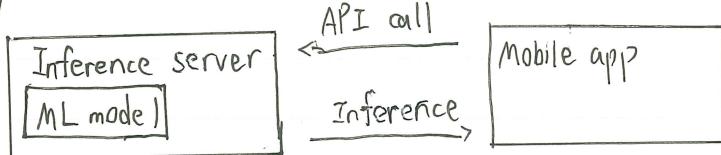
restriction: input type must be same

Full cycle:



Deploy in production
monitor and maintain system

Deployment



Software engineering may be needed.

MLOps - machine learning operation

practice to make the model reliable, scales well, has good laws,
is monitored, can be easily updated

Rare disease classification example — skewed datasets

99.5% accuracy 0.5% error less useful always prints $y=0$

99.0% accuracy 1% error more useful at least make diagnosis

		Actual class	
		1 (在)	0 (不)
predicted class	1	True positive	False positive
	0	False negative	True negative

Precision: $\frac{\text{True positive}}{\text{Predicted positive}}$ detect if 1 all the time.

Recall: $\frac{\text{True positive}}{\text{Actual positive}}$ detect if 0 all the time

algorithm with low precision and recall is useless.

Logistic regression : $0 < f_{\text{wib}}(\vec{x}) < 1$.

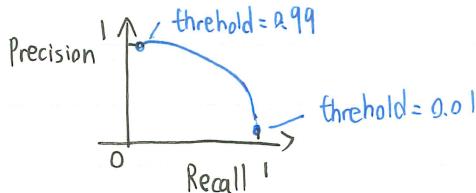
Predict 1 if $f_{\text{wib}}(\vec{x}) \geq 0.5$
0 if $f_{\text{wib}}(\vec{x}) < 0.5$ \rightarrow threshold.

Suppose we want $y=1$ only if very confident

Raise threshold \rightarrow higher precision, lower recall

Suppose we want to avoid missing too many cases

Reduce threshold \rightarrow lower precision, higher recall



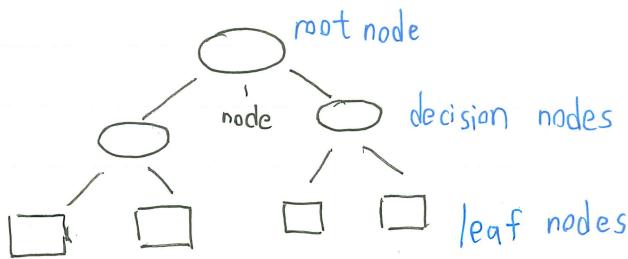
F₁ Score:

$$= \frac{1}{\frac{1}{2}(\frac{1}{P} + \frac{1}{R})} = 2 \frac{PR}{P+R} \quad (\text{harmonic mean})$$

Average $\frac{P+R}{2}$ is bad.

taking an average that
emphasize the smaller values more.

Decision tree

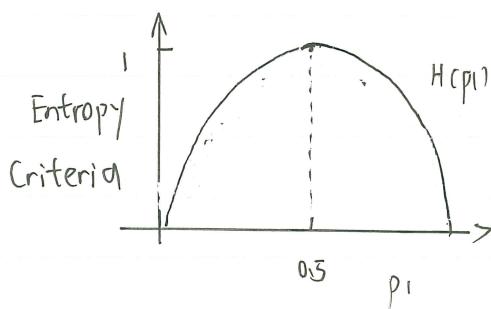


Key decision 1: How to choose what feature to split on at each node?
 Maximize purity (minimize impurity (entropy))

Key decision 2: When do you stop splitting?

- When a node is 100% one class
- When splitting a node will result in the tree exceeding a maximum depth
- When improvements in purity score are below a threshold
- When number of examples in a node is below a threshold
 (avoid overfitting)

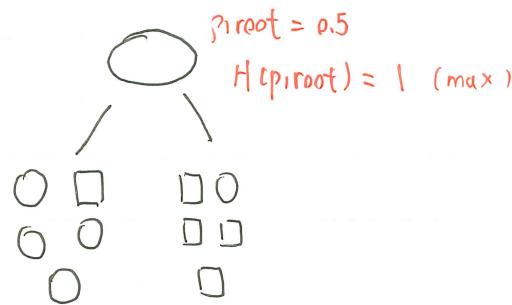
Measure purity - entropy



$$\begin{aligned}
 p_0 &= 1 - p_1 && \text{looks similar?} \\
 H(p_1) &= -p_1 \log_2(p_1) - p_0 \log_2(p_0) \\
 &= -p_1 \log_2(p_1) - (1-p_1) \log_2(1-p_1) \\
 \text{Note } 0 \log_2(0) &= 0
 \end{aligned}$$

Gini criteria is similar

Choosing a split — Information Gain.



$$p_1^{\text{left}} = \frac{4}{5} \quad p_1^{\text{right}} = \frac{1}{5}$$

$$w_{\text{left}} = \frac{5}{10} \quad w_{\text{right}} = \frac{5}{10}$$

weight

Information gain =

$$= H(p_1^{\text{root}}) - (w_{\text{left}} H(p_1^{\text{left}}) + w_{\text{right}} H(p_1^{\text{right}}))$$

- We measure the reduction (IG) to decide whether it's worth splitting
- We add weight because entropy with smaller samples is better than bigger samples

Decision Tree Learning

- Start at root node
- Pick the feature with highest information gain
- Split dataset, create left and right branches of the tree.
- Keep repeating splitting process until stopping criteria
 - When a node is 100% one class
 - When splitting a node will result in the tree exceeding max depth
 - Information gain from additional splits is less than threshold.
 - When number of examples in a node is below a threshold

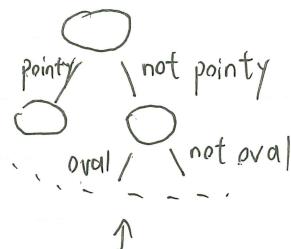
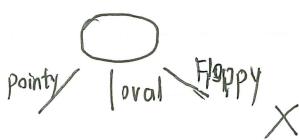
Cover other subtrees and check every node!

Algorithm that calls itself

Recursive splitting —

Build the overall decision tree by build smaller sub-decision trees and then putting them together.

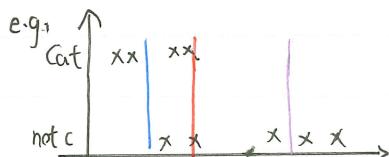
When encounter more than two features.



One hot encoding:

If a categorical feature can take on k values, create k binary features (0 or 1)

When encounter continuous features (numerical)



choose midpoints as possible splits

$$\text{weight} \leq 8 \quad H(0.5) - \left(\frac{2}{10}H\left(\frac{2}{5}\right) + \frac{8}{10}H\left(\frac{3}{8}\right)\right) = 0.24$$

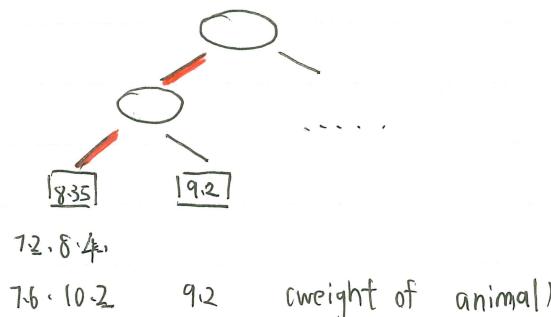
$$\text{weight} \leq 9 \quad H(0.5) - \left(\frac{4}{10}H\left(\frac{4}{7}\right) + \frac{6}{10}H\left(\frac{1}{6}\right)\right) = 0.61 \quad \checkmark$$

$$\text{weight} \leq 16 \quad H(0.5) - \left(\frac{7}{10}H\left(\frac{5}{7}\right) + \frac{3}{10}H\left(\frac{0}{3}\right)\right) = 0.40$$

(1bs)

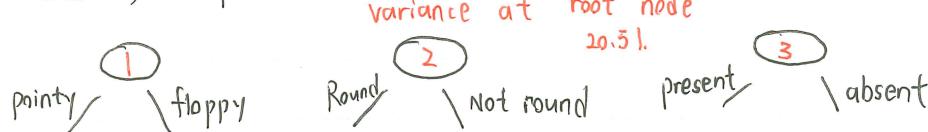


Regression Trees



When given a new example, follow the decision nodes down as usual until it gets to a leaf node and predict the example's weight based on the average of the weights in leaf node

choosing a split:



Variance:	1.47	21.87	27.80	1.37	0.75	23.32
weight	$\frac{5}{10}$	$\frac{5}{10}$	$\frac{7}{10}$	$\frac{3}{10}$	$\frac{4}{10}$	$\frac{6}{10}$

Reduction in variance: ①: $20.51 - (\frac{5}{10} \times 1.47 + \frac{5}{10} \times 21.87) = 8.84 \quad \checkmark$

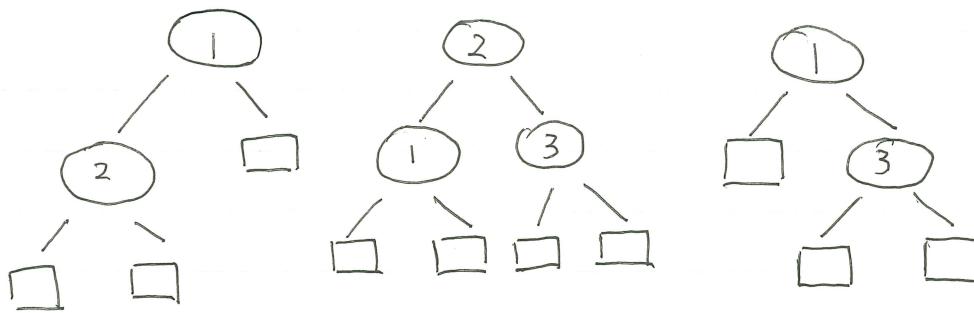
②: $20.51 - (\frac{7}{10} \times 27.80 + \frac{3}{10} \times 1.37) = 0.64$

③: $20.51 - (\frac{4}{10} \times 0.75 + \frac{6}{10} \times 23.32) = 6.22$

Pick the one with highest variance reduction

- Trees are highly sensitive to small changes of data.

Tree ensemble
have a lot of decision trees that do the same thing



Prediction : 1

0

1

Final Prediction : 1.

vote for the final decision.

Sampling with replacement : can be put back

- Random Forest Algorithm:

Given training set of size m

For $b=1$ to B set B larger never hurt performance, but too large slows computation

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

Train a decision tree on the new dataset

Bagged decision tree

Randomizing feature choice, often $k = \sqrt{n}$

For n features available, pick random subset of $k n$ and allow the algorithm to only choose from that subset of features

Campus It causes the algorithm to explore a lot of small changes and average it to reduce impact

Boosted trees intuition.

while applying Random Forest Algorithm instead of picking from all examples with equal ($\frac{1}{m}$) probability, make it more likely to pick misclassified examples from previously trained trees.

Build Tree ensemble on 1, 2 ... b-1
when in the bth, apply boosting.

XGBoost (extreme Gradient Boosting)

- open source
- efficient
- good splitting criteria
- built in regularization
- competitive in competitions such as Kaggle
- assign different ways to different training examples.

from xgboost import XGBClassifier / XGBRegressor

model = XGBClassifier() / XGBRegressor()

model.fit(X_train, y_train)

y-pred = model.predict(X_test)

sklearn.GridSearchCV can refit parameters to best combination

Decision Tree and Tree ensembles

- works well on tabular (structured) data. (can be stored in spread sheet format)
- not recommended for unstructured data (image, audio, text)
- fast.
- small decision trees may be human interpretable.

Neural networks

- works well on all types of data,
- maybe slower.
- works with transfer learning
- when building a system of multiple models working together, it might be easier to string together multiple neural networks

Unsupervised model

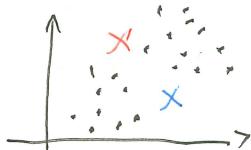
Clustering

K-means

1: Assign points to cluster centroids

2: Move cluster centroids by taking average.

> Repeat



If not data is assigned to a centroid,
delete it.

K-means algorithm.

Randomly initialize k cluster centroid $\mu_1, \mu_2, \dots, \mu_k$

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

$\mu_k :=$ average of points in cluster k .

}

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

K-means always converges

Random initialization:

choose $K < m$ (if $K > m$, there won't be enough data).

Randomly pick K training Examples

Set μ_1, \dots, μ_K equal to these examples

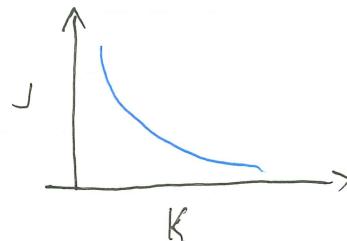
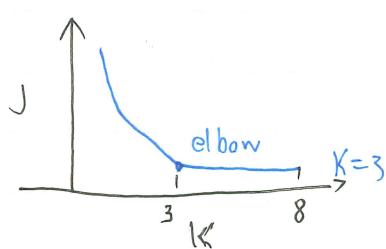
Might end up in local minimum, so compute cost and select.

For $i=1$ to 100

50~1000

Choose the value of K

Elbow method



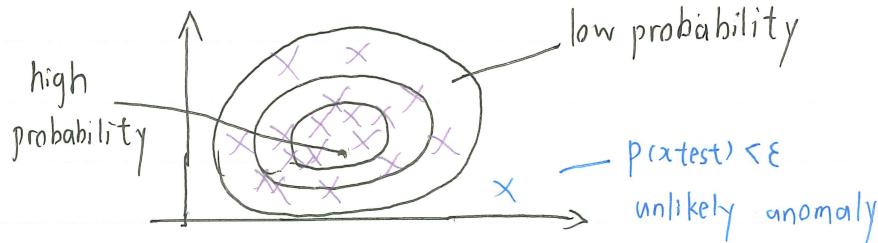
But often
it has no clear
elbow

The right " K " is ambiguous. don't choose K just to minimize cost J

Evaluate k -means based on how well it performs on that later purpose
Based on how well it makes sense in practise.
e.g. image compression

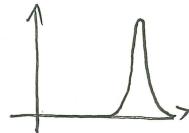
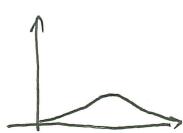
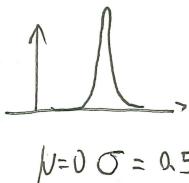
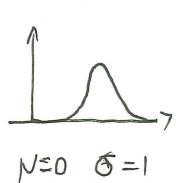
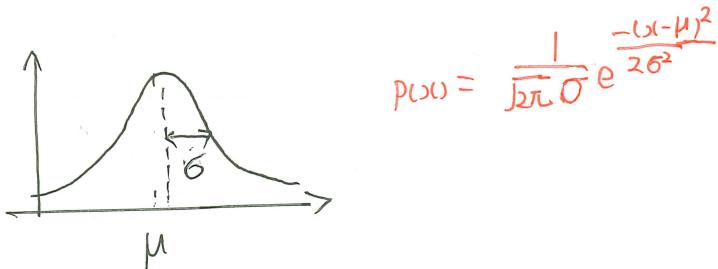
Anomaly detection

Density estimation



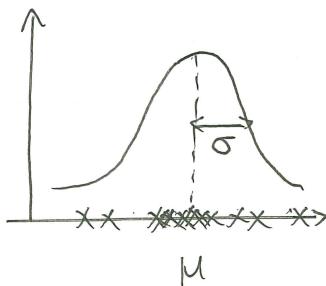
e.g.: Fraud detection - / airplane engine manufacturing

Gaussian(Normal) Distribution



Parameter estimation.

Dataset: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ (x has only one feature)



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

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

maximum likelihood for μ, σ ,
but fine here, no need $\frac{1}{m-1}$

Density estimation.

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

Each example $\vec{x}^{(i)}$ has n features

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

Assume x_1, \dots, x_n are independent, but actually algorithm work even not
 $p(\vec{x}) = p(x_1; \mu_1, \sigma_1^2) * p(x_2; \mu_2, \sigma_2^2) * \dots * p(x_n; \mu_n, \sigma_n^2)$

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

Algorithm.

1. choose n features x_i

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

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

$$\text{(vectorized)} : \vec{\mu} = \frac{1}{m} \sum_{i=1}^m \vec{x}^{(i)} \quad \vec{\sigma}^2 = \begin{bmatrix} \sigma_1^2 \\ \sigma_2^2 \\ \vdots \\ \sigma_n^2 \end{bmatrix}$$

3. Given new examples x_i , compute $p(x)$

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

4. Anomaly if $p(x) < \epsilon$

e.g.

