

## Import Tools

```
import os
import json
import random
import numpy
import random
import flearn
import warnings
import tensorflow
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sb
import plotly.js as js
import plotly.express as px
import matplotlib.pyplot as plt
from sklearn.preprocessing import LabelEncoder
from sklearn.metrics import confusion_matrix, classification_report, accuracy_score
#used as a shortcut to encode
encoder = LabelEncoder()
#reduces amount of warning errors on my mac
warnings.filterwarnings("ignore")
#available that calls nltk function for stemming
stemmer = LancasterStemmer
#working on tensorflow environment
#turned on to process on H2 chip architecture
os.environ['KMP_DUPLICATE_LIB_OK']='True'

WARNING:tensorflow:From /opt/anaconda3/envs/itf/lib/python3.8/site-packages/tensorflow_core/python/compat/v2_compat.py:65: disable_resource_variables (from tensorflow.python.ops.variable_scope) is deprecated and will be removed in a future version.
Instructions for updating:
non-resource variables are not supported in the long term
```

## Load Data Set

```
In [3]: #loading pnvcdiabetes data set
df = pd.read_csv("diabetes_data.csv")
#creating a list of column names
used as utility
feature_cols = df.columns.values.tolist()
feature_cols.remove('Outcome')
```

```
In [3]: #Utility list used for holding value in pre-processing
words = []
labels = []
sample_x = []
sample_y = []
training = []
output = []

#open data file for training chatbot
with open("data.json") as file:
    data = json.load(file)

#iterate through first dimension of training data
for intent in data["intents"]:
    #move data the second dim. to now tokenize language
    for pattern in intent["patterns"]:
        #tokenize our examples per tokenmarker
        words = nltk.word_tokenize(pattern)
        #extend our list of words with token version
        words.extend(words)
        #saving corresponding word
        sample_x.append(words)
        #save the tag that corresponds of sample x
        #used for identification in training
        sample_y.append(intent["tag"])
        #if cases were tag has not been seen
        if intent["tag"] not in labels:
            #we add it as part of our labels
            labels.append(intent["tag"])
#now know how nltk to stem our tokens
words = [stemmer.stem(w.lower()) for w in words if w != "?"]
#quickly sort them
words = sorted(list(set(words)))
#sorting our labels/features
labels = sorted(labels)
#used as slicing utility
out_empty = [0] for i in range(len(labels))
#iterate through samples of tokenized words
for x, doc in enumerate(sample_x):
    #iteration of utility list
    bag = []
    #processing each case inside of our already stemmed examples
    words = [stemmer.stem(w.lower()) for w in doc]
    #where we assign assistance of tokens
    #every similar to one hot encoding
    for w in words:
        if w in words:
            bag.append(1)
        else:
            bag.append(0)
    #creating output data
    output_row = out_empty[:]
    output_row[labels.index(sample_y[x])] = 1
    #creating training data array based on encoded tokens
    training.append(bag)
    #corresponding labels for training data
    output.append(output_row)

#converting to numpy arrays
#this must be done for neural network
training = numpy.array(training)
output = numpy.array(output)
```

## Computational Model: Pre-Processing

```
In [4]: #Utility function
#Parameters(dataframe, result, feature)
#Returns dataframe grouped by relations
def find_avg(df, target, var):
    #Returns dataframe grouped by feature all non-empty values
    temp = df[df[target].notnull()]
    #median created with median on relation between feature and result
    temp = temp[[var, target]].groupby([target])[var].median().reset_index()
    return temp
```

```
In [5]: #Utility function
#Parameters(dataframe, result, feature, median of features with result False, median of features with result True)
#modifies the passed dataframe
def impute_avg(df, target, var, var_0, var_1):
    #iterate through dataframe
    for i in range(len(df)):
        #assign correspondence
        if df.loc[i, target] == 0 and df.loc[i, var] == 0:
            df.loc[i, var] = var_0
        if df.loc[i, target] == 1 and df.loc[i, var] == 0:
            df.loc[i, var] = var_1
```

```
In [7]: #Utility function
#Parameters(dataframe)
#Returns labels and corresponding vectors
def compute_relations(df):
    #Suggested points by kaggle solution
    #have been modified to fit mayo clinic guides
    #used to separate pre-diabetes metrics with high risk diabetic metrics
    #all are true unless
    df.loc[:, "N1"] = 1
    #combination of metrics
    df.loc[(df["Age"]<=30) & (df["Glucose"]<=120), "N1"] = 0
    df.loc[(df["Age"]<=30) & (df["Age"]<=45) & (df["Glucose"]<=88), "N1"] = 0
    df.loc[(df["Age"]>=60) & (df["Glucose"]<=140), "N1"] = 0
    #tracks Obesity in people over 30
    df.loc[:, "N2"] = 1
    df.loc[(df["BMI"]<=30), "N2"] = 0
    #tracks relationship between older women with high rate of pregnancies
    df.loc[:, "N3"] = 1
    df.loc[(df["Age"]<=27) & (df["Pregnancies"]<=6), "N3"] = 0
    df.loc[(df["Age"]>=60) & (df["Pregnancies"]>=8), "N3"] = 0
    #relationship between low glucose levels and bloodpressure
    #high blood pressure can be a sign of many health complications
    df.loc[:, "N4"] = 1
    df.loc[(df["Glucose"]<=105) & (df["BloodPressure"]<=80), "N4"] = 0
    df.loc[(df["Glucose"]<=105) & (df["BloodPressure"]>=83), "N4"] = 0
    #diagnose test results over 20 can be used to track bodyfat percentage
    df.loc[:, "N5"] = 1
    df.loc[(df["SkinThickness"]<=20), "N5"] = 0
    #Relationship between body mass index and bodyfat
    #high bmi and bodyfat are signs of bad health
    df.loc[:, "N6"] = 1
    df.loc[(df["BMI"]<=30) & (df["SkinThickness"]<=20), "N6"] = 0
    df.loc[(df["BMI"]>=31) & (df["SkinThickness"]<=20), "N6"] = 0
    #glucose sensitivity helps us separate high bmi patients
    df.loc[:, "N7"] = 1
    df.loc[(df["Glucose"]<=105) & (df["BMI"]<=30), "N7"] = 0
    df.loc[(df["Glucose"]<=105) & (df["BMI"]>=40), "N7"] = 0
    #Insulin levels over 200 are known to be problematic by mayo clinic
    df.loc[:, "N8"] = 1
    df.loc[(df["Insulin"]<=200), "N8"] = 0
    #medical standards suggest a diastolic blood pressure to be health only if less than 80mm Hg
    df.loc[:, "N9"] = 1
    df.loc[(df["BloodPressure"]<=80), "N9"] = 0
    #numerous pregnancies can result in high risk
    df.loc[:, "N11"] = 1
    df.loc[(df["Pregnancies"]<=4) & (df["Pregnancies"]>=8), "N11"] = 0
    #product of body mass and bodyfat
    df["N0"] = df["BMI"] * df["SkinThickness"]
    #relation between pregnancies and age
    df["N0"] = df["Pregnancies"] / df["Age"]
    #relationship between blood glucose levels and family history
    df["N10"] = df["Glucose"] / df["DiabetesPedigreeFunction"]
    #relationship between age and family history
    df["N12"] = df["Age"] * df["DiabetesPedigreeFunction"]
    #separation of bodyfat per mass index values
    df.loc[:, "N15"] = 1
    df.loc[(df["N0"]<=434), "N15"] = 0
    #separate labels over 200 are known to be problematic by mayo clinic
    y = pd.DataFrame(df["Outcome"])
    x = df.drop("Outcome", axis=1)
    #encode outcomes
    for col in y.columns:
        y[col] = encoder.fit_transform(y[col])
    y = y["Outcome"]
    cols = ["N1", "N2", "N3", "N4", "N5", "N6", "N7", "N8", "N9", "N10", "N11"]
    #new col metrics encoding
    for col in cols:
        x[col] = encoder.fit_transform(x[col])
    return labels and vectors
return x, y
```

## Chatbot Model: Pre-Processing

```
In [8]: #Utility function
#Parameters(user input, words list)
#Returns numpy array of possibility
def bag_of_words(inp, words):
    #create bag of words
    bag = [0 for _ in range(len(words))]
    #tokenize the words
    inp_words = nltk.word_tokenize(inp)
    #stemming words
    inp_words = [stemmer.stem(word.lower()) for word in inp_words]
    #now not encoding for words
    for i, w in enumerate(words):
        if w == inp_words[i]:
            bag[i] = 1
    return numpy.array(bag)
```

```
In [9]: #Functional Function
#Defines Architecture of Neural Net.
#Returns a deep neural network structure
def build_nn():
    #input layer defines shape is [None, 81]
    net = tflearn.input_data(shape=(None, len(training[0])))
    #hidden layer with 8 neurons
    net = tflearn.fully_connected(net, 8)
    #hidden layer with 8 neurons
    net = tflearn.fully_connected(net, 8)
    #activation layer used for probability of correctness
    net = tflearn.fully_connected(net, len(output[0]), activation="softmax")
    #output layer based on probability
    net = tflearn.regression(net)
    #define as deep neural network
    model = tflearn.DNN(net)
    return model
```

```
In [10]: #Functional
#Trains and saves model to current directory
#Returns the trained model
def train_nn():
    #calls the model to be build
    model = build_nn()
    #trains the model based on parameters(arrays training data, labels of training data, num times training data seen, display training iterations)
    model.fit(training, output, n_epoch=1000, batch_size=8, show_metric=True)
    #save model to dir
    model.save("chatbot_diabetes2.tflearn")
    return model
```

```
In [11]: #Utility function
#Parameters(dataframe)
#Returns diagnosis of patient
def format_user_df(user_df):
    #Suggested points by kaggle solution
    #have been modified to fit mayo clinic guides
    #used to separate pre-diabetes metrics with high risk diabetic metrics
    #all are true unless
    df.loc[:, "N1"] = 1
    #combination of metrics
    df.loc[(df["Age"]<=30) & (df["Glucose"]<=120), "N1"] = 0
    df.loc[(df["Age"]<=30) & (df["Age"]<=45) & (df["Glucose"]<=88), "N1"] = 0
    df.loc[(df["Age"]>=60) & (df["Glucose"]<=140), "N1"] = 0
    #tracks Obesity in people over 30
    df.loc[:, "N2"] = 1
    df.loc[(df["BMI"]<=30), "N2"] = 0
    #tracks relationship between older women with high rate of pregnancies
    df.loc[:, "N3"] = 1
    df.loc[(df["Age"]<=27) & (df["Pregnancies"]<=6), "N3"] = 0
    df.loc[(df["Age"]>=60) & (df["Pregnancies"]>=8), "N3"] = 0
    #relationship between low glucose levels and bloodpressure
    #high blood pressure can be a sign of many health complications
    df.loc[:, "N4"] = 1
    df.loc[(df["Glucose"]<=105) & (df["BloodPressure"]<=80), "N4"] = 0
    df.loc[(df["Glucose"]<=105) & (df["BloodPressure"]>=83), "N4"] = 0
    #diagnose test results over 20 can be used to track bodyfat percentage
    df.loc[:, "N5"] = 1
    df.loc[(df["SkinThickness"]<=20), "N5"] = 0
    #Relationship between body mass index and bodyfat
    #high bmi and bodyfat are signs of bad health
    df.loc[:, "N6"] = 1
    df.loc[(df["BMI"]<=30) & (df["SkinThickness"]<=20), "N6"] = 0
    df.loc[(df["BMI"]>=31) & (df["SkinThickness"]<=20), "N6"] = 0
    #glucose sensitivity helps us separate high bmi patients
    df.loc[:, "N7"] = 1
    df.loc[(df["Glucose"]<=105) & (df["BMI"]<=30), "N7"] = 0
    df.loc[(df["Glucose"]<=105) & (df["BMI"]>=40), "N7"] = 0
    #Insulin levels over 200 are known to be problematic by mayo clinic
    df.loc[:, "N8"] = 1
    df.loc[(df["Insulin"]<=200), "N8"] = 0
    #medical standards suggest a diastolic blood pressure to be health only if less than 80mm Hg
    df.loc[:, "N9"] = 1
    df.loc[(df["BloodPressure"]<=80), "N9"] = 0
    #numerous pregnancies can result in high risk
    df.loc[:, "N11"] = 1
    df.loc[(df["Pregnancies"]<=4) & (df["Pregnancies"]>=8), "N11"] = 0
    #product of body mass and bodyfat
    df["N0"] = df["BMI"] * df["SkinThickness"]
    #relation between pregnancies and age
    df["N0"] = df["Pregnancies"] / df["Age"]
    #relationship between blood glucose levels and family history
    df["N10"] = df["Glucose"] / df["DiabetesPedigreeFunction"]
    #relationship between age and family history
    user_df["N13"] = 0
    #separation of bodyfat per mass index values
    df.loc[:, "N15"] = 1
    df.loc[(df["N0"]<=434), "N15"] = 0
    #use machine learning engine to make a prediction of diagnosis
    y_pred = bagging.predict(user_df)
    #return diagnosis
    return y_pred
```

```
In [12]: #Function: Live Interaction
#Must be called with model imported
#Runs by user input
def chat():
    #utility list
    pregnancies = []
    glucose = []
    bloodpressure = []
    skintthickness = []
    insulin = []
    bmi = []
    pedigree = []
    age = []
    #features to save extracted features
    user_features = []
    #information for dataframe below
    marker_means = {'Number of times pregnant':
        'Plasma glucose concentration over 2 hours in an oral glucose tolerance test',
        'Diastolic blood pressure (mm Hg)',
        'Triceps skin fold thickness (mm)',
        '2-Hour serum insulin (mu U/ml)',
        'Please view chart below',
        'Function which scores likelihood of diabetes based on family history',
        'Age (years)'}
    display information to user, may help guide
    df.bio.markers = pd.DataFrame(list(zip(feature_cols, marker_means)), columns = ['Bio Markers', 'How to measure'])
    # greetings
    print("Hello! this is SuperBot, I can people determine if they are in risk of diabetes.")
    print("All you have to do, is provide me with the following information and type done for me to calculate:")
    display(df.bio.markers)
    display(IMAGEFILENAME.bmi.jpg)
    #conversation of chatbot with user
    while True:
        #user input variable
        inp = input("You: ")
        #check for terminating condition
        if inp.lower() == "done":
            break
        #prediction of context by DNN
        results = model.predict(bag_of_words(inp, words))
        results_index = numpy.argmax(results)
        tag = labels[results_index]
        #prediction extractions
        if tag == "Pregnancies":
            if len(pregnancies) == 0:
                pregnancies.append(inp)
            else:
                pregnancies[0] = inp
        elif tag == "Glucose":
            if len(glucose) == 0:
                glucose.append(inp)
            else:
                glucose[0] = inp
        elif tag == "BloodPressure":
            if len(bloodpressure) == 0:
                bloodpressure.append(inp)
            else:
                bloodpressure[0] = inp
        elif tag == "SkinThickness":
            if len(skintthickness) == 0:
                skintthickness.append(inp)
            else:
                skintthickness[0] = inp
        elif tag == "Insulin":
            if len(insulin) == 0:
                insulin.append(inp)
            else:
                insulin[0] = inp
        elif tag == "BMI":
            if len(bmi) == 0:
                bmi.append(inp)
            else:
                bmi[0] = inp
        elif tag == "DiabetesPedigreeFunction":
            if len(pedigree) == 0:
                pedigree.append(inp)
            else:
                pedigree[0] = inp
        elif tag == "Age":
            if len(age) == 0:
                age.append(inp)
            else:
                age[0] = inp
        else:
            pass
        #for tag in intents:
            if tag == inp:
                responses = [r[responses]]
                #message to user about predicted context
                print(random.choice(responses))
                print("If you're not referring to your (s), please rephrase your statement, otherwise ignore this message".format(tag))
        #creating a combined feature list
        features = pregnancies+glucose+bloodpressure+skintthickness+insulin+bmi+pedigree+age
        #one hot encoding for user input metrics
        for feature in features:
            temp = [int(i) for i in feature.split() if i.isdigit()]
            if len(temp) == 0:
                temp = [0]
            else:
                temp = temp[:1]
        #print(temp)
        user_features.append(temp)
        #dataframe of user information
        user_df = pd.DataFrame([user_features], columns = feature_cols)
        #prediction by computational model
        ML_result = format_user_df(user_df[0])
        #display results
        if ML_result == 0:
            print("You are not likely to have diabetes.")
        else:
            print("You are in high risk of diabetes, please contact your medical provider")
```

## Train and Validate Computational Model

```
In [13]: d = find_avg(df, "Outcome", "Insulin")
impute_avg(df, "Outcome", "Insulin", d.loc[0, "Insulin"], d.loc[1, "Insulin"])
x, y = compute_relations(df)
x_train, x_validate, y_train, y_validate = train_test_split(x, y, test_size = .2, random_state=42)
bagging = BaggingClassifier(
    DecisionTreeClassifier(random_state=42),
    n_estimators=500, max_samples=100,
    bootstrap=True, n_jobs=-1, random_state=42,
    oob_score=True)
bagging.fit(x_train, y_train)
print(classification_report(bagging.predict(x_validate), (y_validate)))
```

```
precision    recall  f1-score   support

0   0.94   0.91   0.93   182
1   0.84   0.88   0.86   52

accuracy   0.89   0.90   0.89   154
macro avg   0.89   0.90   0.90   154
```

## Test Chatbot Model

```
In [14]: model = train_nn()
```

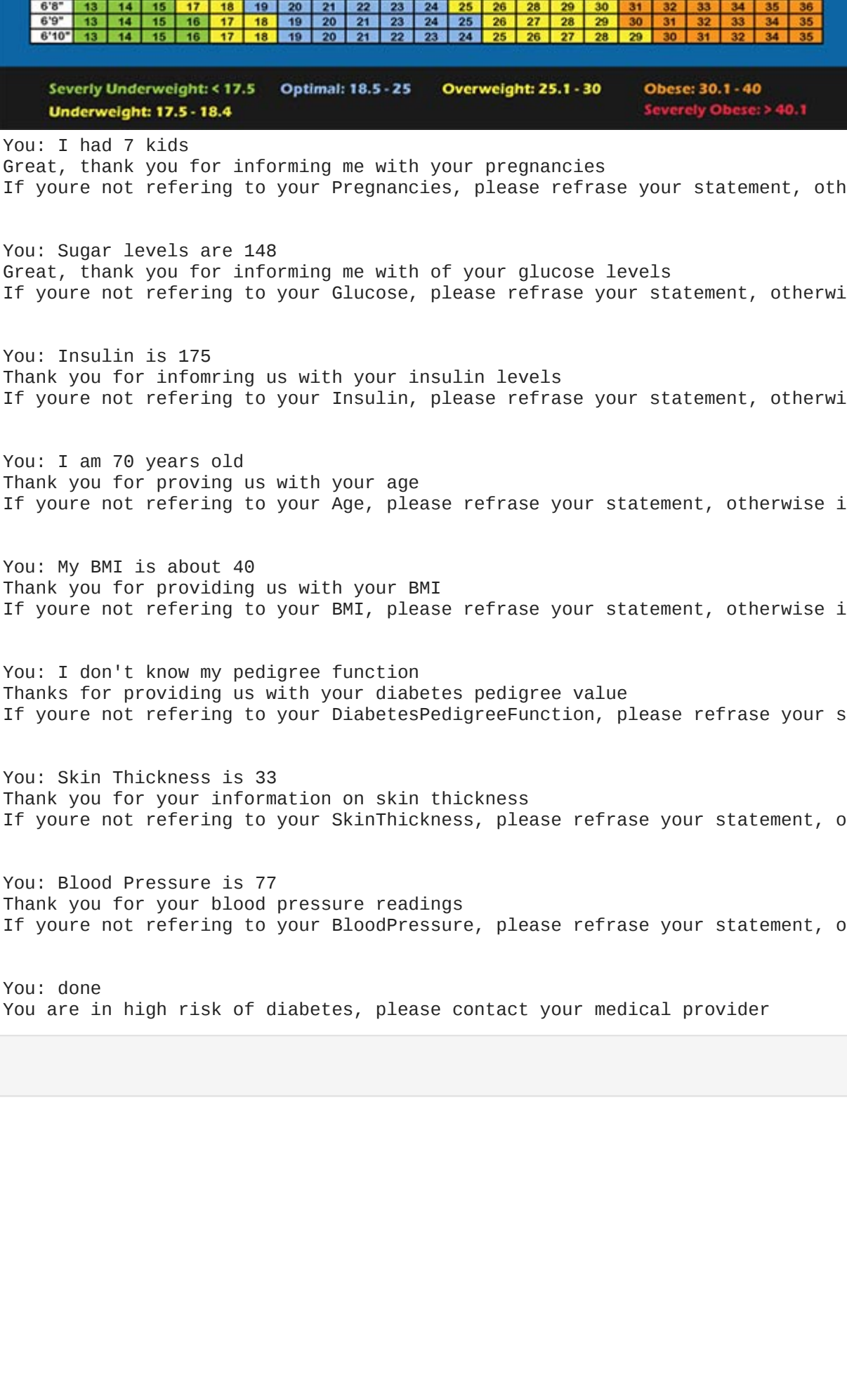
```
Training Step: 6999 | total loss: 0.00913 | time: 0.016s
| Adam | epoch: 1000 | loss: 0.00913 - acc: 0.9997 -- iter: 48/53
Training Step: 7000 | total loss: 0.00885 | time: 0.018s
| Adam | epoch: 1000 | loss: 0.00885 - acc: 0.9997 -- iter: 53/53
```

INFO:tensorflow:Users/Alex/Diabetes Chatbot/diabetes.tflearn is not in all\_model\_checkpoint\_paths. Manually adding it.

```
In [16]: chat()
```

Hello! this is SuperBot, I can people determine if they are in risk of diabetes.  
All you have to do, is provide me with the following information and type done for me to calculate:

	Bio Markers	How to measure
0	Pregnancies	Number of times pregnant
1	Glucose	Plasma glucose concentration over 2 hours in a...
2	BloodPressure	Diastolic blood pressure (mm Hg)
3	SkinThickness	Triceps skin fold thickness (mm)
4	Insulin	2-Hour serum insulin (mu U/ml)
5	BMI	Please view chart below
6	DiabetesPedigreeFunction	function which scores likelihood of diabetes b...
7	Age	Age (years)



You: I had 7 kids  
Great thank you for informing me with your pregnancies  
If you're not referring to your Pregnancies, please rephrase your statement, otherwise ignore this message

You: Sugar levels are 148  
Great thank you for informing me with of your glucose levels  
If you're not referring to your Glucose, please rephrase your statement, otherwise ignore this message

You: Insulin is 176  
Thank you for informing us with your insulin levels  
If you're not referring to your Insulin, please rephrase your statement, otherwise ignore this message

You: I am 78 years old  
Thank you for providing us with your age  
If you're not referring to your Age, please rephrase your statement, otherwise ignore this message

You: My BMI is about 40  
Thank you for providing us with your BMI  
If you're not referring to your BMI, please rephrase your statement, otherwise ignore this message

You: I don't know my pedigree function  
Thanks for providing us with your diabetes pedigree value  
If you're not referring to your DiabetesPedigreeFunction, please rephrase your statement, otherwise ignore this message

You: Skin Thickness is 33  
Thank you for your information on skin thickness  
If you're not referring to your SkinThickness, please rephrase your statement, otherwise ignore this message

You: Blood Pressure is 77  
Thank you for your blood pressure readings  
If you're not referring to your BloodPressure, please rephrase your statement, otherwise ignore this message

You: done  
You are in high risk of diabetes, please contact your medical provider

```
In [ ]:
```