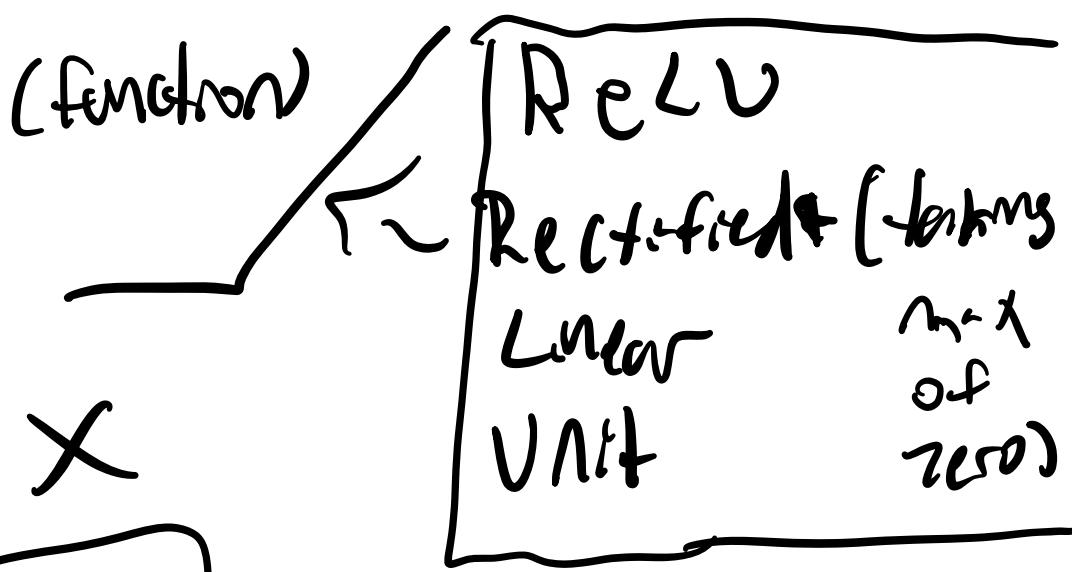


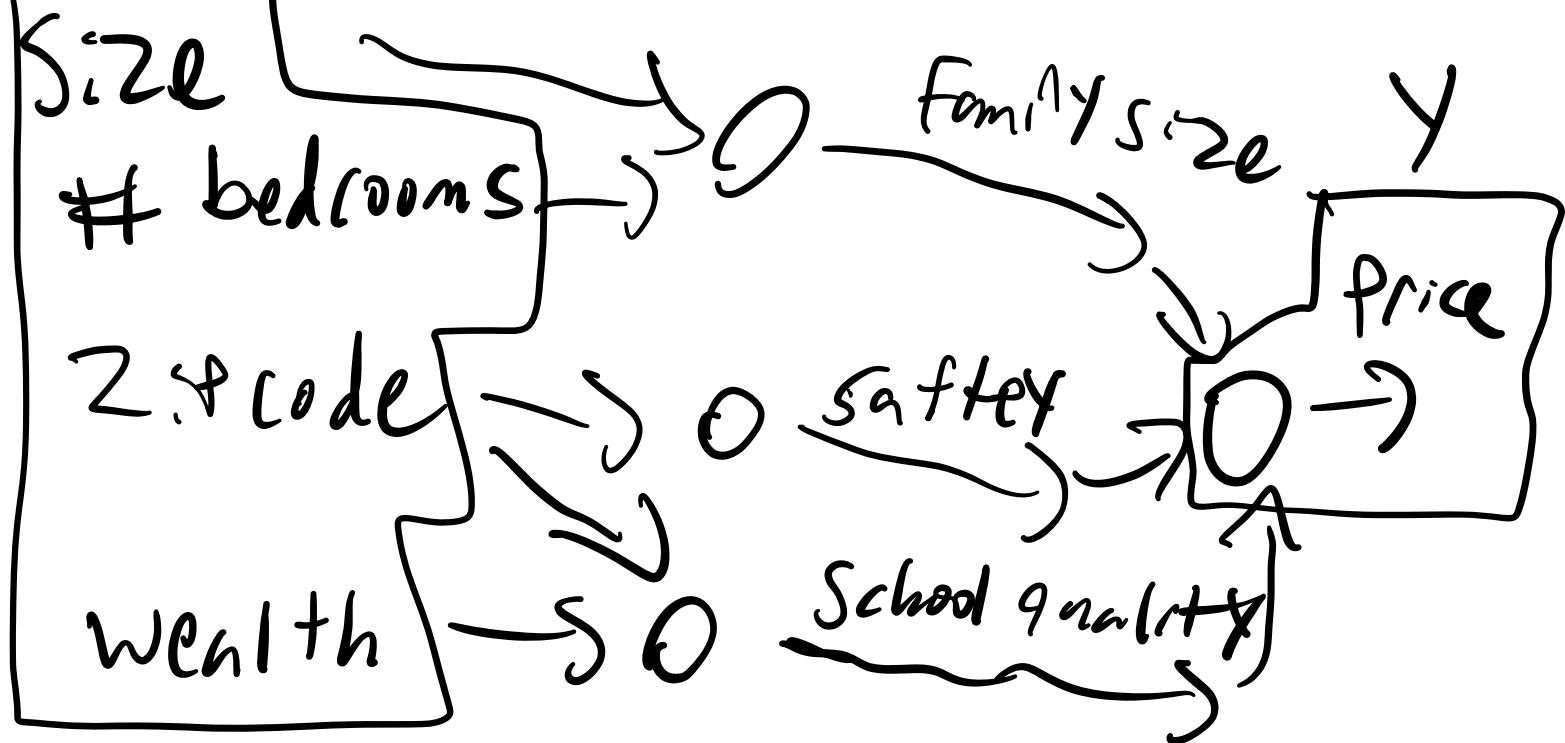
Neural Networks



$\text{size} \rightarrow 0 \rightarrow \text{Price}$



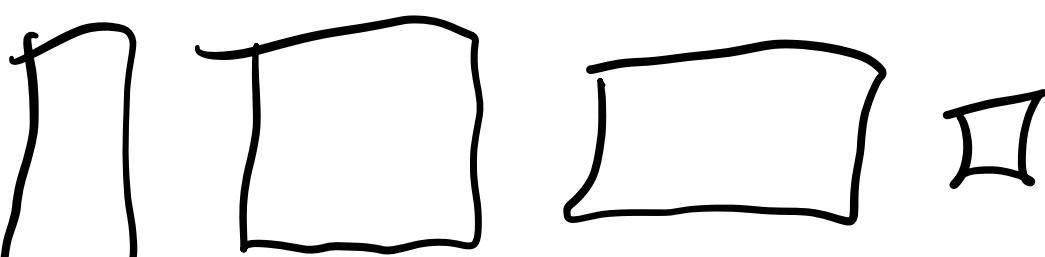
"Neuron"
"inputs are
computes
linear
function
* takes max
of zero
* outputs
price"



↑
Neural Network
will decide what
this will be, all you
need X and Y

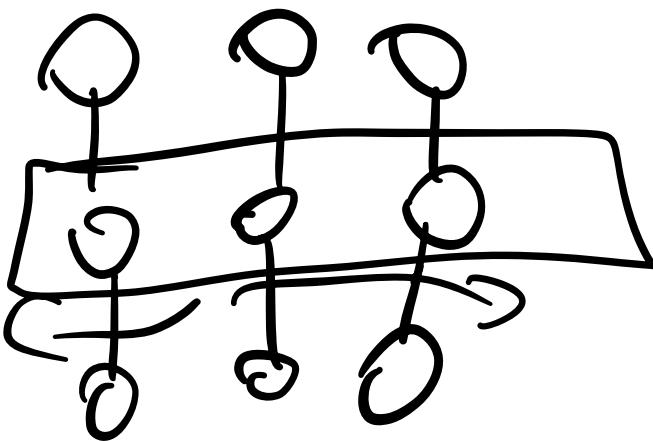
Supervised learning

Application
Standard NN → Real estate, online ads
CNN → Photo tagging
RNN → Speech recognition, Machine translation
Custom hybrid → Autonomous driving

CNN → 
Image data

~~using~~
RNN's

on 1 dimensional
sequence
data

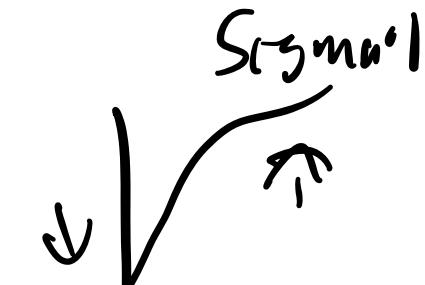


Structured data vs Unstructured data

well defined
features
such as
a number
representing
user age

Audio and
images
or random
text
depending
on more
recognition
pixels,
individual words,

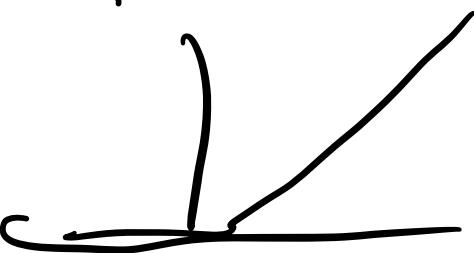
Scale drives deep learning progress *



Slow
learning at
slopes



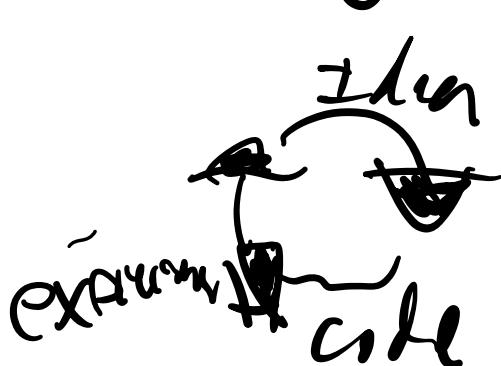
ReLU



Gradient is equal to 1
so it's faster
since it doesn't go
to 0

Data
computation

Algorithms

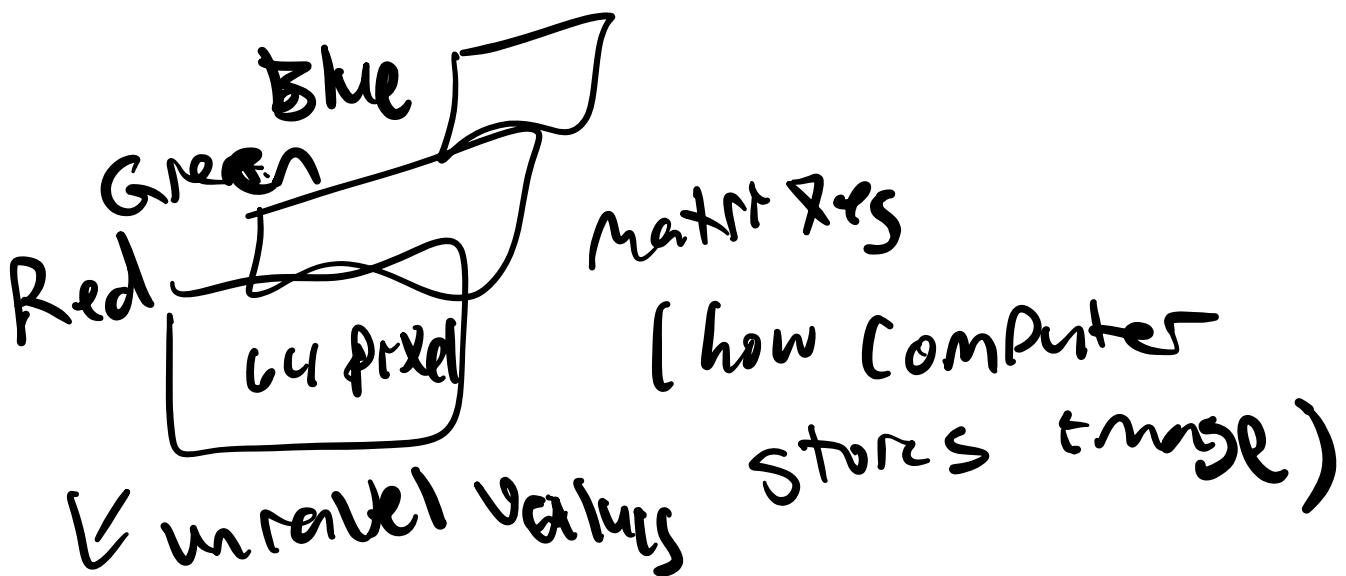
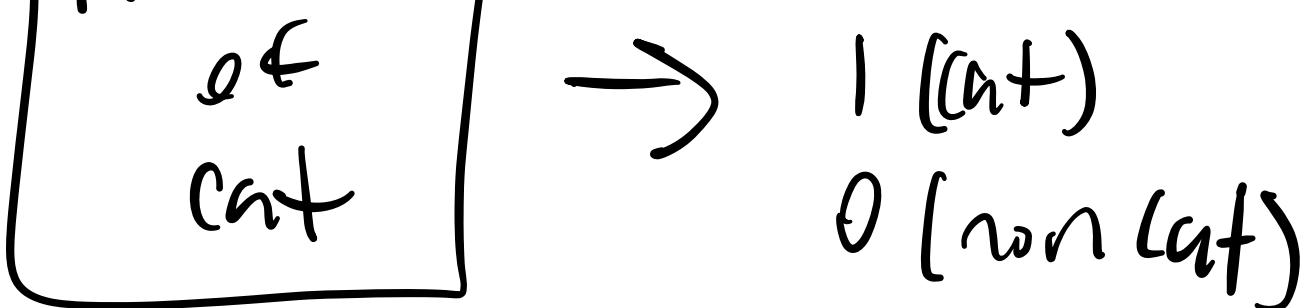


makes
it faster
"convergence"
less
steps
epoch

(M) = size of training set

Binary Classification

Picture



$$X = \begin{bmatrix} a_{11} \\ \vdots \\ a_{1K} \end{bmatrix} \quad 64 \times 64 \times 3 = 12288$$

$$n = n_K = 12288$$

$$X \rightarrow Y$$

$$(X, Y) \quad X \in \mathbb{R}^{n_x}, \quad Y \in \{0, 1\}$$

n training example: $\{(x^1, y^1), (x^{(2)}, y^{(2)})$
 $\dots (x^n, y^n)\}$

$$n = n_{\text{train}} \quad n_{\text{test}} = n_{\text{test}}$$

$$X = \begin{bmatrix} | & | & | \\ x_1 & x_2 & \dots & x_m \\ | & | & | \end{bmatrix}^T$$

$\nwarrow n_x$

$\swarrow m$

$$X \in \mathbb{R}^{n_x \times m} \quad X \cdot \text{Shape} = (n_x, m)$$

$$Y = [y^1, y^2, \dots, y^m]$$

$$Y \in \mathbb{R}^{1 \times m}$$

$$Y \cdot \text{Shape} = (1, m)$$

R_{Real}^2
 R_{Aum}^2

LOGISTIC REGRESSION

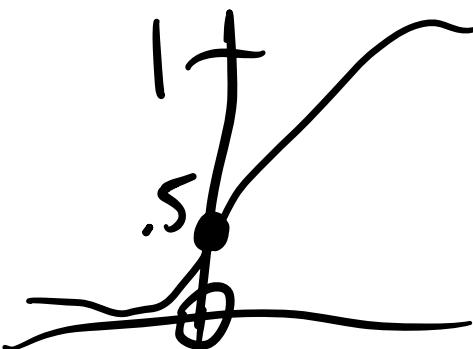
Given X want $\hat{Y} = P(Y=1|X)$

$$X \in \mathbb{R}^{n_x} \quad 0 \leq \hat{Y} \leq 1$$

Parameters: $w \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$

Output: $\hat{y} = \sigma(w^T x + b)$
Sigmoid

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



If z is large

$$G(z) \approx \frac{1}{1+0} = 1$$

If z is large negative

$$G(z) \underset{z \rightarrow -\infty}{\approx} \frac{1}{1+bz} \underset{b \ll 1}{\approx} 0$$

Logistic Regression

Cost function

$\hat{y}^{(i)} = \sigma(w^T x^{(i)} + b)$ where $\sigma(z) = \frac{1}{1+e^{-z}}$

Given $\{(x^1, y^1), \dots, (x^m, y^m)\}$ want $\hat{y}^{(i)} \approx y^{(i)}$

$z^{(i)} = w^T x^{(i)} + b \quad x^{(i)}, y^{(i)}, z^{(i)}$
;ⁱth example

Loss (error) function:

$$L(\hat{y}, y) = -(y \log \hat{y} + (1-y) \log(1-\hat{y}))$$

IF $y=1$: ~~$L(\hat{y}, y) = -\log \hat{y}$~~

IF $y=0$: $L(\hat{y}, y) = -\log(1-\hat{y})$

want $\log(1-\hat{y})$ to
be large, want
 \hat{y} small

want \hat{y}
to be as
big as
poss. b/c
so $\log \hat{y}$ is
big

(cost function):

$$J(w, b) = \frac{1}{m} \sum_i f(\hat{y}^{(i)}, y^{(i)}) = J$$

$$-\frac{1}{n} \sum_{i=1}^n \left[y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)}) \right]$$

Prediction output

want to find w and b that

minimise J (cost function)

Loss - single training example

Cost - average loss entire data set

Gradient Descent



Repeat \mathcal{E}

$w := w - \alpha$ $\frac{d J(w)}{dw} \rightarrow$ derivative
learning rate "dw" → update;
controls slope
size of step
you take at function

3

$$w := w - \alpha dw$$

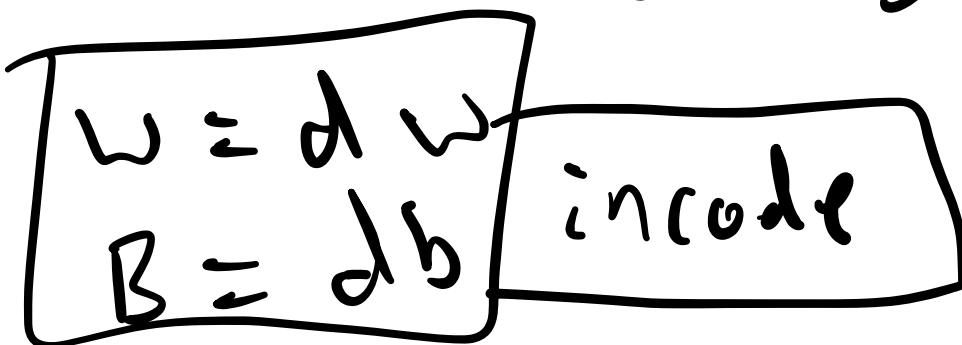
If you start off with low
value of w , gradient will make
you decrease slowly and positive
slope

If you start off with small number and negative slope it increases

Both of these '(decreasing)' tend to move towards which is best (most)

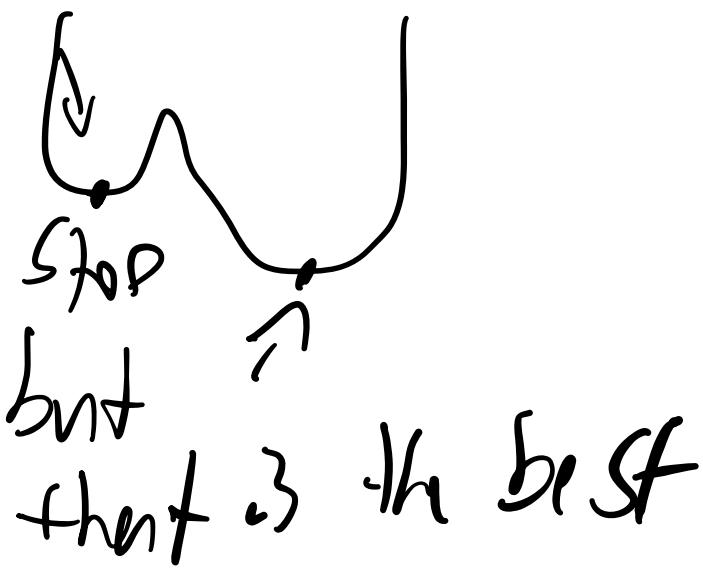
$$J(w, b) \quad w := w - \alpha \underbrace{\frac{dJ(w, b)}{dw}}$$

$$b := b - \alpha \underbrace{\frac{dJ(w, b)}{db}}$$



A convex function does not always have multiple local optima.*

Adag vs Adam



Derivatives

i.e.:

↳ these goes one

$$\frac{d f(a)}{da} = 3 = \frac{d}{da} f(a)$$

↳ these goes

3 times

$$\frac{d}{da} a^2 = 2a$$

0.001

$(2a) \times .0001$

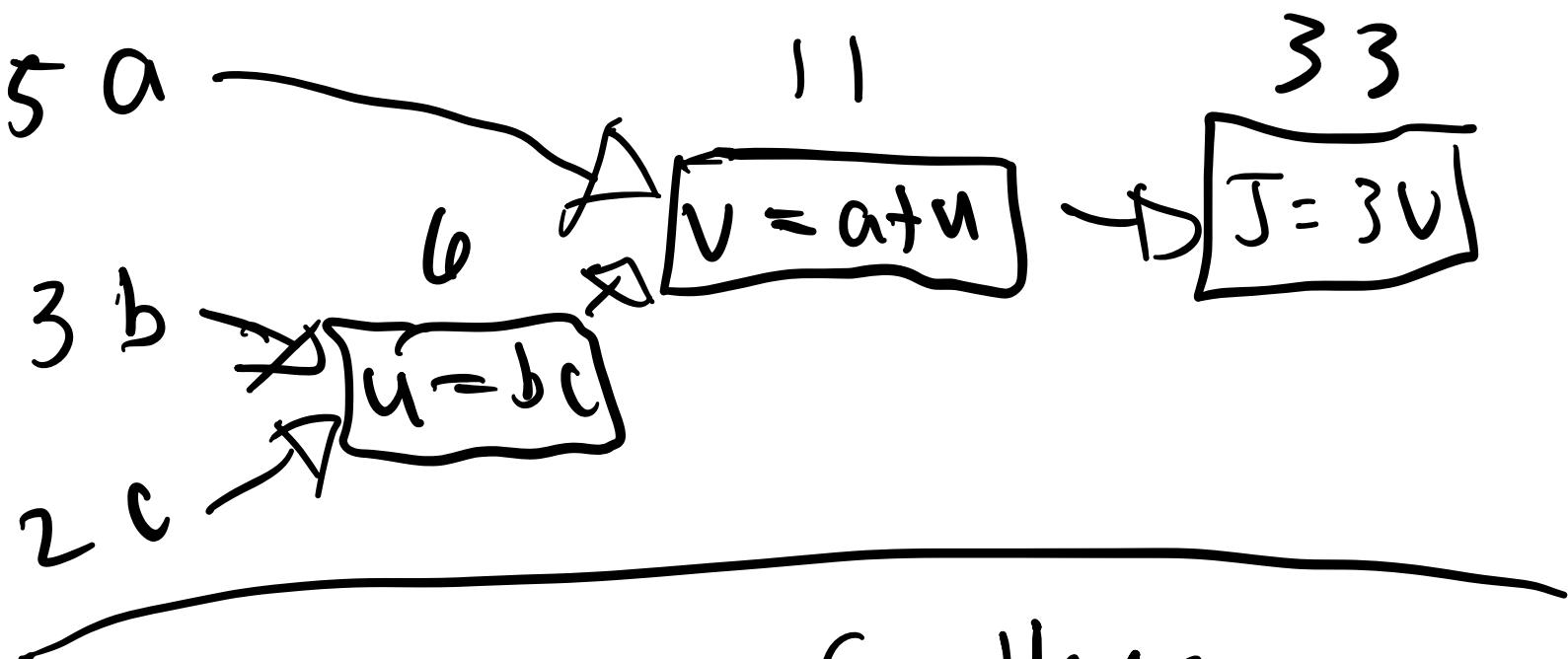
Computation Graph

$$J(a, b, c) = 3(a + bc)$$

$$u = bc$$

$$v = a + u$$

$$J = 3v$$



Combining last few things

Linear Regression

$$z = w^T x + b$$

$$\hat{y} = a = \sigma(z)$$

$$L(a, y) = -(y \log(a) + (1-y) \log(1-a))$$

$$x_1 \rightarrow \\ w_1 \rightarrow \\ x_2 \rightarrow \\ w_2 \rightarrow \\ b \rightarrow z = w_1 x_1 + w_2 x_2 + b \rightarrow \hat{y} = a = \sigma(z) \rightarrow L(a, y)$$

$$dz = \frac{dL}{dz} = \frac{dL(a,y)}{da} da = \frac{dL(a,y)}{da}$$

$$= a - y = -\frac{y}{a} + \frac{1-y}{1-a}$$

$$\frac{\partial L}{\partial w_1} = "dw_1" = x_1 \cdot dz$$

$$\frac{\partial L}{\partial w_2} = "dw_2" = x_2 \cdot dz$$

$$w_1 = w_1 - \alpha dw_1$$

$$w_2 = w_2 - \alpha dw_2$$

$$"db" \approx dz$$

$$b = b - \alpha db$$

$$J=0 \quad dw=0 \quad dw_2=0 \quad db=0$$

for $i=1 \text{ to } N$

$$z^{(i)} = w^T x^{(i)} + b$$

$$a^{(i)} - \sigma(z^{(i)})$$

$$J = -\left[y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log (1-a^{(i)}) \right]$$

$$dZ^{(i)} = a^{(i)} - y^{(i)}$$

$$dw_1^+ = X_1^{(i)} dZ^{(i)}$$

$$dw_2^+ = X_2^{(i)} dZ^{(i)}$$

$$db^+ = dZ^{(i)}$$

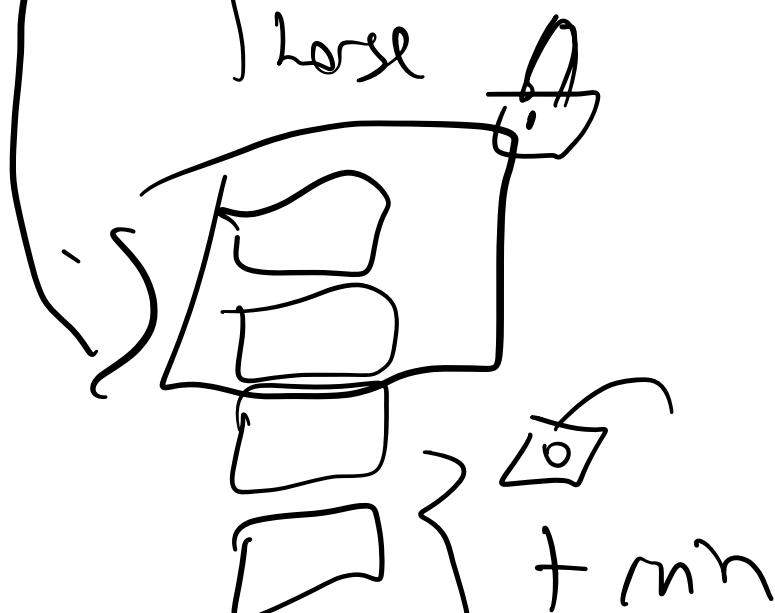
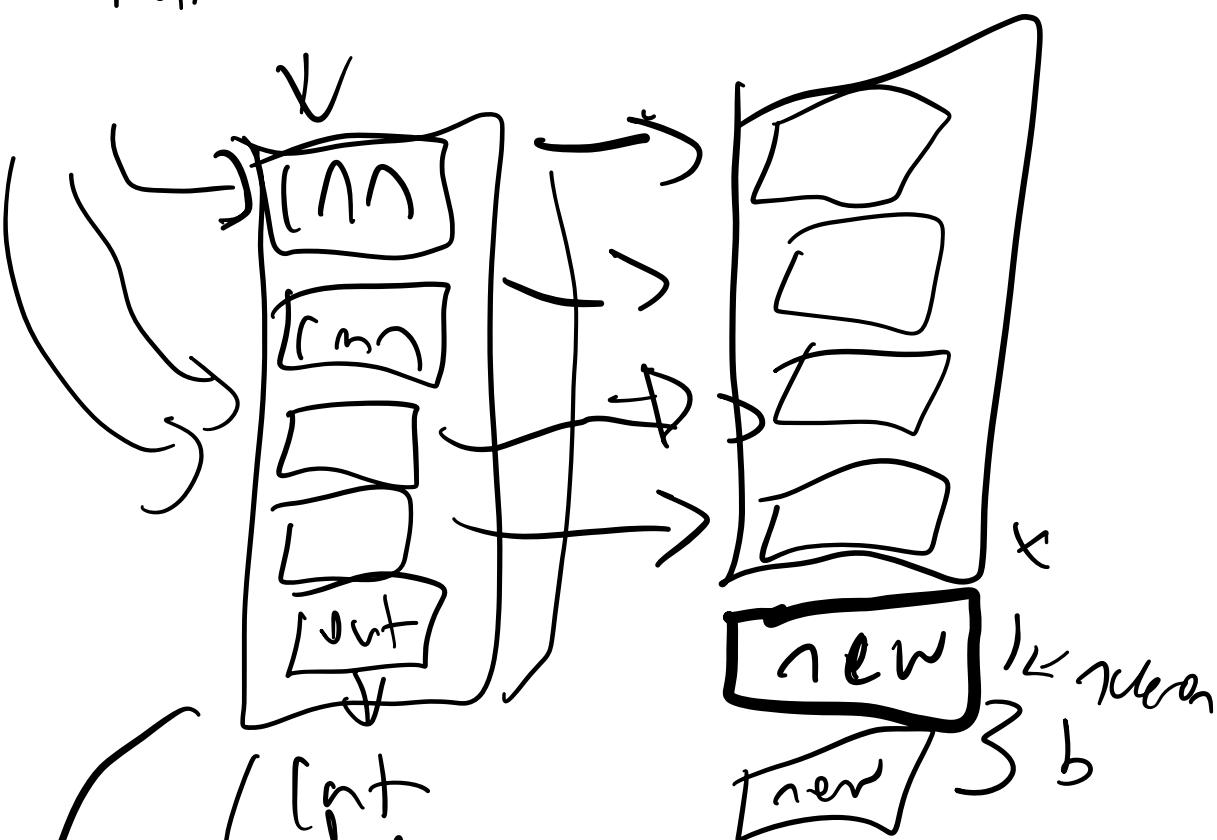
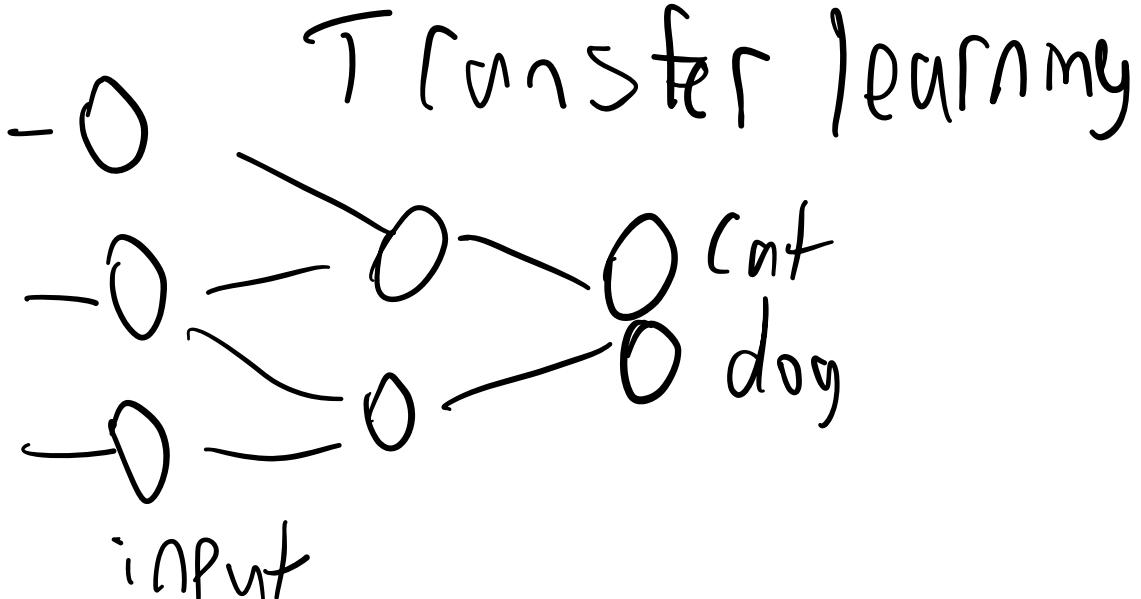
$$J = m dw_1^+ + m dw_2^+ + m db^+ = m$$

Vectorization (Getting rid of loops)

Back Propagation

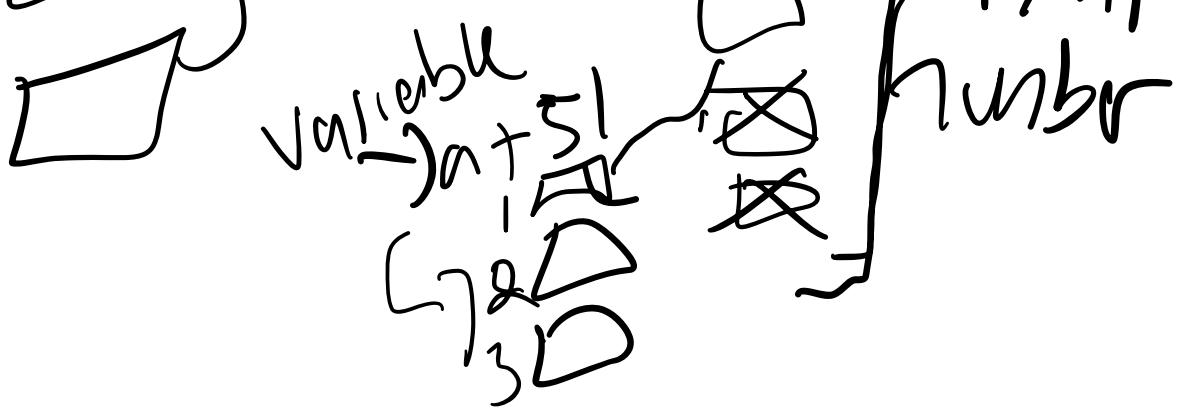
using chain rule from gradient

(all derived)



Res 52
PyTorch





Sample classification

Tensorflow Sample codes

Transfer learning.

Data augmentation
ie

Rotated image of
Cats workshop

Weights Pre-trained set

Sigmoid - Logistic

$$\text{Sigmoid}(x) = \frac{1}{1+e^{-x}}$$

Regression function

use "np. whatever()"

if $x = \text{np. array}[1, 2, 3]$ then
does \exp for every y val in array
same if its $\text{print}(x+3)$ $\rightarrow [4, 5, 6]$

$$\text{For } x \in \mathbb{R}^n, \text{ Sigmoid}(x) = \text{Sigmoid} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} \frac{1}{1+e^{-x_1}} \\ \frac{1}{1+e^{-x_2}} \\ \vdots \\ \frac{1}{1+e^{-x_n}} \end{pmatrix}$$

Coding Implementation

```
def Sigmoid(x):  
    s = 1 / (1 + np.exp(-x))  
    return s
```

Sigmoid Gradient

Optimization using back propagation

Sigmoid derivative:

Proposition

$$\sigma' = \sigma(x)(1 - \sigma(x))$$

Coding Implementation:

```
def sigmoid_derivative(x):
```

$$s = \text{sigmoid}(x)$$

$$ds = s * (1 - s)$$

```
return ds
```

Reshaping

NP. Shape - tells you shape

NP. Reshape - this is how you turn
a for example image input of
(length, height, width) into a 1D array
So it simple for the AI training
use reshape(-1,1) - it just takes
everything if for images it takes
length, height, and width and multiplies.

Coding Implementation:

```
def image2vector(image):
```

$$v = v.reshape(-1, 1).T$$

```
return v
```

.T is transpose which is
needed to make the training example a column

Norm of Matrix

Keep gradients (derivatives) more consistent which helps it converge faster during training

$$X \rightarrow \frac{X}{\|X\|} \quad \|X\| = \sqrt{x_1^2 + x_2^2 + \dots + x_n^2}$$

i.e. $X = \begin{bmatrix} 0 & 3 & 4 \\ 2 & 6 & 1 \end{bmatrix}$

$$\|X\| = \text{np.linalg.norm}(X, \text{axis}=1, \text{keepdims}=\text{True})$$
$$= \begin{bmatrix} \sqrt{50} \end{bmatrix} \quad X_{\text{Normalized}} = \begin{bmatrix} 0 & \frac{3}{\sqrt{50}} & \frac{4}{\sqrt{50}} \\ \frac{2}{\sqrt{50}} & \frac{6}{\sqrt{50}} & \frac{1}{\sqrt{50}} \end{bmatrix}$$

$\text{keepdims} = \text{True}$: keeps it a 2D array so it can be broadcasted back to original matrix during $\alpha X \beta = 1$; calculates norm row by row

Coding Implementation:

```
def normalize_rows(x):
```

```
    norm = np.linalg.norm(x, axis=1, keepdims=True)
```

```
    x = x / norm
```

```
return x
```

SoftMax (Broadcasting)
evidence

```
def softmax(x):
```

```
x_exp = np.exp(x)
```

```
x_sum = np.sum(x, axis=1, keepdims=True)
```

```
s = x_exp / x_sum
```

```
return s
```

Vectorization

Dot Product

NP. dot(x_1, x_2)

$$\sum_{i=1}^n x_1[i] \cdot x_2[i]$$

could be used for weighted sum of inputs

$$x_1 = [1, 2, 3] \rightarrow \text{dot} = (1 \times 4) + (2 \times 5) + (3 \times 6)$$

$$x_2 = [4, 5, 6]$$

Outer Product

NP. outer(x_1, x_2)

$$\text{outer}(x_1, x_2)[i,j] = x_1[i] \cdot x_2[j]$$

Used for constructing matrices from vectors

$$x_1 = [1, 2] \rightarrow \begin{bmatrix} [1 \times 3, 1 \times 4] \\ [2 \times 3, 2 \times 4] \end{bmatrix} \rightarrow \begin{bmatrix} [3, 4] \\ [6, 8] \end{bmatrix}$$

$$x_2 = [3, 4] \rightarrow \text{outer} \quad \text{result}$$

Element wise Multiplication

NP. multiply(x_1, x_2) elementwise(x_1, x_2)[i] =

$$x_1[i] \cdot x_2[i]$$

used for feature wide scaling

$$x_1 = [1, 2, 3] \rightarrow [1 \times 4, 2 \times 5, 3 \times 6] \text{ Mult}$$

$x2 = \{4, 5, 6\} [4, 10, 18] \text{ result}$

General Dot Product

NP. dot (W, X1) $\text{gdot}(W, X1)[i] = \sum_j W_{i,j} \cdot X1[j]$

Core operation in Linear Regression

(calculating weighted sums in hidden layers)

$$W = \begin{pmatrix} [1, 2] \\ [3, 4] \\ [5, 6] \end{pmatrix}, \quad \begin{matrix} [1 \times 5 + 2 \times 6, 3 \times 5 + 4 \times 6] \\ \downarrow \text{dot} \\ [17, 39] \end{matrix} \text{ result}$$

$$X1 = ([5, 6])$$

LOSS FUNCTIONS

sparse data
heavy outliers

def L1(Yhat, Y):

"Y hat = Predicted"

"Y = true"

$$L1(\hat{y}, y) = \sum_{i=0}^{M-1} |y^{(i)} - \hat{y}^{(i)}|$$

$$\text{LOSS} = \text{NP.SUM}(Y - Y\text{hat})$$

return loss

def L2(yhat, y):
$$L_2(\hat{y}, y) = \sum_{i=0}^{M-1} (y^{(i)} - \hat{y}^{(i)})^2$$

loss = np.sum(np.square(y - yhat))

return loss

Minimize loss
Balanced data
more sensitive
to outliers

CNN Logistic Regression Example

We have

1*train_set_X_org which is numpy array
of shape (M-train, num_px, num_px, 3)

*train_set_X_org[0] which is M-train is
the number of training examples

*train_set_X_org[1] which is num_px is

the height of the image

2) Flatten it for the model to process

$\text{train_set_X_flattened} = \text{train_set_X_orig.reshape}(\text{trainset_X_orig.shape[0]}, -1).\text{T}$

-1 : flattens all other dimensions

.T : transposes to a column required by model

3) Standardize the dataset

$\text{train_set_X} = \text{train_set_flattened}/255$

Steps to build Neural Network

0. Preprocess data

1. Define Model Structure (ie: # of input features)

2. Initialize model's parameters

3. Loop

* calculate current loss (forward propagation)

* calculate current gradient (backward propagation)

* update parameters (gradient descent)

Helper function

Sigmoid(z):

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

Return S

Initialize Parameters

initialize-with-zeros(dim):

w = np.zeros([dim, 1]) # column vector
b = 0.0 # bias scalar of zeros

Return w, b

Forward & Backward Propagation

Propagate(w, b, X, Y):

cost function & gradient

Arguments

W - weights (numpy array of size [num_px * num_px * 3, 1])

b - bias scalar

X - data of size [num_px * num_px * 3, # of examples]

Y - true label vector (0 if not cat 1 if cat)

```
# Return  
grads - dictionary containing gradients  
(dw - gradient of loss with respect to w)  
(db - gradient of loss with respect to b)  
cost - negative log-likelihood cost for logistic  
regression
```

|||||

```
m = X.shape[1]
```

```
# Forward propagation
```

```
# Sigmoid activation function
```

```
A = sigmoid(np.dot(w.T, X) + b))
```

```
# instead of there is linear function
```

```
# now comes log loss function
```

```
cost = -(1/m) * np.sum(y * np.log(A) +  
(1-y) * np.log(1-A))
```

```
# Backward propagation
```

```
dw = (1/m) * np.dot(X.T, (y - A))
```

```
dw = (1/m) * np.dot(X, (A - Y))
# Gradient of w
# Gradient of b
db = (1/m) * np.sum(A - Y)

cost = np.squeeze(np.array(cost))

grads = {"dw": dw, "db": db}

return grads, cost
```



Optimization

```
# runs a gradient descent algorithm
def optimize(w, b, X, Y, num_iterations=100, learning_rate=0.009, print_cost=False)
```

avoid modifying original ones

```
w = copy.deepcopy(w)
b = copy.deepcopy(b)

# list all costs during optimization
# to plot learning curve
costs = []
for i in range(num_iterations):
    # cost & gradient calculation
    grads, cost = propagate(w, b, X, Y)
    # get derivatives
    dw = grads["dw"]
    db = grads["db"]
    # update rule
    w = w - learning_rate * dw
    b = b - learning_rate * db
    costs.append(cost)
```

```
# record cost
if i%100 == 0:
    costs.append(cost)

# print cost every 100 iterations
for i in range(1000):
    if print_cost:
        print("cost after {} iterations".format(i))
        print("  : {}%".format(cost))

params = {"w": w, "b": b}

grads = {"dw": dw, "db": db}

return params, grads, cost
```

Predict whether label is 0 or 1

```
def Predict(w, b, X):
```

```
M = X.shape[1]
```

```
Y_Prediction = np.zeros((1, n))
# initialization
# make sure w is right shape
w = w.reshape(x.shape[0], 1)
# activation
A = sigmoid(np.dot(w.T, x) + b)
# prediction
for i in range(A.shape[1]):
    if A[0, i] > .5:
        Y_Prediction[0, i] = 1
    else:
        Y_Prediction[0, i] = 0
return Y_Prediction
```

Merge all functions into mode)

```
def Model(X_train, Y_train, X_test,  
Y_test, num_iterations=2000, learning_  
rate=0.5, print_cost=False):  
    # Initialize  
    W, b = initialize_w_with_zeros(X_train.shape[0])  
    # Gradient descent  
    Params, grads, costs = optimize(W, b,  
        X_train, Y_train, num_iterations, learning_  
        rate, print_cost)  
    # Retrieve parameters  
    w = Params["w"]  
    b = Params["b"]  
    # Predict on training set  
    Y_prediction_train = predict(w, b, X_train)  
    # Predict on test set  
    Y_prediction_test = predict(w, b, X_test)
```

```
# Print results  
print("train accuracy: {}%".format(100 - np.mean(np.abs(Y_Prediction_train - Y_train)) * 100))  
print("test accuracy: {}%".format(100 - np.mean(np.abs(Y_Prediction_test - Y_test)) * 100))
```

Store relevant outputs

d = {
 "costs": cost,
 "Y_Prediction-test": Y_Prediction_test,
 "Y_Prediction-train": Y_Prediction_train,
 "w": w,
 "b": b,
 "learning_rate": learning_rate,
 "num_iterations": num_iterations}

return d

Neural Network

How to tell Shape: `numpyarray.shape`
of training examples: `n = nparray.shape[0]`
Methodology:

1. Define neural network structure (# of inputs, # of hidden units, etc.)

2. Initialize model's parameters

3. Loop

- Implement forward propagation

- compute loss

- Implement backward propagation [get gradients]

- update parameters (gradient descent)

1. Define Structure

```
def layer_sizes(X, Y):
```

X - data set size (input size, # of examples)

Y - labels of shape (output size, # of examples)

$n_x = X.shape[0]$ # input layer size

$n_h = 4$ # hidden layer size

$n_y = Y.shape[0]$ # output layer size

return n_x, n_h, n_y

2. Initialize Parameters

def initialize_parameters(n_x, n_h, n_y):

Returns: Params (dictionary)

w_1 - weight matrix of shape (n_h, n_x)

b_1 - bias vector of shape ($n_h, 1$)

w_2 - weight matrix of shape (n_y, n_h)

b_2 - bias vector of shape ($n_y, 1$)

$w_1 = np.random.rand(n_h, n_x) * 0.01$

$b_1 = np.zeros((n_h, 1))$

$w_2 = np.random.rand(n_y, n_h) * 0.01$

$b_2 = np.zeros((n_y, 1))$

Parameters = { w_1 ''
etc.

return Parameters

3. The loop

```
def forward_propagation(X, parameters):
```

X input data of size (1 x, m)

W1 = parameters["W1"]

etc. W2, b1, b2

Z1 = np.dot(W1, X) + b1

A1 = np.tanh(Z1)

Z2 = np.dot(W2, A1) + b2

A2 = sigmoid(Z2)

assert(A2.shape == (1, X.shape[1]))

cache = {"Z1": Z1
etc.

return A2, cache

```
def compute_cost(A2, Y):
```

M = Y.shape[1] # of examples

Cross entropy cost

logProbs = -np.multiply(np.log(A2), Y) +
-np.multiply(np.log(1-A2), (1-Y))

Apply full cost function formula

Cost = -np.sum(logProbs) / M

Cost = float(np.mean(Cost))

Makes sure dimension is how
we expect eg: [57] > 17

return Cost

def back_propagation(parameters,
cache, X, Y):

M = X.shape[1]

W1 = parameters["W1"]

W2 = parameters["W2"]

A1 = cache["A1"]

A2 = cache["A2"]

dZ2 = A2 - Y # compute error at output layer

dW2 = (1/M) * np.dot(dZ2, A1.T) # gradient of W2

gradient of b2

db2 = (1/M) * np.sum(dZ2, axis=1, keepdims=True)

Back Propagate through tanh activation

dZ1 = np.dot(W2.T, dZ2) * (1 - np.tanh(A1)**2)

dW1 = (1/M) * np.dot(dZ1, X.T) # gradient of W1

gradient of b1

```
db2 = (1/m) * np.sum(dz1, axis=1, keepdims=True)
```

```
grads = {dw1: db1, dw2: db2}
```

return grads

```
def update_parameters(parameters, grads, learning_rate=1.2):
```

```
w1 = copy.deepcopy(parameters["w1"])  
b1, w2, b2 = (None,) * 3
```

deep copy of parameters

retrieve each gradient

```
dw1 = grads["dw1"]
```

```
b1, dw2, db2 = (None,) * 3
```

update rule

```
w1 = w1 - learning_rate * dw1
```

```
b1, w2, b2 = (None,) * 3
```

```
parameters = {w1, j1, w2, b2}
```

return parameters

Integration into model

```
def nn_model(X, Y, n_h) num_iterations = 10000, print_cost=False):
    np.random.seed(3)
    n_x = layer_sizes(X, Y)[0]
    n_y = layer_sizes(X, Y)[2]
    # initialize parameters
    parameters = initialize_parameters(n_x, n_h, n_y)
    # gradient descent loop
    for i in range(num_iterations):
        # forward propagation
        A2, cache = forward_propagation(X, parameters)
        # compute cost
        cost = compute_cost(A2, Y)
        # back propagation
        grads = backward_propagation(parameters, cache, X, Y)
        # update parameters
```

Update parameters
parameters = update_parameters(parameters, grads)
if print_cost and i % 1000 == 0:
print("cost after iteration " + str(i) + " : " + str(cost))

return parameters

Test Model

```
def predict(Parameters, X):  
# forward propagation to get prediction  
A2, cache = forwardPropagation(X, Parameters)  
# convert probabilities into binary predictions  
Predictions = (A2 > .5).astype(int)  
return Predictions
```