

Neural Networks

Inference = prediction

A neuron is a computational unit that takes a vector of inputs, computes a weighted sum plus a bias, then applies a non-linear activation function to produce a single output.

For inputs x_1, x_2, \dots, x_n

$$a = g(w_1x_1 + w_2x_2 + \dots + w_nx_n + b)$$

where w are the weights, b is the bias and g is the activation function.

An activation function is a non-linear function applied to the output of a neuron's linear combination. It allows the network to learn non-linear relationships. Common examples are sigmoid, tanh, ReLU, Softmax.

A layer is a collection of neurons that all receive the same input and compute their outputs in parallel. A layer with n neurons produces a vector of n outputs.

Input layer is the layer that receives the raw input features. It does not perform any computation, it simply represents the input.

Hidden layer is any layer between the input and output layers. It is called hidden because its values are not observed in the training data; they are intermediate computations learned by the network.

Output layer is the final layer of the network. It produces the prediction \hat{y} .

A neural network is a function composed of multiple layers of neurons, where the output of each layer becomes the input to the next layer. The function is parameterized by the weights and biases of all layers, which are learned from the data.

A deep neural network is a neural network with more than one hidden layer. The "deep" refers to the depth or number of layers.

Forward propagation is the process of computing the output of a neural network for a given input. It involves applying each layer's transformation in sequence, starting from the input and moving toward the output.

For layer l :

$$z^{(l)} = W^{(l)} a^{(l-1)} + b^{(l)}$$

$$a^{(l)} = g^{(l)}(z^{(l)})$$

Backward propagation (backprop) is the process of computing the gradients of the loss function with respect to every weight and bias in the network.

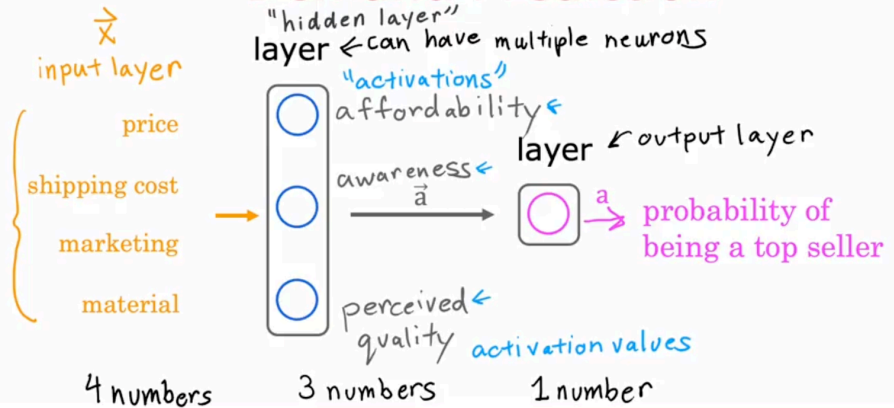
Training is the process of adjusting the parameters of the network to minimize the loss function on a training dataset, typically using gradient descent and back propagation.

Inference is the process of using a trained network to make predictions on new data. Parameters are frozen, and only forward propagation is performed.

Vectorization is the practice of expressing computations as operations on vectors and matrices rather than explicit loops. It enables use of optimized linear algebra libraries and GPUs.

Epoch is one complete pass through the entire training dataset during training.

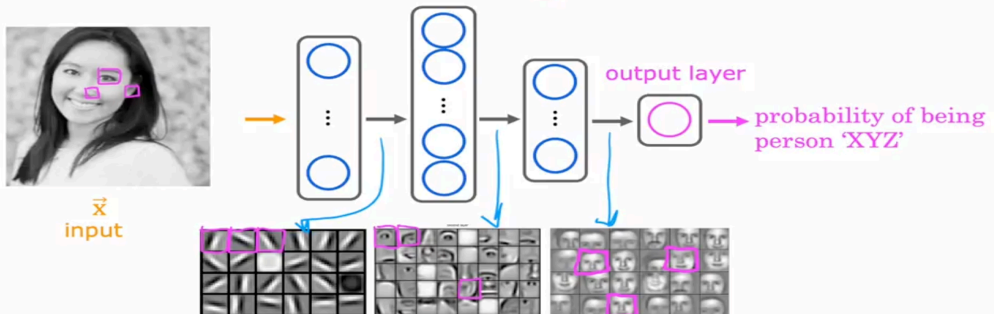
Demand Prediction



It is more like the network is doing feature engineering by combining previous layer to produce the next layer containing features that are more representative of the output.

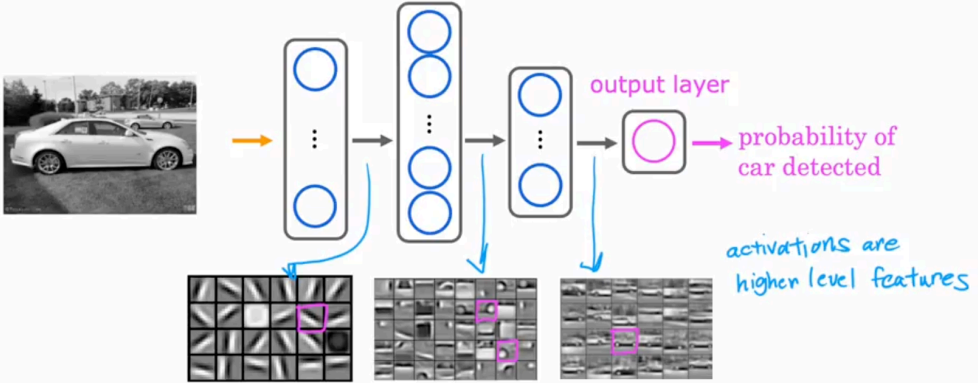
Recognizing images

Face recognition



source: Convolutional Deep Belief Networks for Scalable Unsupervised Learning of Hierarchical Representations by Honglak Lee, Roger Grosse, Ranganath Andrew Y. Ng

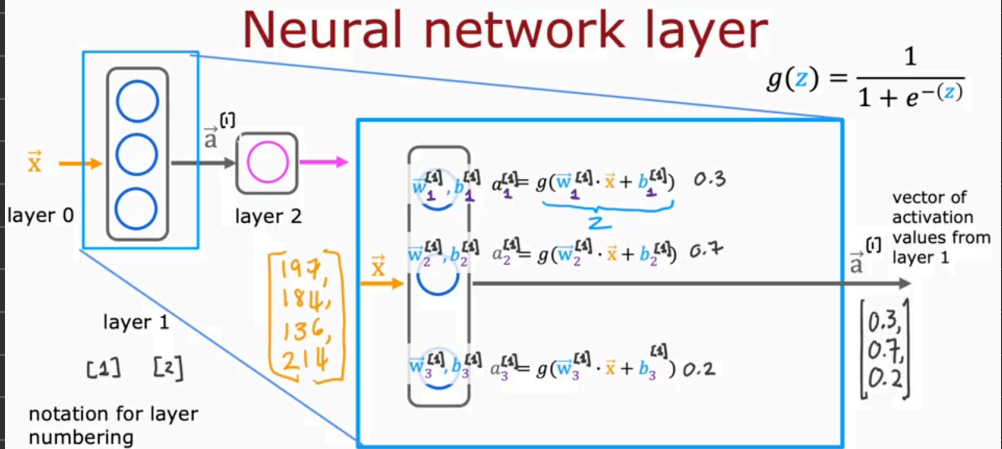
Car classification



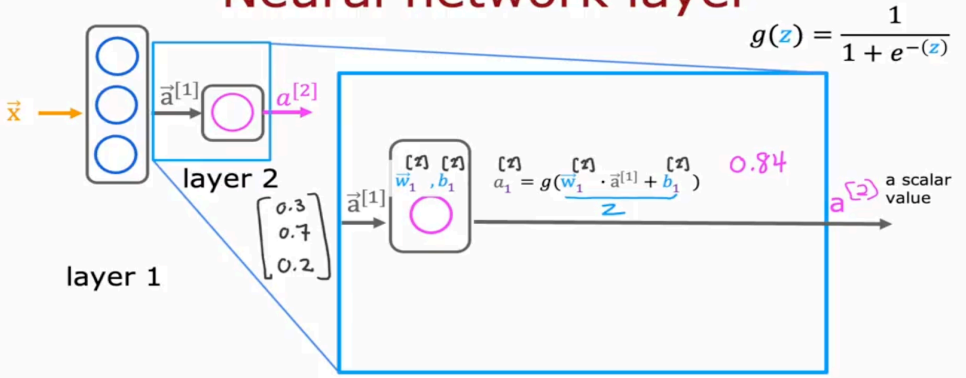
source: Convolutional Deep Belief Networks for Scalable Unsupervised Learning of Hierarchical Representations by Honglak Lee, Roger Grosse, Ranganath Andrew Y. Ng

Nobody taught the network what to look for in each layer, it just learns all of that by itself

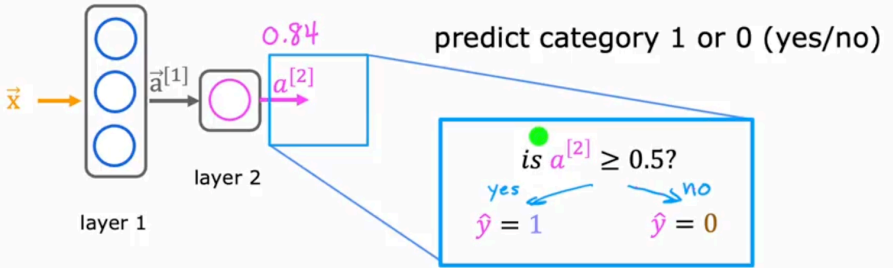
Neural Network Layer:



Neural network layer



Neural network layer



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

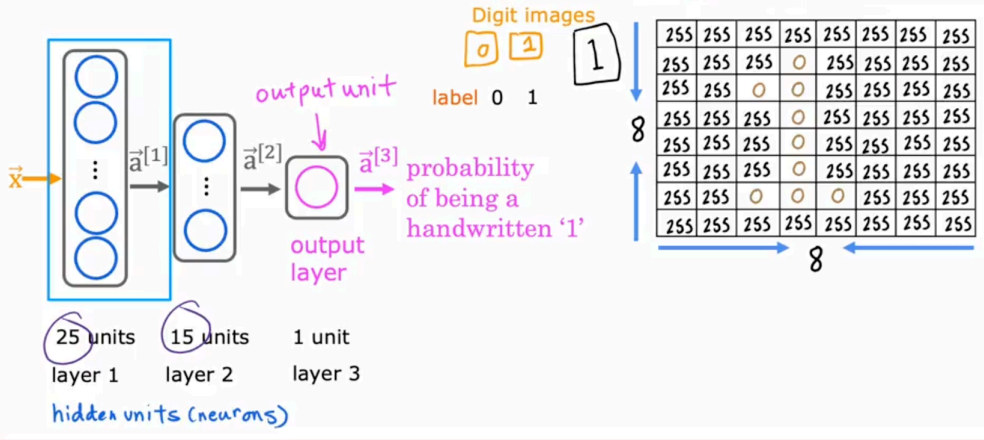
↑
activation function

↑
parameters w and b of layer L , unit j

$$\vec{x} = \vec{a}^{[0]}$$

Inference: making predictions (forward propagation)

Handwritten digit recognition



$$\vec{a}^{[1]} = \begin{bmatrix} g(\vec{w}_1^{[1]} \cdot \vec{x} + b_1^{[1]}) \\ \vdots \\ g(\vec{w}_{25}^{[1]} \cdot \vec{x} + b_{25}^{[1]}) \end{bmatrix}$$

$$\vec{a}^{[2]} = \begin{bmatrix} g(\vec{w}_1^{[2]} \cdot \vec{a}^{[1]} + b_1^{[2]}) \\ \vdots \\ g(\vec{w}_{15}^{[2]} \cdot \vec{a}^{[1]} + b_{15}^{[2]}) \end{bmatrix}$$

$$\vec{a}^{[3]} = \begin{bmatrix} g(\vec{w}_1^{[3]} \cdot \vec{a}^{[2]} + b_1^{[3]}) \end{bmatrix}$$

is $a_1^{[3]} \geq 0.5$?

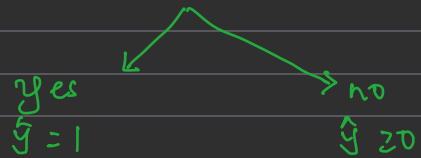


Image is digit 1

Image is not digit 1

Tensorflow

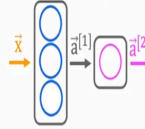
```
import numpy as np
import tensorflow as tf
from tensorflow.keras.layers import Dense
```

```
x = np.array([[200.0, 17.0]])
layer_1 = Dense(units=3,
activation='sigmoid')
a1 = layer_1(x)
```

```
layer_2 = Dense(units=1,
activation='sigmoid')
a2 = layer_2(a1)
```

```
if a2 > 0.5:
    y_hat = 1
else:
    y_hat = 0
print(f"y_hat: {y_hat}")
```

Building a neural network architecture

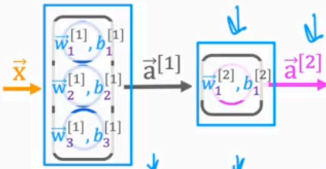


```
model = Sequential([
    Dense(units=3, activation="sigmoid"),
    Dense(units=1, activation="sigmoid")])
```

```
x = np.array([[200.0, 17.0],
              [120.0, 5.0],
              [425.0, 20.0],
              [212.0, 18.0]])
targets = np.array([1,0,0,1])
model.compile(...)
model.fit(x,y)
model.predict(x_new)
```

Forward prop implementation in a single layer

forward prop (coffee roasting model)



```
x = np.array([200, 17])
```

$$a_1^{[1]} = g(\bar{w}_1^{[1]} \cdot \bar{x} + b_1^{[1]})$$

1D arrays

$$a_2^{[1]} = g(\bar{w}_2^{[1]} \cdot \bar{x} + b_2^{[1]})$$

$$a_3^{[1]} = g(\bar{w}_3^{[1]} \cdot \bar{x} + b_3^{[1]})$$

```
w1_1 = np.array([1, 2])
```

```
w1_2 = np.array([-3, 4])
```

```
w1_3 = np.array([5, -6])
```

```
b1_1 = np.array([-1])
```

```
b1_2 = np.array([1])
```

```
b1_3 = np.array([2])
```

```
z1_1 = np.dot(w1_1, x) + b1_1
```

```
z1_2 = np.dot(w1_2, x) + b1_2
```

```
z1_3 = np.dot(w1_3, x) + b1_3
```

```
a1_1 = sigmoid(z1_1)
```

```
a1_2 = sigmoid(z1_2)
```

```
a1_3 = sigmoid(z1_3)
```

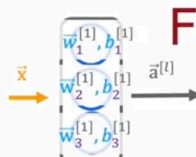
```
a1 = np.array([a1_1, a1_2, a1_3])
```

```
a1^{[2]} = g(\bar{w}_1^{[2]} \cdot \bar{a}^{[1]} + b_1^{[2]})
w2_1 = np.array([-7, 8, 9])
b2_1 = np.array([3])
z2_1 = np.dot(w2_1, a1) + b2_1
a2_1 = sigmoid(z2_1)
```

$w_1^{[2]}$ w_{2-1}

General Implementation of forward propagation.

Forward prop in NumPy



\vec{x} → $\vec{a}^{[l]}$

$\vec{w}_1^{[1]} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$ $\vec{w}_2^{[1]} = \begin{bmatrix} -3 \\ 4 \end{bmatrix}$ $\vec{w}_3^{[1]} = \begin{bmatrix} 5 \\ -6 \end{bmatrix}$

$W = \text{np.array}(\begin{bmatrix} 1, -3, 5 \\ 2, 4, -6 \end{bmatrix})$ 2 by 3

$b_1^{[1]} = -1$ $b_2^{[1]} = 1$ $b_3^{[1]} = 2$

$b = \text{np.array}([-1, 1, 2])$

$\vec{a}^{[0]} = \vec{x}$

$a_in = \text{np.array}([-2, 4])$

```
def dense(a_in,W,b):  
    3 units = W.shape[1] [0,0,0] a[1]  
    a_out = np.zeros(units)  
    for j in range(units): 0,1,2 a[2]  
        w = W[:,j]  
        z = np.dot(w,a_in) + b[j]  
        a_out[j] = g(z)  
    return a_out
```

```
def sequential(x):  
    a1 = dense(x,W1,b1)  
    a2 = dense(a1,W2,b2)  
    a3 = dense(a2,W3,b3)  
    a4 = dense(a3,W4,b4)  
    f_x = a4  
    return f_x
```

Note: g() is defined outside of dense().
(see optional lab for details)

capital W refers to a matrix

Could be done using dot product for matrices and makes it simpler and more efficient.

Vectorization is the practice of expressing computations as operations on vectors and matrices rather than explicit loops. It enables use of optimized linear algebra libraries and GPUs.

How are neural networks implemented efficiently?

$$z = w_1x_1 + w_2x_2 + \dots + w_nx_n + b$$

Non vectorized version:

```
z = 0
for i in range(n):
    z += w[i] * x[i]
z += b
```

This loops n times, doing one multiplication and one addition per iteration

Vectorized version:

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

This computes the dot product as a single optimized operation.

What vectorization buys us?

- * **Speed!** Modern CPU have SIMD (Single instruction, multiple data) units that can perform the same operation on multiple values simultaneously. GPUs take this further with thousands of parallel cores. Vectorized code uses these capabilities, loop code does not.
- * **Cleaner code**
- * **Mathematical clarity!** Linear algebra notation maps to vectorized code. Reading $z = Wx + b$ reads exactly like the math.

Then we can see a pattern:

$$\frac{\partial \ell}{\partial z_k} = (\hat{y} - y) \cdot \prod_{j=k+1}^n w_j \cdot \prod_{j=k}^{n-1} A_j (1 - A_j)$$

and

$$\frac{\partial \ell}{\partial w_k} = \frac{\partial \ell}{\partial z_k} \cdot A_{k-1}$$

$$\frac{\partial \ell}{\partial b} = \frac{\partial \ell}{\partial z_k}$$

We can see the closed form has overlapping terms. Each gradient contains products that earlier gradients also contain. If we always compute from scratch, we will always redo some multiplications many times.

We can compute dz and then reuse it in the next layer.

$$dz_k = (\hat{y} - y) \cdot \prod_{j=k+1}^n w_j \cdot \prod_{j=k}^{n-1} A_j (1 - A_j)$$

$$dz_{k+1} = (\hat{y} - y) \cdot \prod_{j=k+2}^n w_j \cdot \prod_{j=k+1}^{n-1} A_j (1 - A_j)$$

element wise
mul, not dot

$$\frac{dz_k}{dz_{k+1}} = w_{k+1} \cdot A_k (1 - A_k) \Rightarrow dz_k = \left(w_{k+1}^T dz_{k+1} \right) \odot A_k (1 - A_k)$$

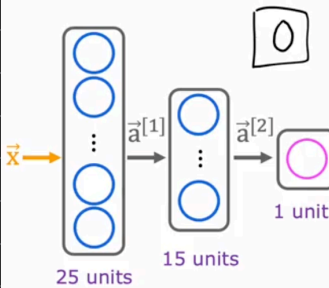
Implementation here:

<https://colab.research.google.com/drive/1hHemUJ3s-m8jnZ4jlawQ2C2XKZbwAUKX?usp=sharing>

Neural Network Training

Train a neural network in TensorFlow

Train a Neural Network in TensorFlow



Given set of (x, y) examples
How to build and train this in code?

```
import tensorflow as tf
from tensorflow.keras import Sequential
from tensorflow.keras.layers import Dense

model = Sequential([
    Dense(units=25, activation='sigmoid'),
    Dense(units=15, activation='sigmoid'),
    Dense(units=1, activation='sigmoid'),
])
from tensorflow.keras.losses import BinaryCrossentropy
model.compile(loss=BinaryCrossentropy())

model.fit(X, Y, epochs=100)
```

①

②

③ epochs: number of steps in gradient descent

Model training Steps

Model Training Steps TensorFlow

- ① specify how to compute output given input x and parameters w, b (define model)

$$f_{\bar{w}, b}(\bar{x}) = ?$$

- ② specify loss and cost

$$L(f_{\bar{w}, b}(\bar{x}), y) \quad \text{1 example}$$

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

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

logistic regression

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

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

logistic loss

$$\text{loss} = -y * \text{np.log}(f_x) - (1-y) * \text{np.log}(1-f_x)$$

$$w = w - \text{alpha} * \text{dj_dw}$$

$$b = b - \text{alpha} * \text{dj_db}$$

neural network

```
model = Sequential([
    Dense(...),
    Dense(...),
    Dense(...)])
```

binary cross entropy

```
model.compile(
    loss=BinaryCrossentropy())
```

```
model.fit(X, y, epochs=100)
```

Training - Definitions and concepts.

Training any neural network follows the same three steps as logistic regression, just at a larger scale:

- * Specify the model: define how to compute the output given input and parameters.
- * Specify the loss and cost function: define what "good" means
- * Train on data to minimize the cost: adjust parameters to reduce loss

```
# Step 1: Specify the model
```

```
model = Sequential([  
    Dense(units=25, activation='sigmoid'),  
    Dense(units=15, activation='sigmoid'),  
    Dense(units=1, activation='sigmoid')  
])
```

```
# Step 2: Specify the loss
```

```
model.compile(loss=BinaryCrossentropy())
```

```
# Step 3: Train
```

```
model.fit(X, y, epochs=100)
```

Loss function: measures how wrong the model is on a single training example.

Cost function: Average of the loss across all training examples.

Activation functions: non-linear function applied to the output of each neuron. Without it, stacking layers would collapse into one linear transformation.

Sigmoid

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

- * Outputs between 0 and 1
- * Used for binary classification at the output layer
- * Suffers from vanishing gradients in deep networks.
- * Rarely used in hidden layers today.

ReLU (Rectified Linear Unit)

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

- * Outputs 0 for negative inputs, the input itself for positive
- * Default choice for hidden layer
- * Much faster to compute than sigmoid
- * Does not suffer from vanishing gradients
- * Can have "dead neurons" if many inputs are negative.

Linear (no activation)

$$g(z) = z$$

- * Used for regression outputs that can be any real number.
- * Equivalent to no activation function.

Softmax

$$g(z)_k = \frac{e^{z_k}}{\sum_j e^{z_j}}$$

- * Used for multi-class classification at the output layer.
- * Converts a vector of scores into probabilities summing to 1.

Choosing the Output activation

The output activation depends on the type of problem:

Problem type	Output activation
Binary classification	Sigmoid
Multiclass classification	Softmax
Regression (any real number)	Linear
Regression (non-negative)	ReLU

Choosing hidden layer activations

For hidden layers, the default is ReLU. Reasons are:

- * Faster to compute than sigmoid (just a comparison vs exponential)
- * Doesn't saturate on the positive side, so gradients don't vanish
- * In practice, networks with ReLU hidden layers train faster and reach better accuracy.

Why we need activation functions

Without non-linear activation, a deep network reduces to a single linear transformation:

$$Z_2 = W_2 Z_1 + b_2 = W_2 (W_1 X + b_1) + b_2 = (W_2 W_1) X + (W_2 b_1 + b_2)$$

Which is just another linear function, a single layer can represent it. The network gains no expressive power from depth.

Non-linear activations break this collapse. They let the network learn arbitrarily complex functions by composing many simple non-linear pieces.

Multiclass classification

When we have more than two classes, we switch from sigmoid + binary cross-entropy to softmax + categorical cross-entropy.

Softmax converts logits into a probability distribution over K classes:

$$P(y=k|x) = \frac{e^{z_k}}{\sum_{j=1}^K e^{z_j}}$$

Categorical cross-entropy measures how far the predicted distribution is from the true (one-hot) distribution:

$$L = - \sum_{k=1}^K y_k \log(\hat{y}_k)$$

```
model = Sequential([
    Dense(units=25, activation='relu'),
    Dense(units=15, activation='relu'),
    Dense(units=10, activation='linear') # 10 classes
])

model.compile(loss=SparseCategoricalCrossentropy(from_logits=True))
```

Note: using linear activation at the output and `from_logits = True` in the loss is numerically more stable than computing softmax separately. The loss function handles both internally.

Less stable (computes softmax then log separately)

```
Dense(units=10, activation='softmax')
loss=SparseCategoricalCrossentropy()
```

The output of the network is now in "logit" form (raw scores), and the loss handles conversions to probabilities internally. To get probabilities at prediction time, apply softmax manually:

```
logits = model(X)
probabilities = tf.nn.softmax(logits)
```

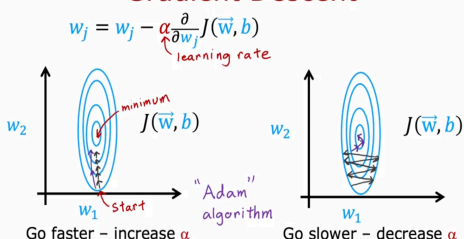
Adam Optimizer (Adaptive moment estimation)
 Instead of plain gradient descent with a fixed learning rate, Adam automatically adapts the learning rate per parameter based on past gradients.

```
model.compile(
    optimizer=Adam(learning_rate=0.001),
    loss=BinaryCrossentropy()
)
```

Why Adam?

- * Adapts the learning rate for each weight individually
- * Larger learning rate for weights with consistent gradient direction
- * Smaller learning rate for noisy or oscillating gradients.
- * Generally converges faster and more reliably than vanilla SGD

Gradient Descent



MNIST Adam

```
model = Sequential([
    tf.keras.layers.Dense(units=25, activation='sigmoid'),
    tf.keras.layers.Dense(units=15, activation='sigmoid'),
    tf.keras.layers.Dense(units=10, activation='linear')
])

compile
model.compile(optimizer=tf.keras.optimizers.Adam(learning_rate=1e-3),
              loss=tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True))

fit
model.fit(X, Y, epochs=100)
```


Backpropagation

Computing the gradient of the cost with respect to every weight in the network. We use chain rule from calculus to propagate the error backward from the output to the input.

$$dz^{[L]} = \left(W^{[L+1]T} dz^{[L+1]} \right) \odot g'^{[L]}(z^{[L]})$$

$$dW^{[L]} = \frac{1}{m} dz^{[L]} (A^{[L+1]})^T$$

$$db^{[L]} = \frac{1}{m} \sum dz^{[L]}$$

if we have N nodes and P parameters, backprop computes derivatives in $N+P$ steps rather than $N \times P$ steps.

Weight Initialization

How we set the initial values of weights matters enormously.

If we use zeros, all neurons compute the same thing, so they all update identically.

If we use $W = np.random.randn(\dots) * 0.01$. This works for shallow networks but can cause vanishing or exploding gradients in deeper ones.

Xavier initialization (for sigmoid/tanh)

$$W = \mathcal{N}(0,1) \times \sqrt{\frac{1}{n^{[L-1]}}}$$

He initialization (Best for ReLU)

$$W = \mathcal{N}(0,1) \times \sqrt{\frac{2}{n^{[L-1]}}}$$

Building a binary classifier

```
import tensorflow as tf
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense
from tensorflow.keras.losses import BinaryCrossentropy
from tensorflow.keras.optimizers import Adam

# 1. Specify the model
model = Sequential([
    Dense(units=25, activation='relu'),
    Dense(units=15, activation='relu'),
    Dense(units=1, activation='sigmoid')
])

# 2. Specify the loss
model.compile(
    optimizer=Adam(learning_rate=0.001),
    loss=BinaryCrossentropy(),
    metrics=['accuracy']
)

# 3. Train
model.fit(X_train, y_train, epochs=100,
          validation_data=(X_val, y_val))

# 4. Predict
predictions = model.predict(X_test)
```

Multilabel Classification

In multi-label classification, each example can belong to multiple classes at the same time. Each label is independent, we predict yes or no for each one.

Contrast this with multi-class classification, where each example belongs to exactly one of k classes.

- * Multi-class uses softmax because outputs must sum to 1, picking one class means rejecting the others.
- * Multi-label uses sigmoid on each output independently, saying yes to cat doesn't affect whether we say yes to dog.

For k labels, the output layer has k neurons, each with Sigmoid activation:

$$\hat{y}_k = \sigma(z_k) \quad \text{for } k=1, 2, \dots, k$$

Each output is interpreted as the probability that label k is present, independent of others.

```
model = Sequential([
    Dense(units=64, activation='relu'),
    Dense(units=32, activation='relu'),
    Dense(units=K, activation='sigmoid') # K sigmoid
    outputs, not softmax
])
```

Binary cross-entropy is applied independently to each label then averaged

$$L = -\frac{1}{k} \sum_{k=1}^k [y_k \log(\hat{y}_k) + (1-y_k) \log(1-\hat{y}_k)]$$

```

model.compile(
    optimizer=Adam(learning_rate=0.001),
    loss=BinaryCrossentropy(),
    metrics=['accuracy']
)

```

The labels are multi-hot encoded instead of one-hot encoded for multiclass.

Additional Layer Types:

In dense layers, each neuron output is a function of all the activation outputs of the previous layer.

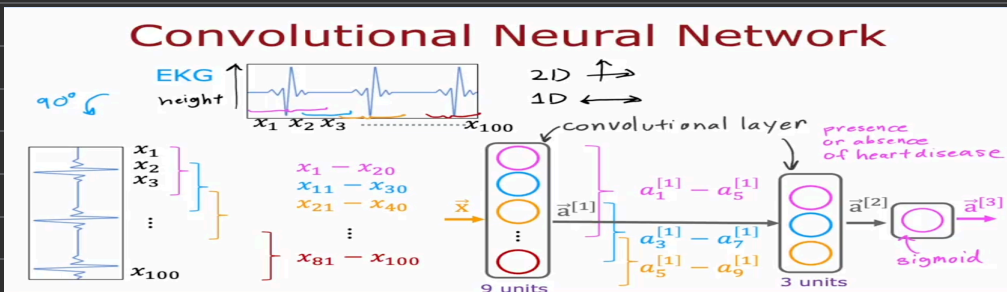
$$\vec{a}_i^{[2]} = g(\vec{w}_i^{[2]} \vec{a}^{[1]} + b^{[2]})$$

Convolutional layers, each neuron only looks at part of the previous layer's output.

But why would we want to do this?

- * Faster computation
- * Need less training data (less prone to overfitting)

If we have multiple convolutional layers in a neural network we call that network a convolutional neural network (CNN).



Advice for Applying Machine Learning.

When the model performs poorly, what should we actually do?
There are many possible next steps:

- * Collecting more training data.
- * Try a smaller set of features.
- * Add more features.
- * Add polynomial features.
- * Decrease the regularization parameter λ
- * Increase λ
- * Make the model bigger
- * Train longer.

Train/Validation/test split:

We can split the data into three sets:

- * Training Set - the data the model learns from
- * Validation set (also called dev set or cross-validation set) - used to tune hyperparameters and pick between models.
- * Test set - Used only once at the very end to estimate true generalization.

Typical splits:

Dataset size	Train	Val	Test
Small (~10k)	60%	20%	20%
Large (millions)	98%	1%	1%

We use three sets and not two because if we use the test set to pick hyperparameters, we have effectively trained on it. It's accuracy is no longer an honest measurement of generalization. The validation set is the data we are allowed to look at when making decisions.

Bias and Variance

* Bias is how much the model underfits

* Variance is how much the model overfits

High variance means the model has memorized training-set noise and fails on new data. It overfits.

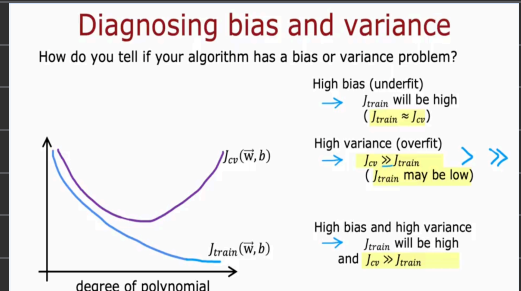
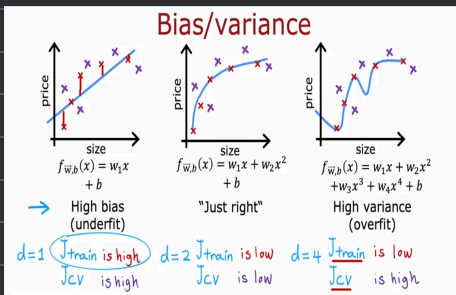
We diagnose these by comparing training error and validation error.

Training error	Validation error	Diagnosis
High	High	High bias (underfitting)
Low	High	High variance (overfitting)
Low	Low	Good fit
High	Even higher	High bias and high variance

Baseline on performance! We can't judge if training error is high or low without context. We compare against:

- * Human-level performance - how well do people do this task?
- * Simple baseline - What does a basic model achieve
- * Best published performance - What's the state of the art (SOTA)?

For example, 10% error on speech recognition sounds bad, but if humans also get 10% on the same audio (because of background noise), the model is fine. The gap that matters is between the model and achievable best.



Diagnostic procedure

Given training error, validation error and a baseline

$$\text{gap_bias} = \text{training_error} - \text{baseline}$$

$$\text{gap_variance} = \text{validation_error} - \text{training_error}$$

- * Large gap_bias \rightarrow underfitting
- * Large gap_variance \rightarrow overfitting

Both can happen at once.

Fixing high bias (underfitting)

- * Add more features
- * Add polynomial features
- * Use a more complex model (more layers, more units)
- * Decrease the regularization parameter λ
- * Train longer.

Note that adding more training data does not fix high bias. The model isn't expressive enough to fit even what it has.

Fixing high variance (overfitting)

- * Get more training data
- * Reduce the number of features
- * Increase the regularization parameter
- * Use a simpler model
- * Use dropout (for neural networks)
- * Use early stopping

Quick reference table

Symptom	Diagnosis	What to try
High train error, high val error	High bias	Bigger model, more features, less regularization
Low train error, high val error	High variance	More data, fewer features, more regularization
Both low	Healthy	Done, possibly try a smaller model for efficiency

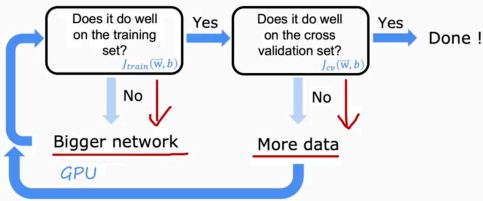
Bias/Variance and neural networks

For classical ML, bias and variance are in tension. We trade one against the other by adjusting model complexity.

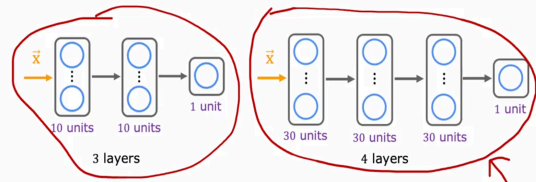
Neural networks change this trade-off, with enough data and proper regularization, a large network can reduce both bias and variance simultaneously

Neural networks and bias variance

Large neural networks are low bias machines



Neural networks and regularization



A large neural network will usually do as well or better than a smaller one so long as regularization is chosen appropriately.

Neural network regularization

$$J(\mathbf{W}, \mathbf{B}) = \frac{1}{m} \sum_{i=1}^m L(f(\vec{x}^{(i)}), y^{(i)}) + \frac{\lambda}{2m} \sum_{\text{all weights } \mathbf{W}} (w^2)$$

Unregularized MNIST model

```
layer_1 = Dense(units=25, activation="relu")
layer_2 = Dense(units=15, activation="relu")
layer_3 = Dense(units=1, activation="sigmoid")
model = Sequential([layer_1, layer_2, layer_3])
```

Regularized MNIST model

```
layer_1 = Dense(units=25, activation="relu", kernel_regularizer=L2(0.01))
layer_2 = Dense(units=15, activation="relu", kernel_regularizer=L2(0.01))
layer_3 = Dense(units=1, activation="sigmoid", kernel_regularizer=L2(0.01))
model = Sequential([layer_1, layer_2, layer_3])
```

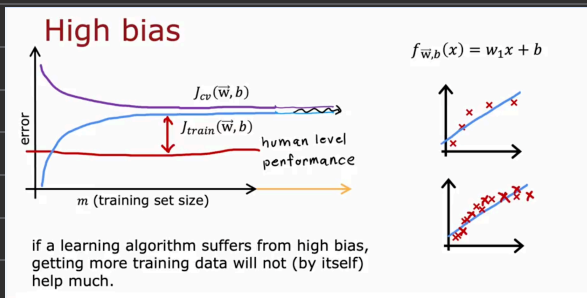
Learning Curves

One downside is that it is computationally expensive to train so many different models using different size subset of training set. It isn't done that much in practice.

We can plot the training error and validation error as a function of training set size.

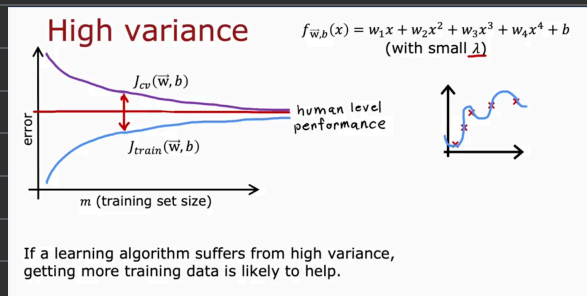
High bias pattern:

- * Both curves plateau at a high error
- * The gap between them is small
- * Adding more data won't help, both curves are already flat.



High variance pattern:

- * Training error is low and stays low
- * Validation error decreases as you add data.
- * The gap closes as data grows, more data will help.



The curves are diagnostic and we should run them before deciding to collect more data.

Regularization

Adding a penalty to large weights to discourage overfitting.
Cost function with L2 regularization

$$J(w, b) = \frac{1}{m} \sum L(y, \hat{y}) + \frac{\lambda}{2m} \sum w^2$$

Tuning λ :

- * $\lambda = 0 \rightarrow$ no regularization. Model can overfit freely. High variance.
- * λ small \rightarrow slight regularization. Doesn't restrict the model much.
- * λ moderate \rightarrow Sweet spot. Model fits data without overfitting noise.
- * λ very large \rightarrow Over regularized. All weights forced near zero. High bias.

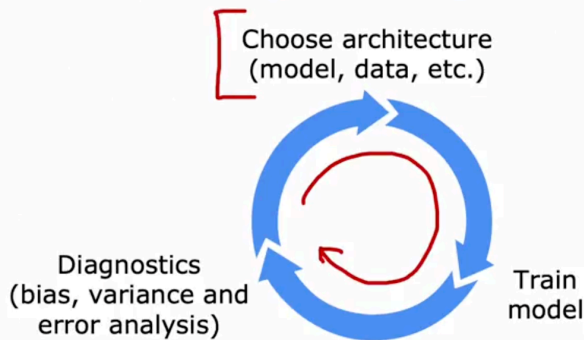
Iterative loop of ML development.

The whole development process is iterative!

- * Choose architecture (model type, features, hyper parameters)
- * Train the model
- * Diagnose errors (high bias? high variance? specific failure modes?)
- * Use diagnostics to decide what to change.
- * Go back to the first step.

We don't get everything right on the first try. The bias/variance make iteration efficient.

Iterative loop of ML development



Error Analysis

After training, we can manually look at examples the model gets wrong. We can categorize them, for instance, if we are building a spam filter and we have 500 misclassified examples in the validation set:

- * 50 are spam but classified as not spam (false negatives)
- * 200 are not spam but classified as spam (false positives)

Then we can look at false positives:

- * 30% are pharmaceutical ads
- * 20% are stolen passwords/phishing
- * 15% are unusual punctuation patterns
- * 10% are deliberately misspelled words.

We can see phishing detection is more impactful than misspellings. Then we know the problem and we can try fixing the biggest categories.

Error analysis is hard for tasks humans are not good at, for example when we try to predict clicks on a website.

Adding more data - Selectively:

If error analysis reveals one category dominates errors, we can target data collection at that category instead of just collecting more of everything.

Data augmentation (for image and audio) creates new training examples from existing ones:

- * Rotate, crop or flip images
- * Add noise to audio
- * Apply distortions that don't change the label

Data synthesis uses artificial data inputs to create new training data (synthetic data).

Let's say for OCR we need synthetic data, we can use different fonts to create random texts and screenshot it using different colors, contrast and fonts and we get synthetic data like that:

Artificial data synthesis for photo OCR



Real data



Synthetic data

[Adam Coates and Tao Wang]

Data synthesis is mostly used for computer vision tasks.

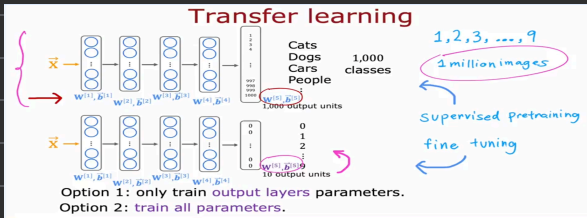
This is much cheaper than collecting genuinely new labeled data.

Transfer Learning

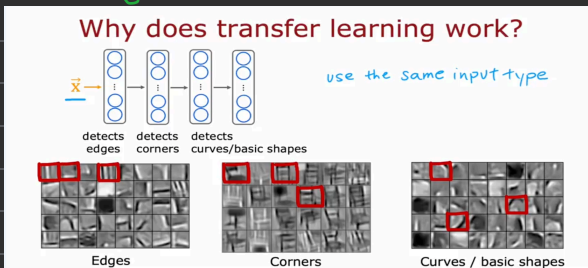
When we have a small dataset, we can train on a related large dataset first (supervised pretraining) then fine-tune on our smaller dataset.

Two approaches!

- * Use a pretrained model as a feature extractor. We replace just the output layer. We train only the new output layer's weights. Useful when we have very little data.
- * Fine-tune the whole network. We initialize with pretrained weights, then train all weights on our data with a small learning rate. Useful when we have more data.



Why this works: Lower layer learns general features (edges in image, common words in text) that transfer across tasks. Higher layer learns task-specific things. By starting with good lower layers, we only need data for the task-specific learning.

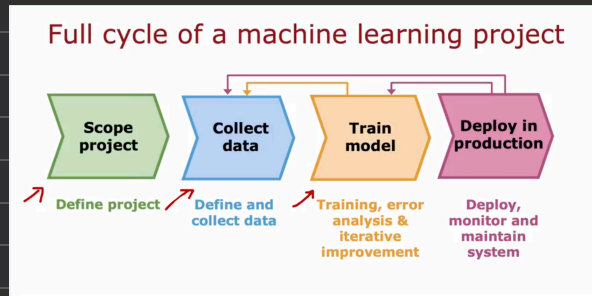


this is one of the most impactful techniques in practical deep learning.

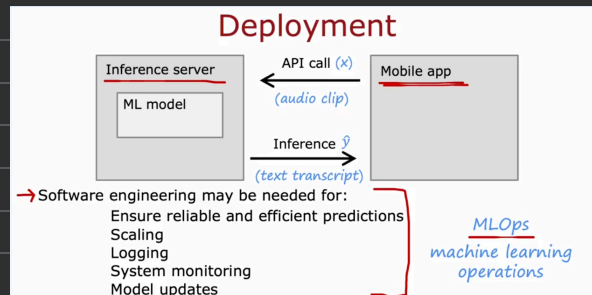
Full ML workflow:

- * Scope the project. What's the problem? What does success look like?
- * Collect data. Gather and label training data.
- * Train a model. Pick architecture, train, evaluate.
- * Deploy and monitor: Get the model into production, watch for drift and failures.

We iterate between these. We can go back to collect more data after a model fails in production.



Deployment



Fairness, Bias and Ethics -

Bias

Hiring tool that discriminates against women.

Facial recognition system matching dark skinned individuals to criminal mugshots.

Biased bank loan approvals.

Toxic effect of reinforcing negative stereotypes.

Adverse use cases

Deepfakes

Spreading toxic/incendiary speech through optimizing for engagement.

Generating fake content for commercial or political purposes.

Using ML to build harmful products, commit fraud etc.

Spam vs anti-spam : fraud vs anti-fraud.

Guidelines

Get a diverse team to brainstorm things that might go wrong, with emphasis on possible harm to vulnerable groups.

Carry out literature search on standards/guidelines for your industry.

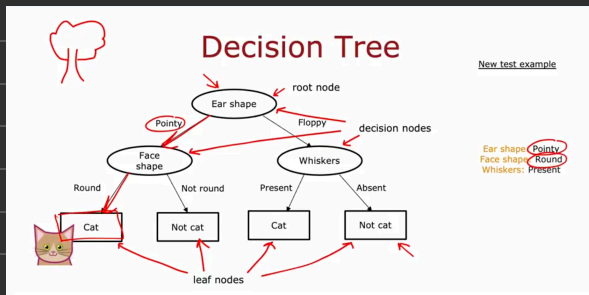
Audit systems against possible harm prior to deployment.



Develop mitigation plan (if applicable), and after deployment, monitor for possible harm.

Decision Trees

A decision tree is a model that makes prediction by asking a series of yes/no questions about the input feature. Each internal node tests one feature; each branch is a possible answer; each leaf is a final prediction.



Learning process:

the tree is built greedily - top down, one split at a time.

1. Start at the root with all training examples
2. Pick the best feature to split on (to max purity or min impurity)
3. Partition the data based on that feature
4. Repeat the same process on each child node.
5. Stop when a criterion is met (pure node, max depth, etc).

The key question is which feature should we split on at each node?

The answer comes down to pick the split that produces the purest child nodes.

When a node is 100% one class

When splitting a node will result in 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.

Measuring Purity (Entropy)

A node is pure when all its examples belong to the same class. A node is impure when classes are mixed.

Entropy measures impurity:

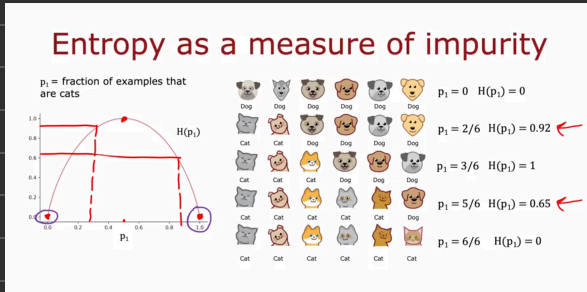
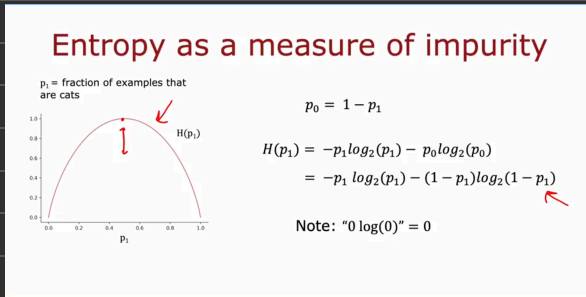
$$H(p_1) = -p_1 \log_2(p_1) - (1-p_1) \log_2(1-p_1)$$

Where p_1 is the fraction of positive examples in the node.

Key properties:

- * $H(0.5) = 1$, maximum entropy, 50/50 split, most impure
- * $H(0) = 0$ and $H(1) = 0$, all examples in one class, pure, zero entropy
- * The curve looks like an inverted U peaking at $p_1 = 0.5$

By convention
 $0 \log_2 0 = 0$



Choosing a split: Information gain

When we split a node, we reduce entropy. The amount of reduction is information gain.

For a split that produces a node with weights w^{left} and a right node with weight w^{right} :

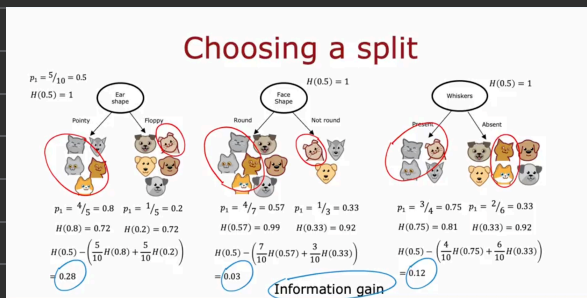
$$\text{Information Gain} = H(p_i^{\text{node}}) - [w^{\text{left}} H(p_i^{\text{left}}) + w^{\text{right}} H(p_i^{\text{right}})]$$

Where:

- * p_i^{root} is the fraction of positive examples at the parent node.
- * w^{left} and w^{right} are the fractions of examples going to each child
- * p_i^{left} and p_i^{right} are the fraction of positives in each child

The split with the highest positive gain is chosen. This is the splitting criterion.

Why did we compute reduction in entropy instead of entropy? The stopping criteria for deciding when to not bother to split any further is if the reduction in entropy is too small. (one of the criteria)



Stopping criteria

Without a stopping rule, the tree keeps splitting until every node is pure, which overfits.

Common Stopping conditions!

- * A node is fully pure (all examples in one class)
- * Reaching a maximum depth
- * The information gain from any further split is below a threshold
- * The number of examples in a node is below a minimum.

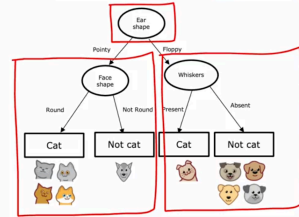
These act as regularization, they prevent the tree from memorizing training data.

Whole learning process =

Decision Tree Learning

- Start with all examples at the root node
- Calculate information gain for all possible features, and pick the one with the highest information gain
- 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

Recursive splitting



Recursive algorithm

`build_tree(node):`

if stopping condition is met:
return leaf with majority class

best_feature = feature with highest information gain
split the node on best_feature

for each child node:
build_tree(child)

Recursive
Algorithm

One-hot encoding for categorical features

If a feature has more than two possible values (e.g. ear shape: pointy, floppy, oval), the tree can't split into 3 branches. We can convert it to multiple binary features instead.

Features with three possible values

Ear shape (x_1)	Face shape (x_2)	Whiskers (x_3)	Cat (y)
Pointy ✓	Round	Present	1
Oval	Not round	Present	1
Oval ✓	Round	Absent	0
Pointy	Not round	Present	0
Oval	Round	Present	1
Pointy	Round	Absent	1
Floppy ✓	Not round	Absent	0
Oval	Round	Absent	1
Floppy	Round	Absent	0
Floppy	Round	Absent	0

3 possible values



One hot encoding →

One hot encoding

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

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

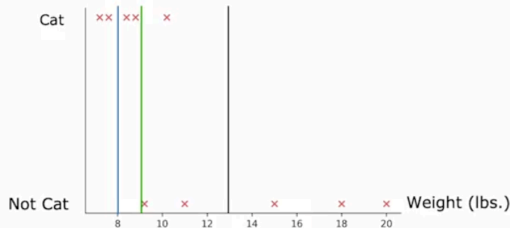
Continuous valued features

For continuous features (e.g. weight = 8.5 kg), the tree tries different threshold values and picks the one with highest information gain. The split becomes $\text{weight} < \text{threshold} \rightarrow \text{yes/no}$.

Procedure:

- 1 Sort examples by the feature value
- 2 Try each midpoint between consecutive distinct values as a candidate threshold
- 3 Compute information gain for each
- 4 Pick the threshold with the best gain.

Splitting on a continuous variable



$$H(0.5) - \left(\frac{2}{10} H\left(\frac{2}{2}\right) + \frac{8}{10} H\left(\frac{3}{8}\right) \right) = 0.24$$

$$H(0.5) - \left(\frac{4}{10} H\left(\frac{4}{4}\right) + \frac{6}{10} H\left(\frac{1}{6}\right) \right) = 0.61$$

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

Regression Trees

Decision trees can predict continuous values too and not just classes. The differences:

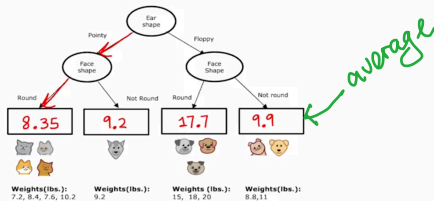
- * Each leaf predicts the average of the target values of examples in that leaf.
- * The splitting criterion is the variance reduction instead of information gain.

For a split:

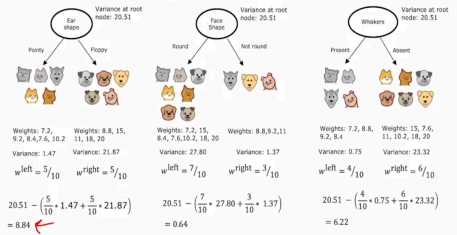
$$\text{Variance Reduction} = \text{Var}(\text{parent}) - [w^{\text{left}} \cdot \text{Var}(\text{left}) + w^{\text{right}} \cdot \text{Var}(\text{right})]$$

The split that reduces the variance the most is chosen.

Regression with Decision Trees



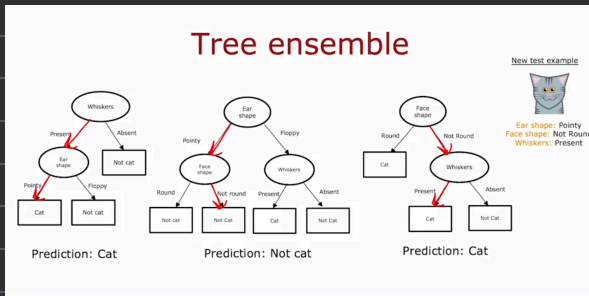
Choosing a split



Tree ensembles

Decision trees are sensitive, small changes in training data can produce very different trees. This high variance is a problem.

Tree ensembles combine many trees to reduce variance. It is basically a collection of trees



We vote and majority carries the vote.

Two main approaches:

- * Random forest! bagging + random feature selection
- * Boosted trees (XGBoost): sequential trees focused on misclassified examples.

Sampling with replacement (bootstrap)

To build an ensemble, we need different trees. We create different training sets by sampling with replacement from the original data, pick a random example, put it back, pick again. Some examples might be picked multiple times, some won't be picked at all.

Each tree in the ensemble gets its own bootstrap sample. This gives the tree diverse training data without needing more data overall.

Random Forest

A random forest is an ensemble of decision tree with two sources of randomness:

- * Bootstrap sampling! Each tree trained on a different bootstrap sample.
- * Random feature! At each node, we only consider a random subset of features (typically \sqrt{n} for classification)

For prediction:

- * Classification \rightarrow majority vote across trees
- * Regression \rightarrow average of trees predictions.

The random feature subset is what distinguishes random from plain bagging. It forces trees to be more diverse, reducing correlation between them. More diversity \rightarrow better variance reduction when averaging.

Bagging

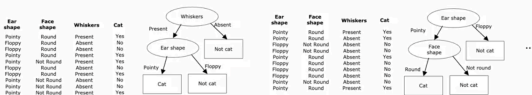
Generating a tree sample

Given training set of size m

For $b = 1$ to $B \Rightarrow 64 - 123$ for 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

Boosted Trees (XGBoost)

Unlike random forest where trees are trained independently in parallel, boosting trains trees sequentially. Each new tree focuses on the example that previous trees got wrong.

Process:

- * Train a tree on the data
- * Identify examples the current ensemble misclassifies
- * Train the next tree with more weight on those misclassified examples.
- * Add new tree to the ensemble
- * Repeat.

XGBoost (Extreme Gradient Boosting) is the most popular implementation. It's fast, regularized and dominates many ML competitions on tabular data.

```
from xgboost import XGBClassifier
model = XGBClassifier()
model.fit(X_train, y_train)
predictions = model.predict(X_test)
```

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

Ear shape	Face shape	Whiskers	Cat
Floppy	Round	Present	Yes
Floppy	Round	Absent	No
Floppy	Round	Absent	No
Floppy	Round	Present	Yes
Floppy	Not Round	Present	No
Floppy	Round	Absent	Yes
Floppy	Round	Present	Yes
Floppy	Not Round	Absent	No
Floppy	Not Round	Absent	No
Floppy	Not Round	Present	Yes



Ear shape	Face shape	Whiskers	Prediction
Floppy	Round	Present	Cat (B)
Floppy	Round	Absent	Not cat (X)
Floppy	Round	Present	Cat (B)
Floppy	Round	Absent	Not cat (X)
Floppy	Not Round	Present	Cat (B)
Floppy	Not Round	Absent	Not cat (X)
Floppy	Round	Absent	Not cat (X)
Floppy	Round	Present	Cat (B)
Floppy	Not Round	Absent	Not cat (X)
Floppy	Not Round	Present	Cat (B)

1, 2, ..., b-1

Using XGBoost

Classification

```
from xgboost import XGBClassifier
model = XGBClassifier()
model.fit(X_train, y_train)
y_pred = model.predict(X_test)
```

Regression

```
from xgboost import XGBRegressor
model = XGBRegressor()
model.fit(X_train, y_train)
y_pred = model.predict(X_test)
```

When to use decision trees?

Decision trees (and ensembles) work well when:

- * Data is tabular / structured (rows and columns)
- * We need interpretability
- * Training time matters (they train faster)
- * Small to moderate dataset sizes
- * Mixed feature types (numerical + categorical)

Neural networks

- * Works on all sorts of data whether structured or not
- * 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