

Your very own neural network

In this notebook we're going to build a neural network using naught but pure numpy and steel nerves. It's going to be fun, I promise!



In [1]:

```
import sys
sys.path.append("..")
import tqdm_utils
import download_utils
```

In [2]:

```
# use the preloaded keras datasets and models
download_utils.link_all_keras_resources()
```

In [3]:

```
from __future__ import print_function
import numpy as np
np.random.seed(42)
```

Here goes our main class: a layer that can do `.forward()` and `.backward()` passes.

In [4]:

```
class Layer:
    """
    A building block. Each layer is capable of performing two things:

    - Process input to get output:          output = layer.forward(input)

    - Propagate gradients through itself:    grad_input = layer.backward(input,
grad_output)

    Some layers also have learnable parameters which they update during layer.backward.
    """
    def __init__(self):
        """Here you can initialize layer parameters (if any) and auxiliary stuff."""
        # A dummy layer does nothing
        pass

    def forward(self, input):
        """
        Takes input data of shape [batch, input_units], returns output data [batch, output_units]
        """
        # A dummy layer just returns whatever it gets as input.
        return input

    def backward(self, input, grad_output):
        """
        Performs a backpropagation step through the layer, with respect to the given input.

        To compute loss gradients w.r.t input, you need to apply chain rule (backprop):

        
$$d \text{ loss} / d x = (d \text{ loss} / d \text{ layer}) * (d \text{ layer} / d x)$$


        Luckily, you already receive  $d \text{ loss} / d \text{ layer}$  as input, so you only need to multiply it by  $d \text{ layer} / d x$ .

        If your layer has parameters (e.g. dense layer), you also need to update them here using  $d \text{ loss} / d \text{ layer}$ 
        """
        # The gradient of a dummy layer is precisely grad_output, but we'll write it more explicitly
        num_units = input.shape[1]

        d_layer_d_input = np.eye(num_units)

        return np.dot(grad_output, d_layer_d_input) # chain rule
```

The road ahead

We're going to build a neural network that classifies MNIST digits. To do so, we'll need a few building blocks:

- Dense layer - a fully-connected layer, $f(X) = W \cdot X + \vec{b}$
- ReLU layer (or any other nonlinearity you want)
- Loss function - crossentropy
- Backprop algorithm - a stochastic gradient descent with backpropagated gradients

Let's approach them one at a time.

Nonlinearity layer

This is the simplest layer you can get: it simply applies a nonlinearity to each element of your network.

In [5]:

```
class ReLU(Layer):
    def __init__(self):
        """ReLU layer simply applies elementwise rectified linear unit to all in
puts"""
        pass

    def forward(self, input):
        """Apply elementwise ReLU to [batch, input_units] matrix"""
        # <your code. Try np.maximum>
        re=np.maximum(input, np.zeros(input.shape))
        return re
    def backward(self, input, grad_output):
        """Compute gradient of loss w.r.t. ReLU input"""
        relu_grad = input > 0
        return grad_output*relu_grad
```

In [6]:

```
# some tests
from util import eval_numerical_gradient
x = np.linspace(-1,1,10*32).reshape([10,32])
l = ReLU()
grads = l.backward(x,np.ones([10,32])/(32*10))
numeric_grads = eval_numerical_gradient(lambda x: l.forward(x).mean(), x=x)
assert np.allclose(grads, numeric_grads, rtol=1e-3, atol=0),\
    "gradient returned by your layer does not match the numerically computed gra
dient"
```

Instant primer: lambda functions

In python, you can define functions in one line using the `lambda` syntax: `lambda param1, param2: expression`

For example: `f = lambda x, y: x+y` is equivalent to a normal function:

```
def f(x,y):  
    return x+y
```

For more information. click here (http://www.secnetix.de/olli/Python/lambda_functions.hawk).

Dense layer

Now let's build something more complicated. Unlike nonlinearity, a dense layer actually has something to learn.

A dense layer applies affine transformation. In a vectorized form, it can be described as:

$$f(X) = W \cdot X + \vec{b}$$

Where

- X is an object-feature matrix of shape `[batch_size, num_features]`,
- W is a weight matrix `[num_features, num_outputs]`
- and b is a vector of `num_outputs` biases.

Both W and b are initialized during layer creation and updated each time backward is called.

In [7]:

```

class Dense(Layer):
    def __init__(self, input_units, output_units, learning_rate=0.1):
        """
        A dense layer is a layer which performs a learned affine transformation:
         $f(x) = \langle W * x \rangle + b$ 
        """
        self.learning_rate = learning_rate

        # initialize weights with small random numbers. We use normal initialization,
        # but surely there is something better. Try this once you got it working:
        # http://bit.ly/2vTlmaJ
        self.weights = np.random.randn(input_units, output_units)*0.01
        self.biases = np.zeros(output_units)

    def forward(self, input):
        """
        Perform an affine transformation:
         $f(x) = \langle W * x \rangle + b$ 

        input shape: [batch, input_units]
        output shape: [batch, output units]
        """
        return np.dot(input, self.weights) + self.biases

    def backward(self, input, grad_output):

        # compute  $df / dx = df / ddense * ddense / dx$ 
        # where  $ddense / dx = weights^T$ 
        grad_input = np.dot(grad_output, self.weights.T) #<your code here>

        # compute gradient w.r.t. weights and biases
        grad_weights = np.dot(input.T, grad_output) #<your code here>
        grad_biases = np.sum(grad_output, axis = 0) #<your code here>

        assert grad_weights.shape == self.weights.shape and grad_biases.shape == self.biases.shape
        # Here we perform a stochastic gradient descent step.
        # Later on, you can try replacing that with something better.
        self.weights = self.weights - self.learning_rate * grad_weights
        self.biases = self.biases - self.learning_rate * grad_biases

        return grad_input

```

Testing the dense layer

Here we have a few tests to make sure your dense layer works properly. You can just run them, get 3 "well done"s and forget they ever existed.

... or not get 3 "well done"s and go fix stuff. If that is the case, here are some tips for you:

- Make sure you compute gradients for W and b as **sum of gradients over batch**, not mean over gradients. Grad_output is already divided by batch size.
- If you're debugging, try saving gradients in class fields, like "self.grad_w = grad_w" or print first 3-5 weights. This helps debugging.
- If nothing else helps, try ignoring tests and proceed to network training. If it trains alright, you may be off by something that does not affect network training.

In [8]:

```
l = Dense(128, 150)

assert -0.05 < l.weights.mean() < 0.05 and 1e-3 < l.weights.std() < 1e-1,\
    "The initial weights must have zero mean and small variance. "\
    "If you know what you're doing, remove this assertion."
assert -0.05 < l.biases.mean() < 0.05, "Biases must be zero mean. Ignore if you
    have a reason to do otherwise."

# To test the outputs, we explicitly set weights with fixed values. DO NOT DO TH
AT IN ACTUAL NETWORK!
l = Dense(3,4)

x = np.linspace(-1,1,2*3).reshape([2,3])
l.weights = np.linspace(-1,1,3*4).reshape([3,4])
l.biases = np.linspace(-1,1,4)

assert np.allclose(l.forward(x),np.array([[ 0.07272727,  0.41212121,  0.75151515
,  1.09090909],
                                         [-0.90909091,  0.08484848,  1.07878788
,  2.07272727]]))
print("Well done!")
```

Well done!

In [9]:

```
# To test the grads, we use gradients obtained via finite differences

from util import eval_numerical_gradient

x = np.linspace(-1,1,10*32).reshape([10,32])
l = Dense(32,64,learning_rate=0)

numeric_grads = eval_numerical_gradient(lambda x: l.forward(x).sum(),x)
grads = l.backward(x,np.ones([10,64]))

assert np.allclose(grads,numeric_grads,rtol=1e-3,atol=0), "input gradient does n
ot match numeric grad"
print("Well done!")
```

Well done!

In [10]:

```
#test gradients w.r.t. params
def compute_out_given_wb(w,b):
    l = Dense(32,64,learning_rate=1)
    l.weights = np.array(w)
    l.biases = np.array(b)
    x = np.linspace(-1,1,10*32).reshape([10,32])
    return l.forward(x)

def compute_grad_by_params(w,b):
    l = Dense(32,64,learning_rate=1)
    l.weights = np.array(w)
    l.biases = np.array(b)
    x = np.linspace(-1,1,10*32).reshape([10,32])
    l.backward(x,np.ones([10,64]) / 10.)
    return w - l.weights, b - l.biases

w,b = np.random.randn(32,64), np.linspace(-1,1,64)

numeric_dw = eval_numerical_gradient(lambda w: compute_out_given_wb(w,b).mean(0)
    .sum(),w )
numeric_db = eval_numerical_gradient(lambda b: compute_out_given_wb(w,b).mean(0)
    .sum(),b )
grad_w,grad_b = compute_grad_by_params(w,b)

assert np.allclose(numeric_dw,grad_w,rtol=1e-3,atol=0), "weight gradient does not match numeric weight gradient"
assert np.allclose(numeric_db,grad_b,rtol=1e-3,atol=0), "weight gradient does not match numeric weight gradient"
print("Well done!")
```

Well done!

The loss function

Since we want to predict probabilities, it would be logical for us to define softmax nonlinearity on top of our network and compute loss given predicted probabilities. However, there is a better way to do so.

If you write down the expression for crossentropy as a function of softmax logits (a), you'll see:

$$loss = -\log \frac{e^{a_{correct}}}{\sum_i e^{a_i}}$$

If you take a closer look, ya'll see that it can be rewritten as:

$$loss = -a_{correct} + \log \sum_i e^{a_i}$$

It's called Log-softmax and it's better than naive log(softmax(a)) in all aspects:

- Better numerical stability
- Easier to get derivative right
- Marginally faster to compute

So why not just use log-softmax throughout our computation and never actually bother to estimate probabilities.

Here you are! We've defined the both loss functions for you so that you could focus on neural network part.

In [11]:

```
def softmax_crossentropy_with_logits(logits,reference_answers):
    """Compute crossentropy from logits[batch,n_classes] and ids of correct answers"""
    logits_for_answers = logits[np.arange(len(logits)),reference_answers]

    xentropy = - logits_for_answers + np.log(np.sum(np.exp(logits),axis=-1))

    return xentropy

def grad_softmax_crossentropy_with_logits(logits,reference_answers):
    """Compute crossentropy gradient from logits[batch,n_classes] and ids of correct answers"""
    ones_for_answers = np.zeros_like(logits)
    ones_for_answers[np.arange(len(logits)),reference_answers] = 1

    softmax = np.exp(logits) / np.exp(logits).sum(axis=-1,keepdims=True)

    return (- ones_for_answers + softmax) / logits.shape[0]
```

In [12]:

```
logits = np.linspace(-1,1,500).reshape([50,10])
answers = np.arange(50)%10

softmax_crossentropy_with_logits(logits,answers)
grads = grad_softmax_crossentropy_with_logits(logits,answers)
numeric_grads = eval_numerical_gradient(lambda l: softmax_crossentropy_with_logits(l,answers).mean(),logits)

assert np.allclose(numeric_grads,grads,rtol=1e-3,atol=0), "The reference implementation has just failed. Someone has just changed the rules of math."
```

Full network

Now let's combine what we've just built into a working neural network. As we announced, we're gonna use this monster to classify handwritten digits, so let's get them loaded.

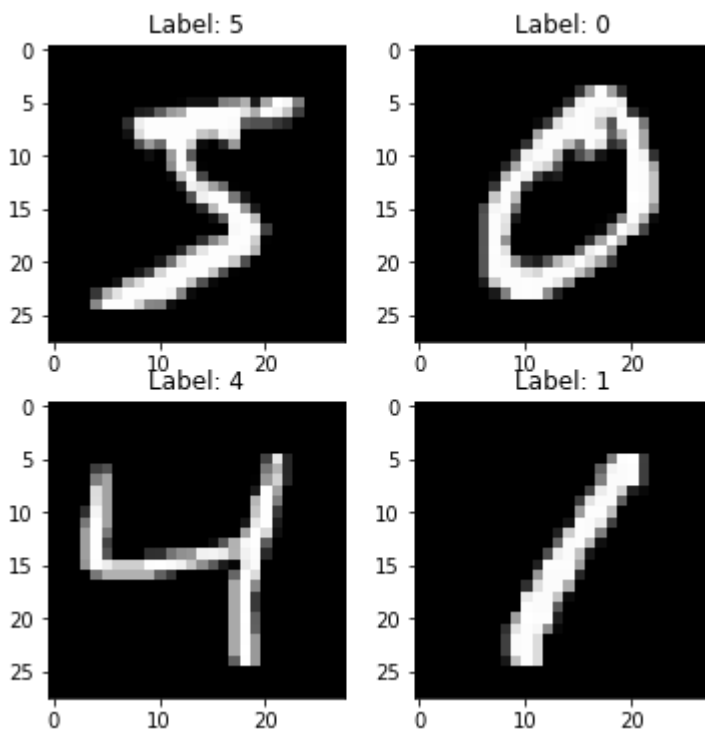
In [13]:

```
import matplotlib.pyplot as plt
%matplotlib inline

from preprocessed_mnist import load_dataset
X_train, y_train, X_val, y_val, X_test, y_test = load_dataset(flatten=True)

plt.figure(figsize=[6,6])
for i in range(4):
    plt.subplot(2,2,i+1)
    plt.title("Label: %i"%y_train[i])
    plt.imshow(X_train[i].reshape([28,28]),cmap='gray');
```

Using TensorFlow backend.



We'll define network as a list of layers, each applied on top of previous one. In this setting, computing predictions and training becomes trivial.

In [14]:

```
network = []
network.append(Dense(X_train.shape[1],100))
network.append(ReLU())
network.append(Dense(100,200))
network.append(ReLU())
network.append(Dense(200,10))
```

In [15]:

```
def forward(network, X):
    """
    Compute activations of all network layers by applying them sequentially.
    Return a list of activations for each layer.
    Make sure last activation corresponds to network logits.
    """
    activations = []
    input = X

    # <your code here>
    for layer in network:
        output = layer.forward(input)
        activations.append(output)
        input = output

    assert len(activations) == len(network)
    return activations

def predict(network,X):
    """
    Compute network predictions.
    """
    logits = forward(network,X)[-1]
    return logits.argmax(axis=-1)

def train(network,X,y):
    """
    Train your network on a given batch of X and y.
    You first need to run forward to get all layer activations.
    Then you can run layer.backward going from last to first layer.

    After you called backward for all layers, all Dense layers have already made
    one gradient step.
    """

    # Get the layer activations
    layer_activations = forward(network,X)
    layer_inputs = [X]+layer_activations #layer_input[i] is an input for network[i]
    logits = layer_activations[-1]

    # Compute the loss and the initial gradient
    loss = softmax_crossentropy_with_logits(logits,y)
    loss_grad = grad_softmax_crossentropy_with_logits(logits,y)

    # <your code: propagate gradients through the network>
    grad_output = loss_grad
    layer_inputs = layer_inputs[:-1]
    for input,layer in zip(layer_inputs[:-1],network[:-1]):
        grad_output = layer.backward(input,grad_output)

    return np.mean(loss)
```

Instead of tests, we provide you with a training loop that prints training and validation accuracies on every epoch.

If your implementation of forward and backward are correct, your accuracy should grow from 90~93% to >97% with the default network.

Training loop

As usual, we split data into minibatches, feed each such minibatch into the network and update weights.

In [16]:

```
def iterate_minibatches(inputs, targets, batchsize, shuffle=False):
    assert len(inputs) == len(targets)
    if shuffle:
        indices = np.random.permutation(len(inputs))
    for start_idx in tqdm_utils.tqdm_notebook_failsafe(range(0, len(inputs) - batchsize + 1, batchsize)):
        if shuffle:
            excerpt = indices[start_idx:start_idx + batchsize]
        else:
            excerpt = slice(start_idx, start_idx + batchsize)
        yield inputs[excerpt], targets[excerpt]
```

In [17]:

```
from IPython.display import clear_output
train_log = []
val_log = []
```

In [18]:

```

for epoch in range(25):

    for x_batch,y_batch in iterate_minibatches(X_train,y_train,batchsize=32,shuffle=True):
        train(network,x_batch,y_batch)

    train_log.append(np.mean(predict(network,X_train)==y_train))
    val_log.append(np.mean(predict(network,X_val)==y_val))

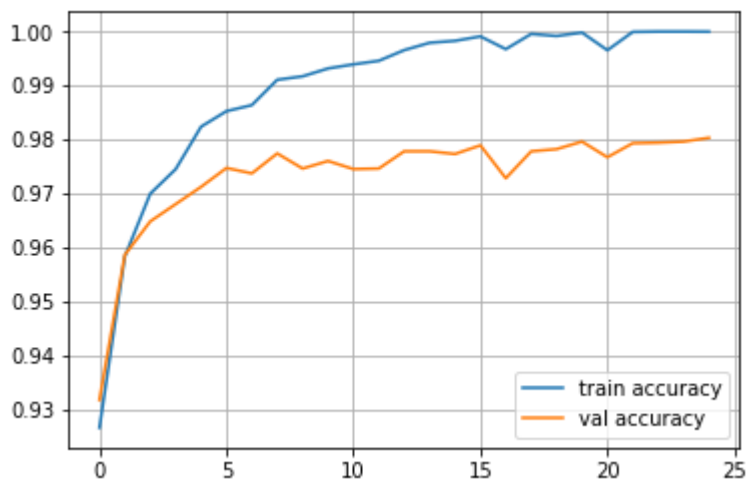
    clear_output()
    print("Epoch",epoch)
    print("Train accuracy:",train_log[-1])
    print("Val accuracy:",val_log[-1])
    plt.plot(train_log,label='train accuracy')
    plt.plot(val_log,label='val accuracy')
    plt.legend(loc='best')
    plt.grid()
    plt.show()

```

Epoch 24

Train accuracy: 0.99998

Val accuracy: 0.9803



Peer-reviewed assignment

Congratulations, you managed to get this far! There is just one quest left undone, and this time you'll get to choose what to do.

Option I: initialization

- Implement Dense layer with Xavier initialization as explained [here](http://bit.ly/2vTlmaJ) (<http://bit.ly/2vTlmaJ>).

To pass this assignment, you must conduct an experiment showing how xavier initialization compares to default initialization on deep networks (5+ layers).

Option II: regularization

- Implement a version of Dense layer with L2 regularization penalty: when updating Dense Layer weights, adjust gradients to minimize

$$Loss = Crossentropy + \alpha \cdot \sum_i w_i^2$$

To pass this assignment, you must conduct an experiment showing if regularization mitigates overfitting in case of abundantly large number of neurons. Consider tuning α for better results.

Option III: optimization

- Implement a version of Dense layer that uses momentum/rmsprop or whatever method worked best for you last time.

Most of those methods require persistent parameters like momentum direction or moving average grad norm, but you can easily store those params inside your layers.

To pass this assignment, you must conduct an experiment showing how your chosen method performs compared to vanilla SGD.

General remarks

Please read the peer-review guidelines before starting this part of the assignment.

In short, a good solution is one that:

- is based on this notebook
- runs in the default course environment with Run All
- its code doesn't cause spontaneous eye bleeding
- its report is easy to read.

Formally we can't ban you from writing boring reports, but if you bored your reviewer to death, there's no one left alive to give you the grade you want.

Bonus assignments

As a bonus assignment (no points, just swag), consider implementing Batch Normalization ([guide](https://gab41.lab41.org/batch-normalization-what-the-hey-d480039a9e3b) (<https://gab41.lab41.org/batch-normalization-what-the-hey-d480039a9e3b>)) or Dropout ([guide](https://medium.com/@amarbudhiraja/https-medium-com-amarbudhiraja-learning-less-to-learn-better-dropout-in-deep-machine-learning-74334da4bfc5) (<https://medium.com/@amarbudhiraja/https-medium-com-amarbudhiraja-learning-less-to-learn-better-dropout-in-deep-machine-learning-74334da4bfc5>)). Note, however, that those "layers" behave differently when training and when predicting on test set.

- Dropout:
 - During training: drop units randomly with probability p and multiply everything by $1/(1-p)$
 - During final prediction: do nothing; pretend there's no dropout
- Batch normalization
 - During training, it subtracts mean-over-batch and divides by std-over-batch and updates mean and variance.
 - During final prediction, it uses accumulated mean and variance.

In [19]:

```

# Implementation of Regularization, L^2 regularization

# First modified the Dense Layer to have regularization
# reg in initial is the coefficient of regularization
class Dense_reg(Layer):
    def __init__(self, input_units, output_units, learning_rate=0.1, reg=0.001):
        """
        A dense layer is a layer which performs a learned affine transformation:
         $f(x) = \langle W * x \rangle + b$ 
        """
        self.learning_rate = learning_rate
        self.reg = reg

        # initialize weights with small random numbers. We use normal initialization,
        # but surely there is something better. Try this once you got it working:
        # http://bit.ly/2vTlmaJ
        self.weights = np.random.randn(input_units, output_units) * 0.01
        self.biases = np.zeros(output_units)

    def forward(self, input):
        """
        Perform an affine transformation:
         $f(x) = \langle W * x \rangle + b$ 

        input shape: [batch, input_units]
        output shape: [batch, output units]
        """
        return np.dot(input, self.weights) + self.biases

    def backward(self, input, grad_output):
        # compute  $df / dx = df / ddense * ddense / dx$ 
        # where  $ddense / dx = weights^T$ 
        # <your code here>
        grad_input = np.dot(grad_output, self.weights.T)

        # compute gradient w.r.t. weights and biases
        grad_weights = np.dot(input.T, grad_output) / input.shape[0] + self.reg * self.weights
        # <your code here>
        grad_biases = grad_output.mean(axis=0) # <your code here>

        assert grad_weights.shape == self.weights.shape and grad_biases.shape == self.biases.shape

        # Here we perform stochastic gradient descent step.
        # later on, you can try replacing that with something better.
        self.weights = self.weights - self.learning_rate * grad_weights
        self.biases = self.biases - self.learning_rate * grad_biases

        return grad_input

```

In [20]:

```
# Define the new network
network_reg = []
network_reg.append(Dense_reg(X_train.shape[1],100))
network_reg.append(ReLU())
network_reg.append(Dense_reg(100,200))
network_reg.append(ReLU())
network_reg.append(Dense_reg(200,10))
```

In [21]:

```
# Modified the train function to include calculation of loss with regularization

def train_reg(network,X,y,reg=0.001):
    """
    Train your network on a given batch of X and y.
    You first need to run forward to get all layer activations.
    Then you can run layer.backward going from last to first layer.

    After you called backward for all layers, all Dense layers have already made
    one gradient step.
    """

    #get layer activations
    layer_activations = forward(network,X)
    layer_inputs = [X]+layer_activations #layer_input[i] is an input for network
    logits = layer_activations[-1]

    #Compute loss and initial gradient
    loss = softmax_crossentropy_with_logits(logits,y)
    loss_grad = grad_softmax_crossentropy_with_logits(logits,y)

    #<your code: propagate gradients through the network>
    for layer_i in range(len(network))[:-1]:
        layer = network[layer_i]

        loss_grad = layer.backward(layer_inputs[layer_i],loss_grad) #grad w.r.t.
        input, also weight updates
        if isinstance(layer, Dense_reg):
            loss += 0.5 * reg * np.sum(layer.weights ** 2)

    return np.mean(loss)
```

In [22]:

```
# create empty array to hold training and testing loss
train_log_reg = []
val_log_reg = []
```


In [23]:

```
# plot
for epoch in range(30):

    for x_batch,y_batch in iterate_minibatches(X_train,y_train,batchsize=32,shuffle=True):
        train_reg(network_reg,x_batch,y_batch)

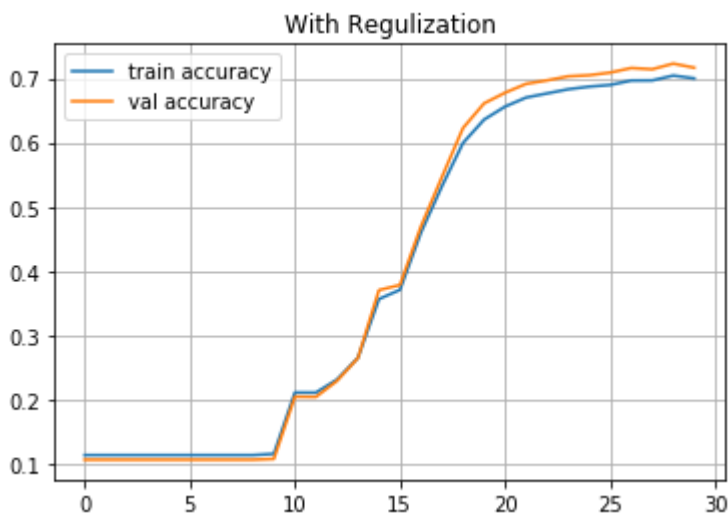
    train_log_reg.append(np.mean(predict(network_reg,X_train)==y_train))
    val_log_reg.append(np.mean(predict(network_reg,X_val)==y_val))

    clear_output()
    print("Epoch",epoch)
    print("Train accuracy:",train_log_reg[-1])
    print("Val accuracy:",val_log_reg[-1])
    plt.title("With Regularization")
    plt.plot(train_log_reg,label='train accuracy')
    plt.plot(val_log_reg,label='val accuracy')
    plt.legend(loc='best')
    plt.grid()
    plt.show()
```

Epoch 29

Train accuracy: 0.69984

Val accuracy: 0.7164



In [24]:

```
# Test set accuracy
np.mean(predict(network_reg,X_test)==y_test)
```

Out[24]:

0.71050000000000002

In [41]:

```
def weight_mag(network):  
    result = 0.0  
    for layer in network:  
        if isinstance(layer, Dense) or isinstance(layer, Dense_reg):  
            result += np.sum(np.abs(layer.weights))  
    return result
```

In [46]:

```
print(f'weight_mag(network):', weight_mag(network))
```

weight_mag(network): 4328.28169504

In [49]:

```
print(f'weight_mag(network):', weight_mag(network_reg))
```

weight_mag(network): 628.912075242

In []: