


Deep Learning for Life Sciences: Bayesian Deep Learning

Nikolay Oskolkov, NBIS SciLifeLab
11.12.2020



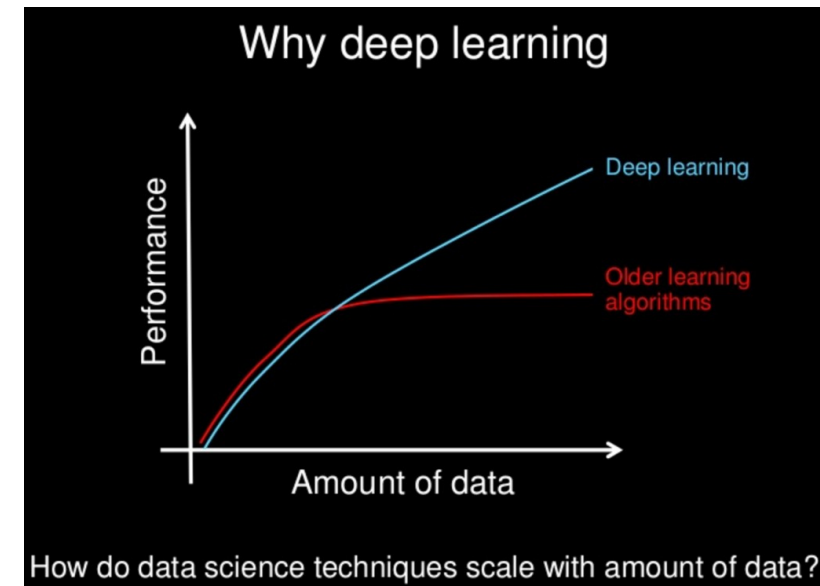
- Apply to real Life Science projects (NGS: tabular data)
- Apply only if Deep Learning better than simpler methods



Occam's Razor: No more things should be presumed to exist than are absolutely necessary, i.e., the fewer assumptions an explanation of a phenomenon depends on, the better the explanation.

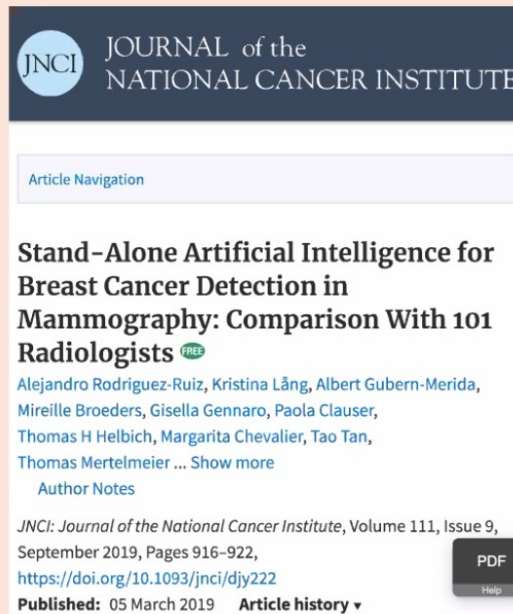
(William of Occam)

izquotes.com



Why don't neural networks always work?

Rodriguez et al compared AI with 101 radiologists – AI was as good as radiologists



JNCI JOURNAL of the NATIONAL CANCER INSTITUTE

Article Navigation

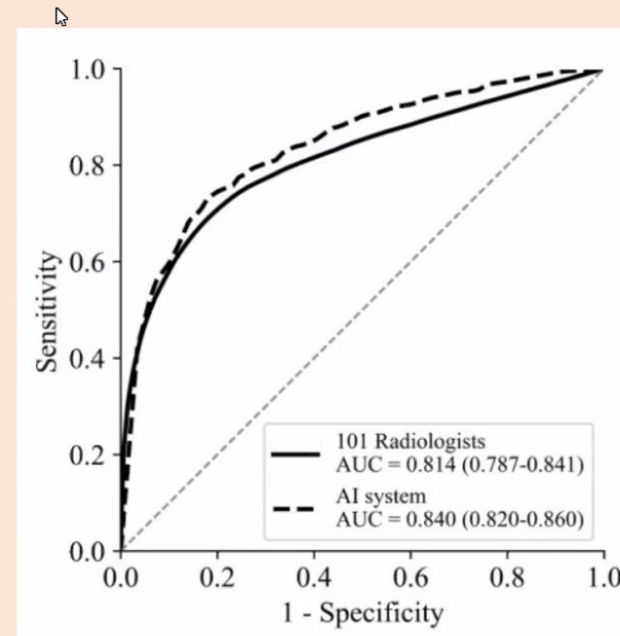
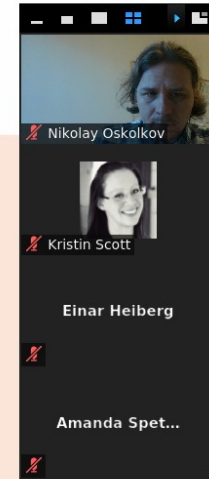
Stand-Alone Artificial Intelligence for Breast Cancer Detection in Mammography: Comparison With 101 Radiologists FREE

Alejandro Rodriguez-Ruiz, Kristina Lång, Albert Gubern-Merida, Mireille Broeders, Gisella Gennaro, Paola Clauser, Thomas H Helbich, Margarita Chevalier, Tao Tan, Thomas Mertelmeier ... Show more

Author Notes

JNCI: *Journal of the National Cancer Institute*, Volume 111, Issue 9, September 2019, Pages 916–922, <https://doi.org/10.1093/jnci/djy222>

Published: 05 March 2019 Article history

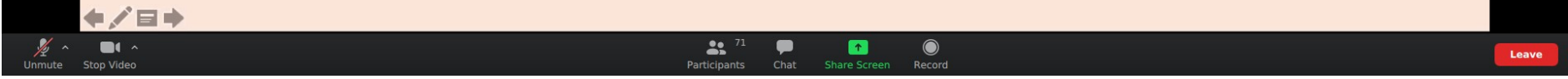



Nikolay Oskolkov

Kristin Scott

Einar Heiberg

Amanda Spet...



Unmute Stop Video Participants 71 Chat Share Screen Record Leave



Why do you compare AI against radiologists?
You should compare it against simpler models

Bayesianism



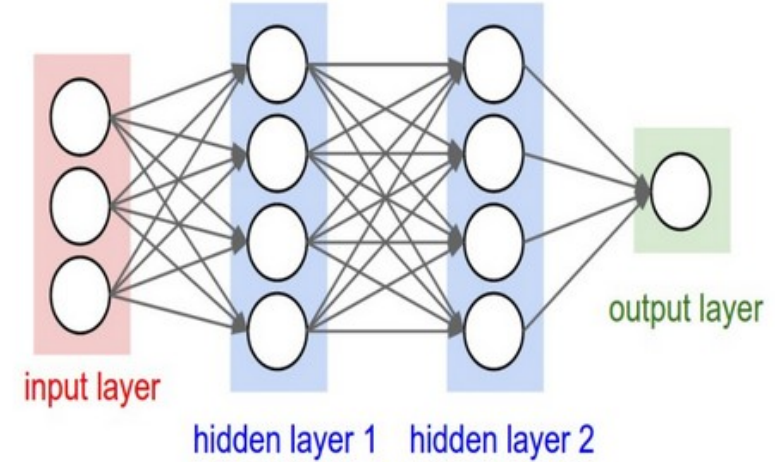
$$P \gg N$$

Frequentism



$$P \sim N$$

Deep Learning



$$P \ll N$$



Amount of Data

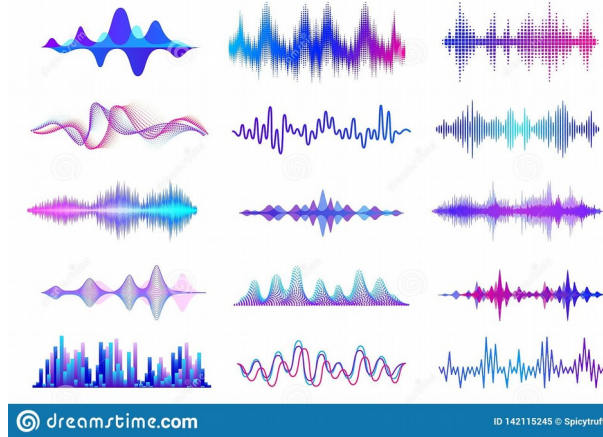
Deep Learning is a yet another tool



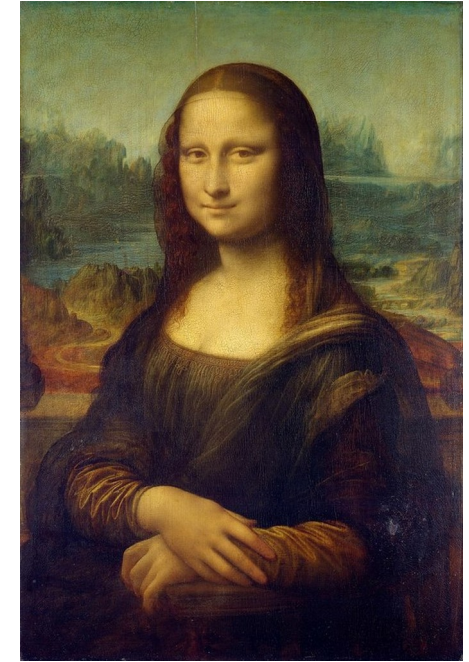
Comparison is important: If you do not compare, your neural network is the best

Tabular

Sound



Image



Text

Editing Wikipedia articles on Medicine



Editing Wikipedia can be daunting for novices, especially if you're contributing to Wikipedia for the first time as a class assignment. This guide is designed to assist students who have been assigned to contribute biomedical-related content to Wikipedia. Here's what other editors will expect you to know.

Be accurate
You're editing a resource millions of people use to make medical decisions, so it's vitally important to be accurate. Wikipedia is used more for medical information than the websites for WebMD, NIH, and the WHO. But with great power comes great responsibility!

Understand the guidelines
Wikipedia editors in the medicine area have developed additional guidelines to ensure that the content on Wikipedia is medically sound. Take extra time to read and understand these guidelines. When you edit an article, ensure your changes meet these special requirements. If not, your work is likely to be undone by other editors as they clean up after you. That takes valuable volunteer time away from creating content. If you're not comfortable working under these guidelines, talk to your instructor about an alternative off-wiki assignment.



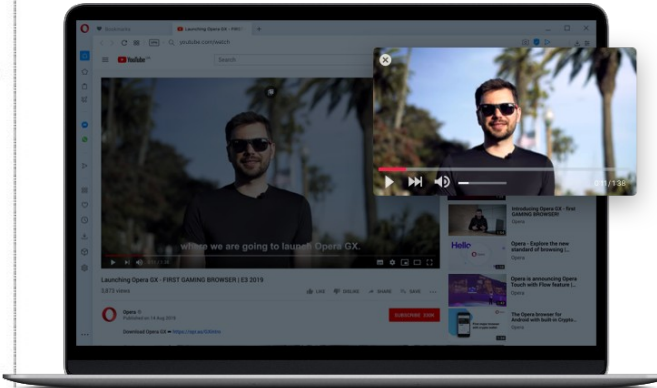
Engage with editors
Part of the Wikipedia experience is receiving and responding to feedback from other editors. Do not submit your content on the last day, then leave Wikipedia! Real human volunteers from the Wikipedia community will likely read and respond to it, and it would be polite for you to acknowledge the time they volunteer to polish your work! Everything submitted to Wikipedia is reviewed by multiple, real humans! You may not get a comment, but if you do, please acknowledge it.

Watch out for close paraphrasing
Plagiarizing or close paraphrasing is never okay on Wikipedia and is a violation of your university's academic honor code. It's even worse on Wikipedia, as valuable volunteer time that could be used to create good content is instead used to clean up plagiarized work.

Note that even educational materials from organizations like the WHO and abstracts of articles in PubMed are under copyright and cannot be copied. Write them in your own words whenever possible. If you aren't clear on what close paraphrasing is, visit your university's writing center.

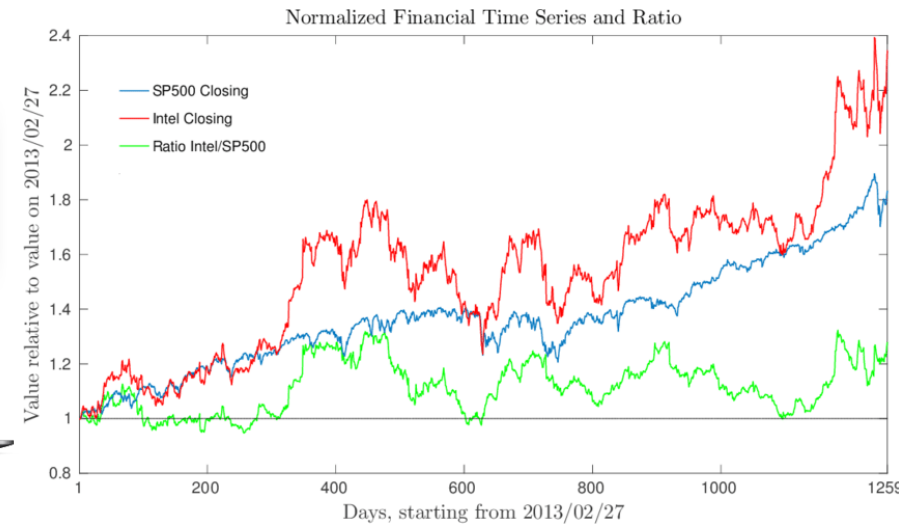
Scared? Don't be!
Everybody on Wikipedia wants to make the best encyclopedia they can. Take the time to understand the rules, and soon you'll be contributing to a valuable resource you use on a daily basis!

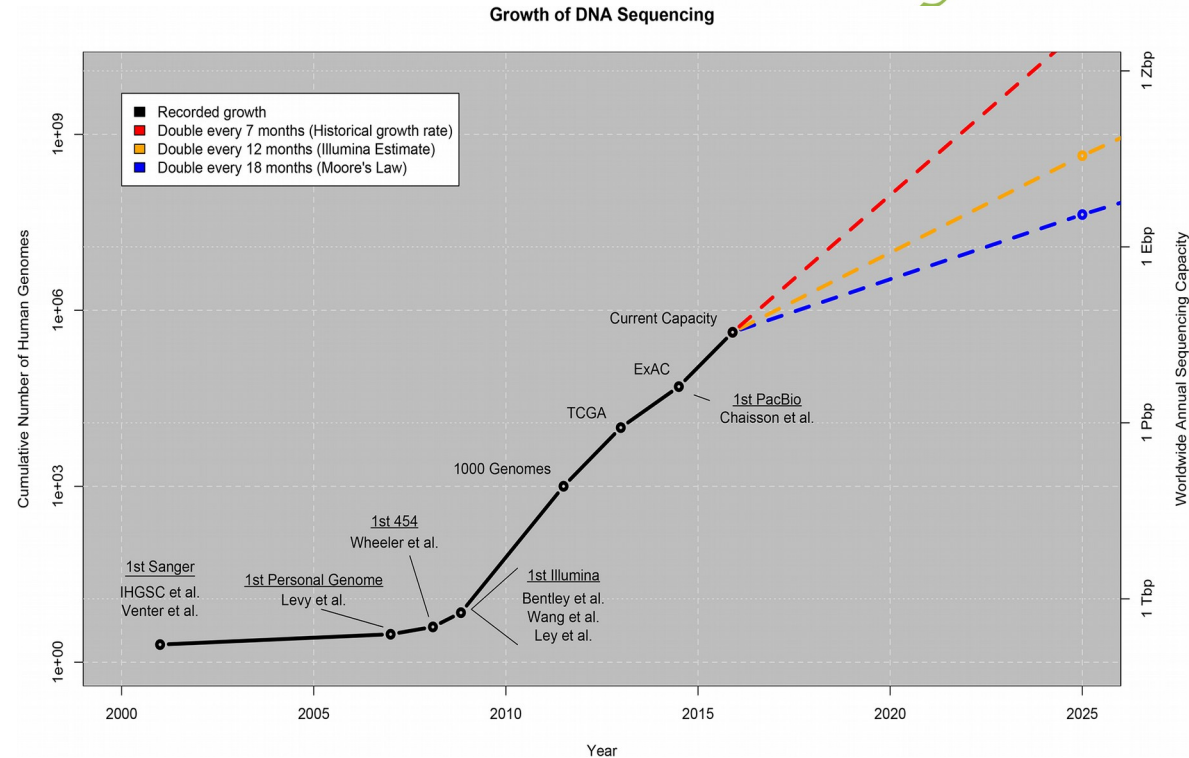
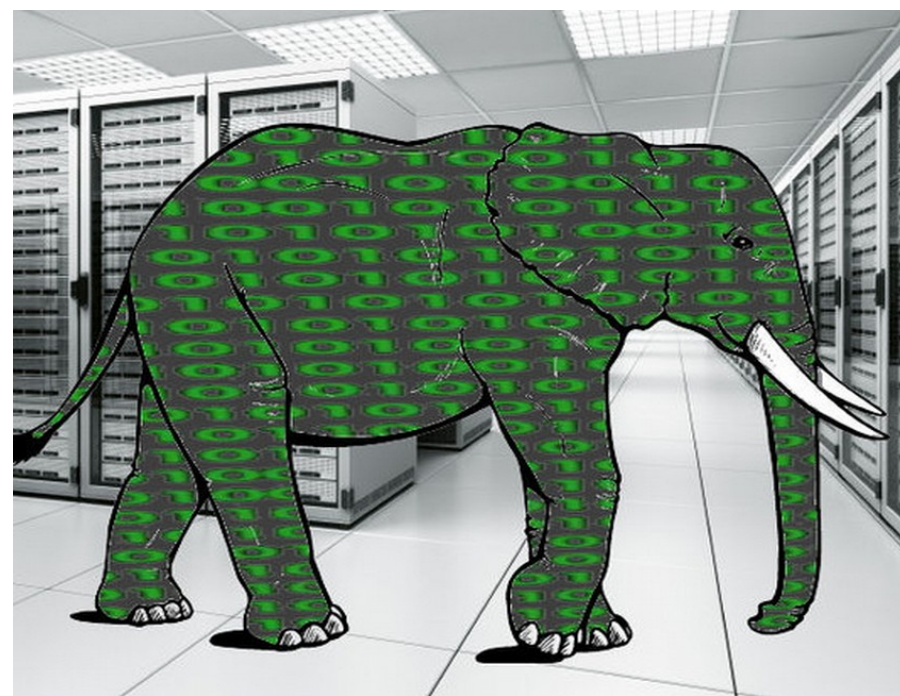
Video



DATA

Time Series



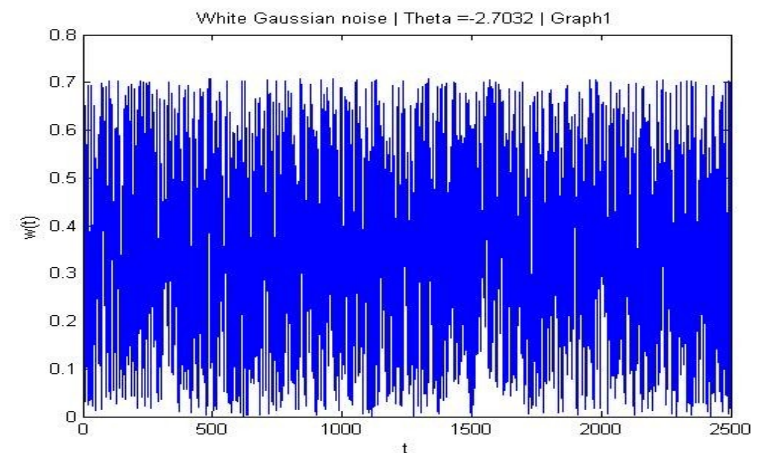


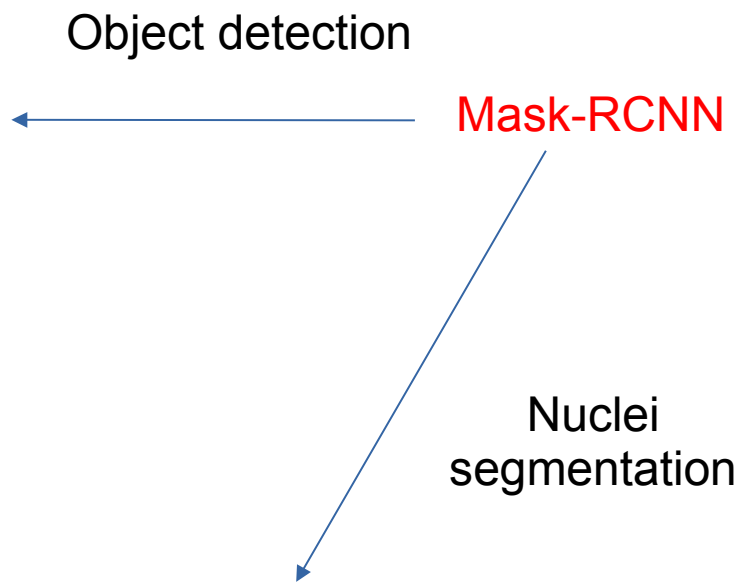
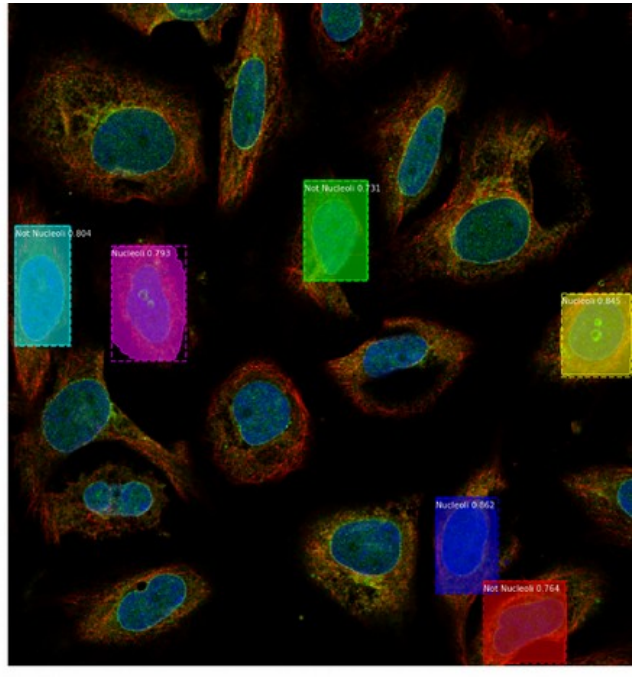
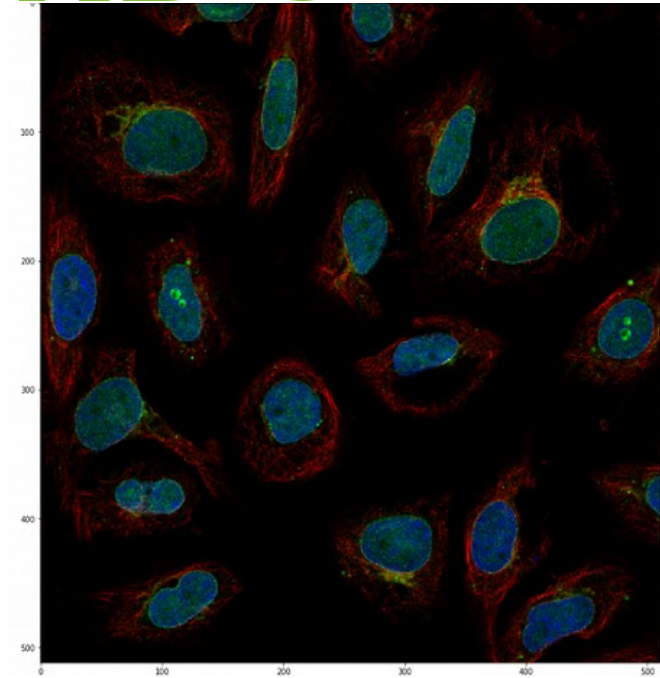
Possible Big Data in Life Sciences:

- Microscopy Imaging
- Single Cell Omics
- Metagenomics (possibly)
- Genomics (sequence is an observation)

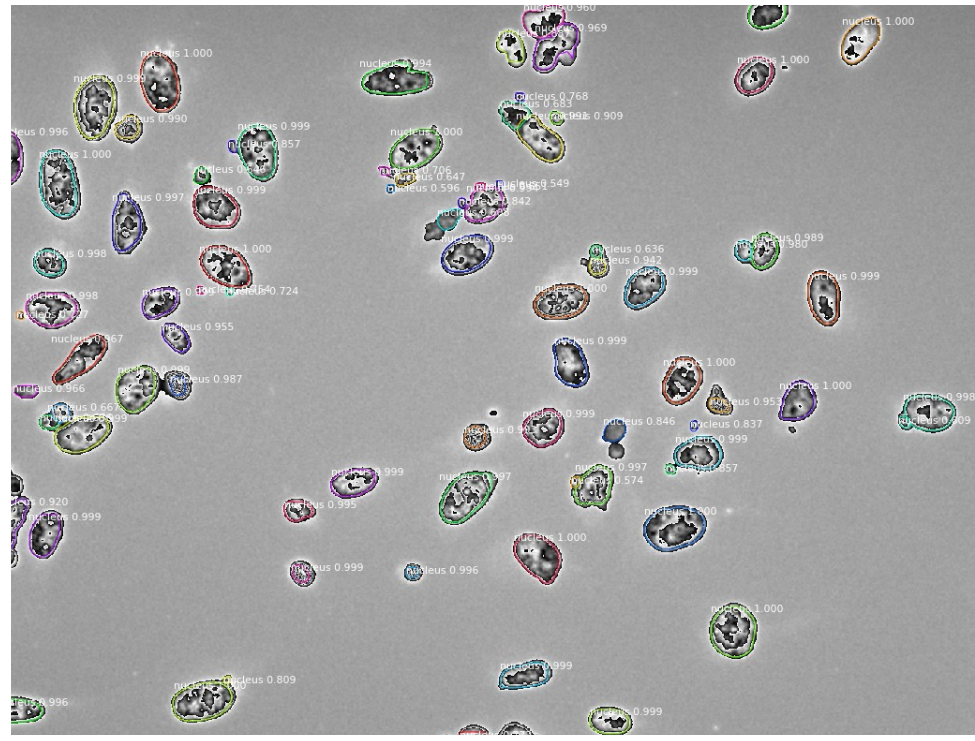
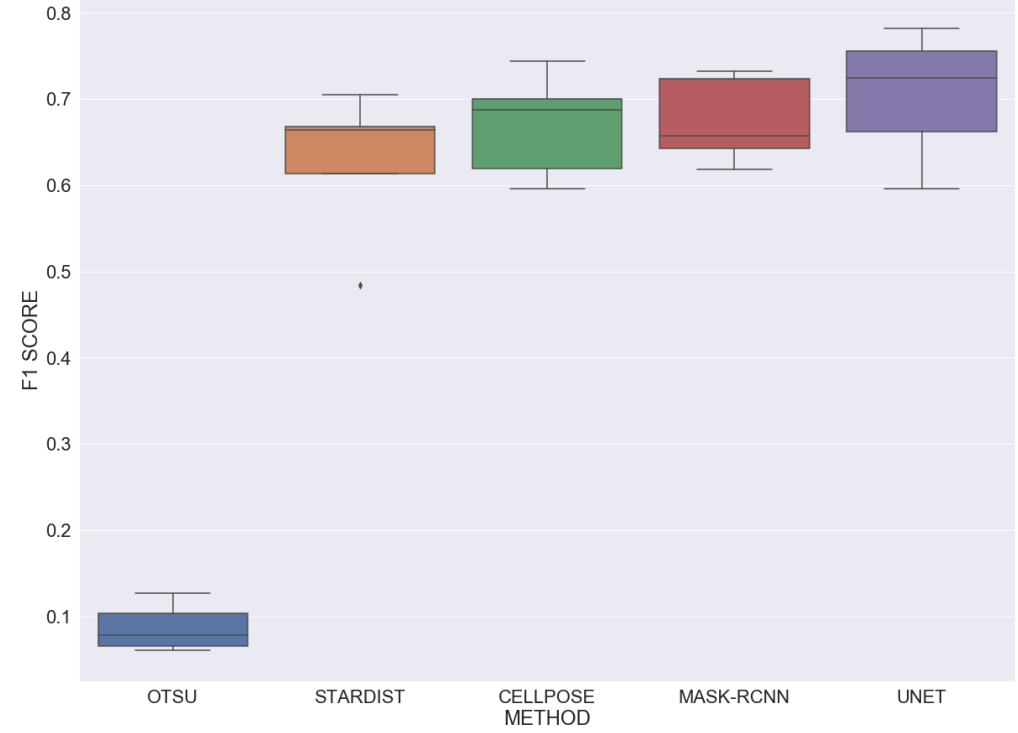
I have 500 TB of data on my disk, this is big.

I have Big Data, I want to run Deep Learning on my Big Data

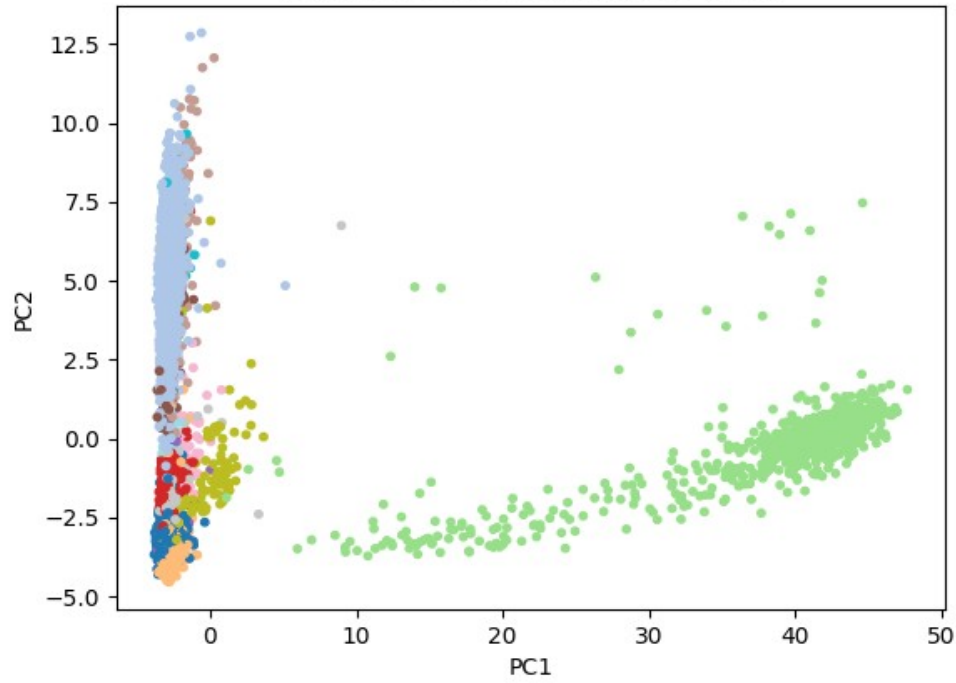




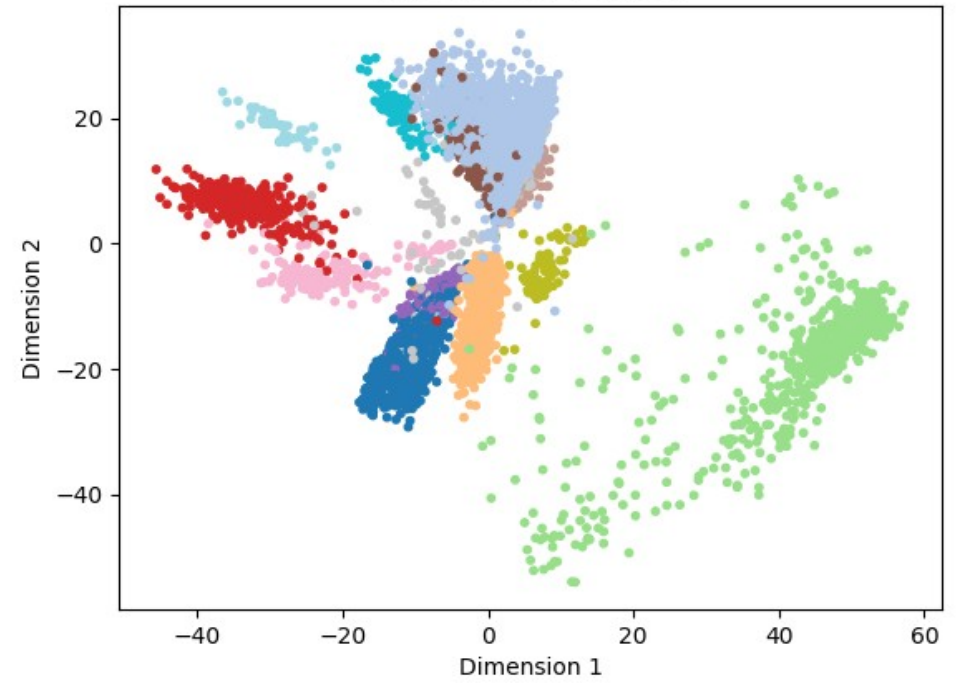
COMPARING DIFFERENT NUCLEI SEGMENTATION STRATEGIES FOR BBBC039 CELLS: IOU = 0.9



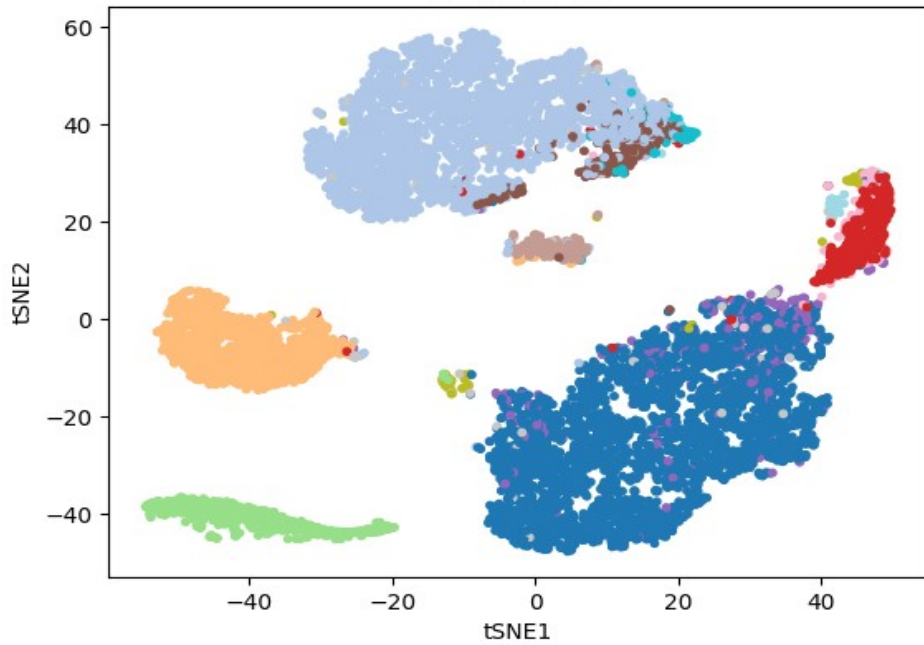
Principal Component Analysis (PCA)



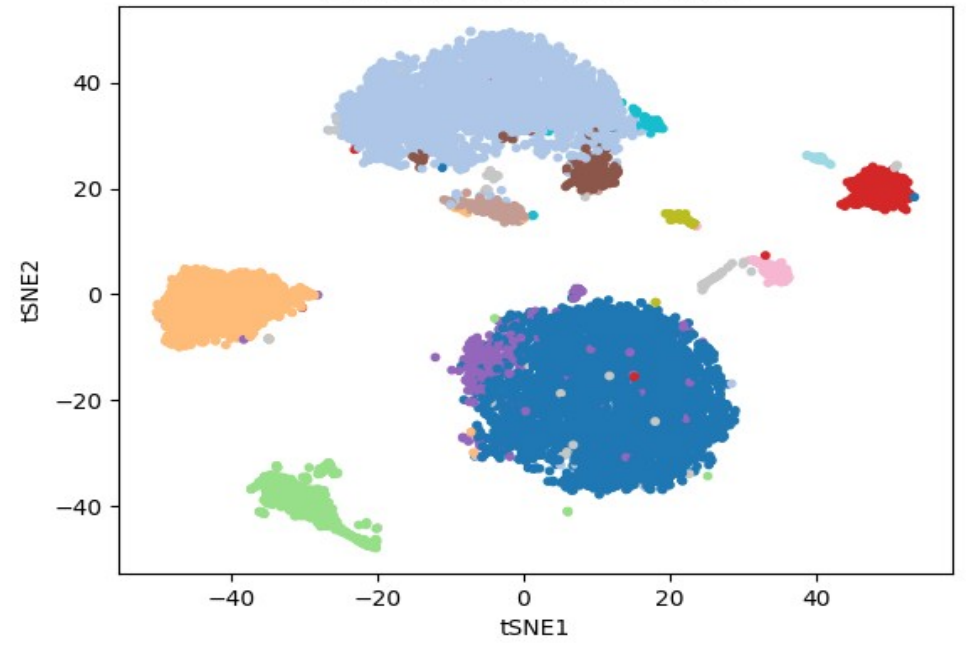
Autoencoder: 8 Layers

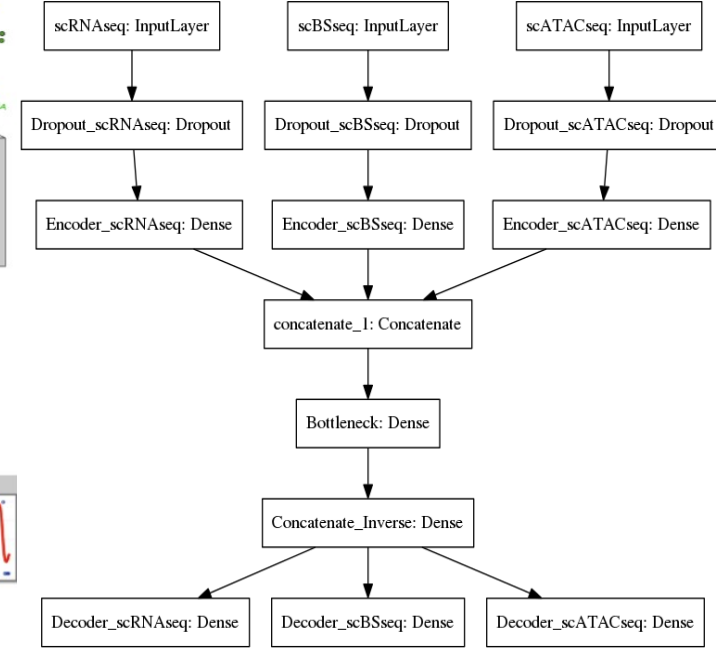
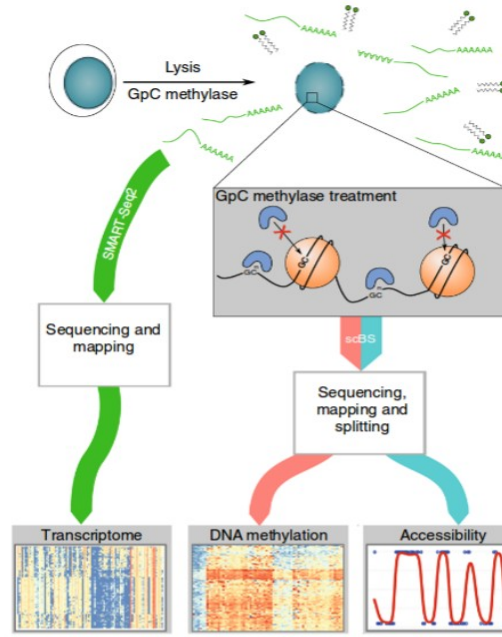
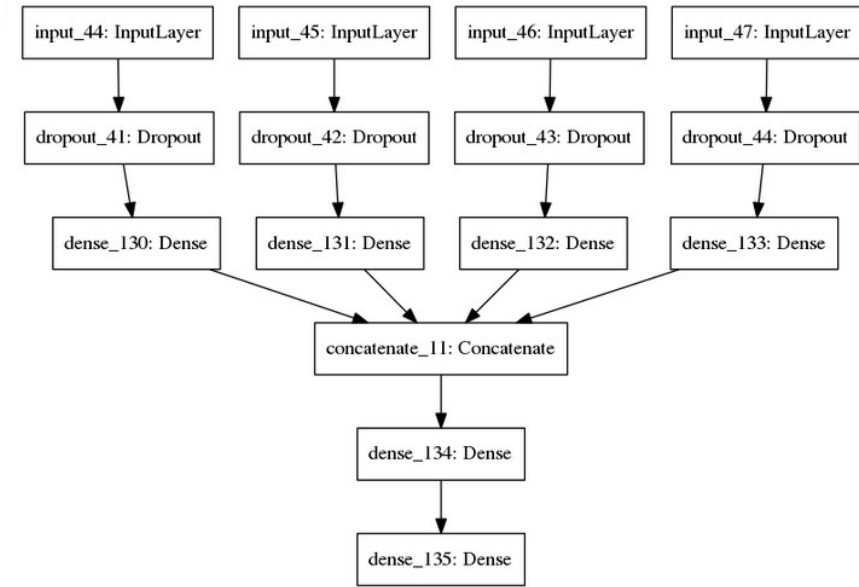


tSNE on PCA

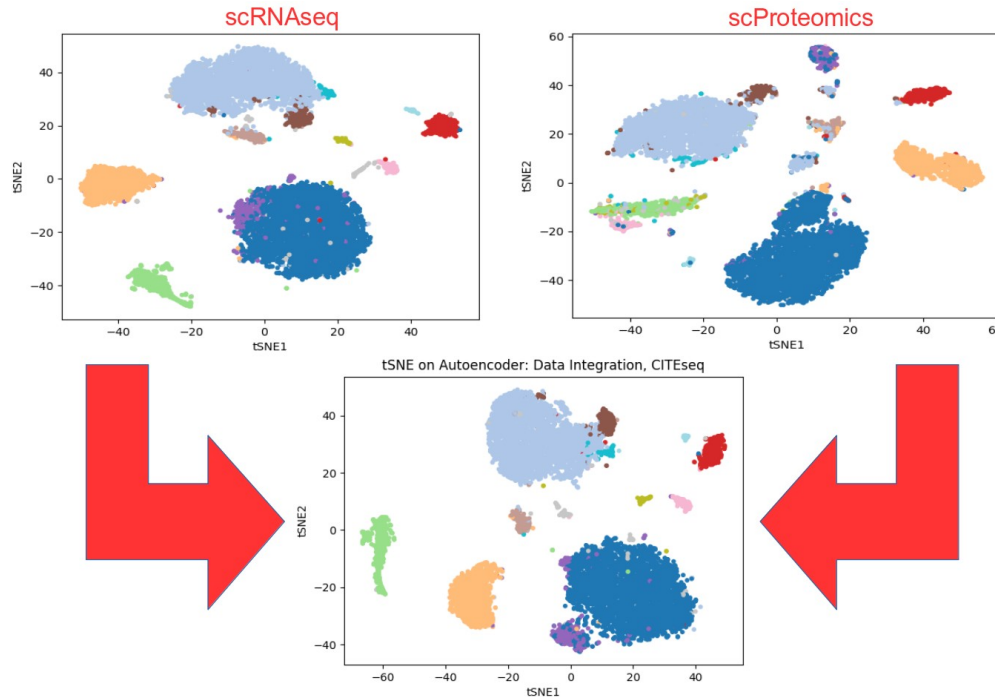


tSNE on Autoencoder: 8 Layers



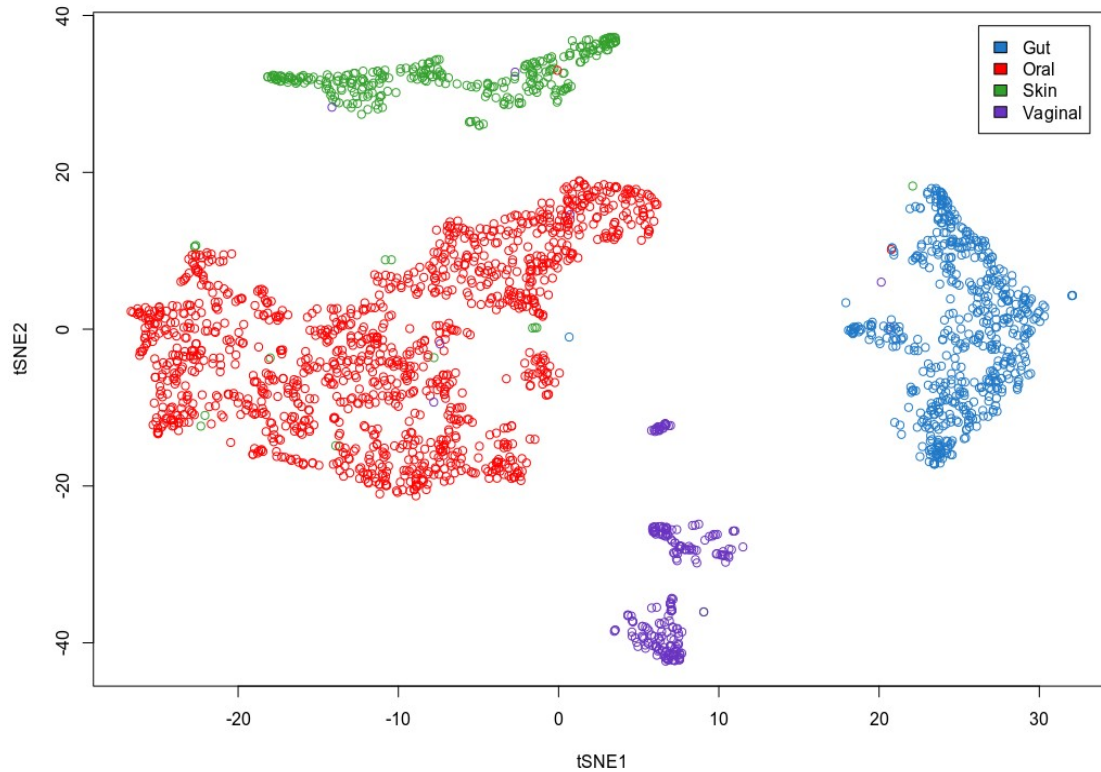


scNMTseq: Clark et al., 2018, Nature Communications 9, 781

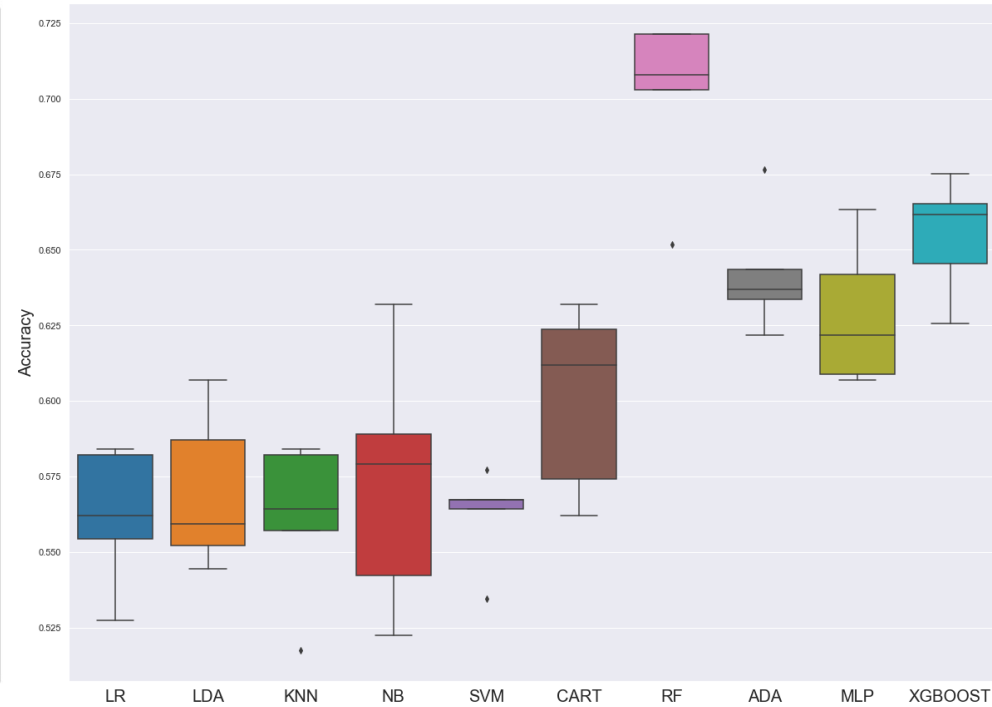


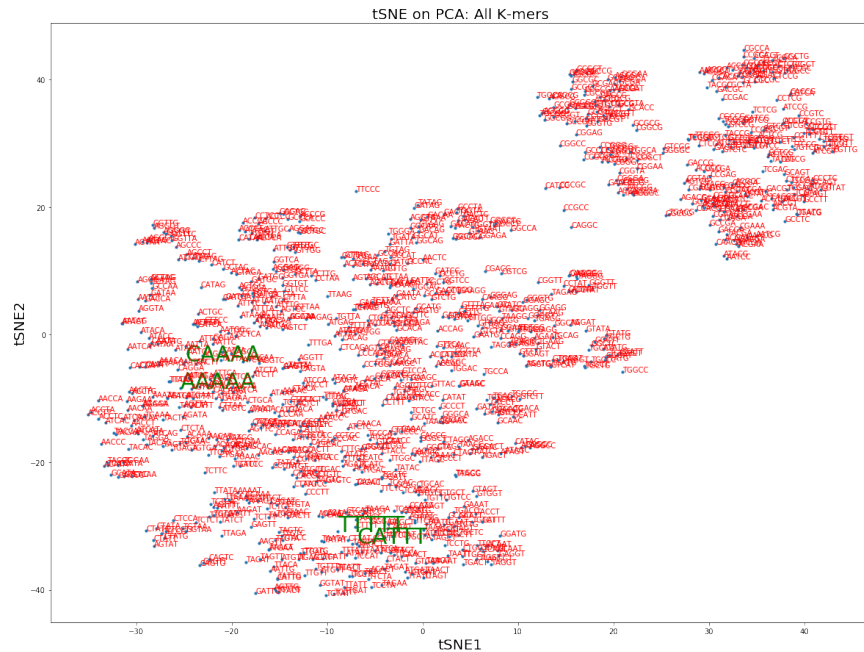


tSNE: Tissue Effect

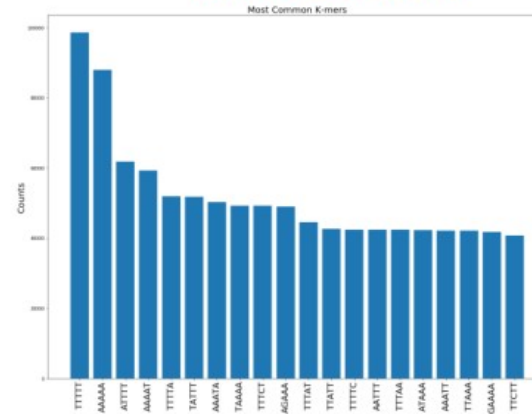


Oral microbiome: female vs. male microbial difference





Neanderthal introgressed vs depleted
NLP: Bag of Words



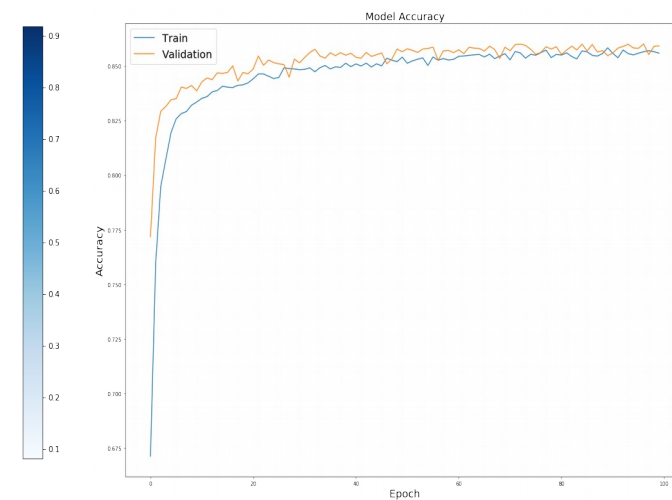
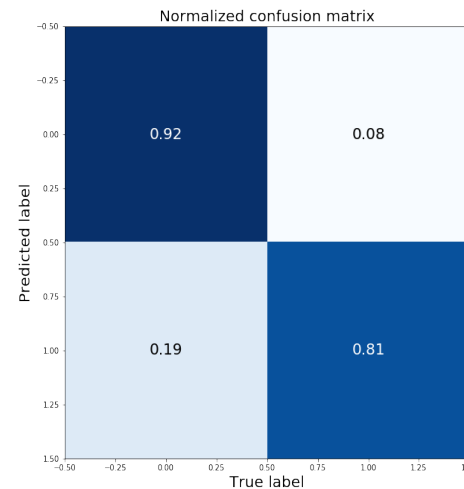
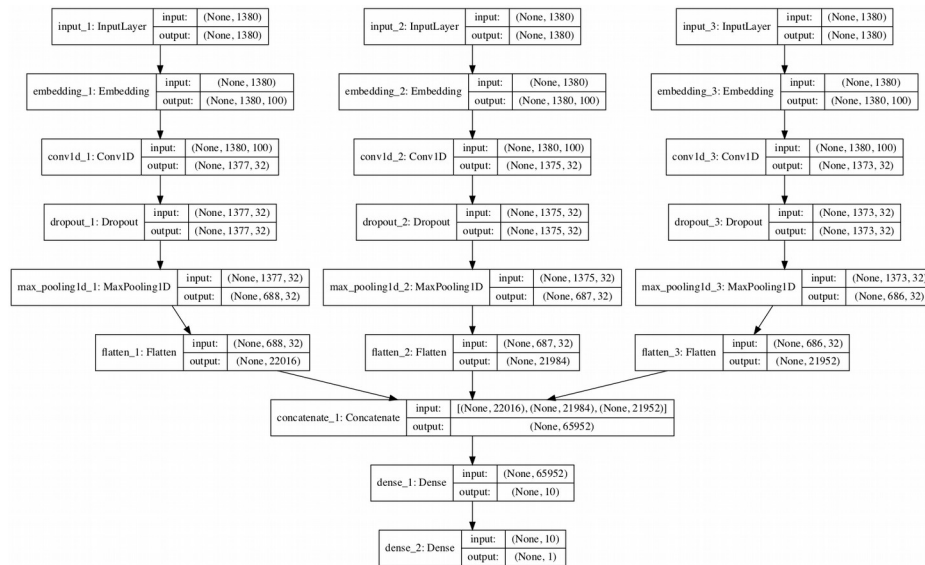
Sequence
GATCGTAC



Sentence / Text

GATCG ATCGT TCGTA CGTAC

word word word word



COVID Symptom Study
LUNDS UNIVERSITET I SAMARBETE MED KING'S COLLEGE LONDON OCH ZOE GLOBAL LTD

Om studien | Vanliga frågor | Info till studieledare | Media | Kartor | Veckorapporter | Sök på denna webbplatsen | SÖK

Kartor med uppskattad förekomst av symtomatisk covid-19 i Sverige

Aktuella kartor

Powered by Lund University and ZOE

NYHETER

2020-06-29
COVID Symptom Study live - Istället för Almedalen
When citizens engage in public health research, pitfalls and prospects

2020-06-23
COVID Symptom Study i TV4:s Nyhetsmorgon
I dag den 23 juni deltog professor Maria Gomez i Nyhetsmorgon på TV4 för att prata om COVID Symptom Study och den aktuella kartläggningen av smittspridningen av den uppskattade förekomsten av symtomatisk covid-19 [...]

2020-06-21
Nya frågor i appen: Levnadsvanor och diabetes
I den senaste uppdateringen av appen COVID Symptom Study har två nya uppdateringar av frågor som ska hjälpa oss förstå covid-19 bättre lagts till.

Lägg 1 minut om dagen på att hjälpa oss bekämpa COVID-19.

Läs mer på <https://www.covid19app.lu.se/>

App store | Google Play

Ett forskningsprojekt vid Lunds universitet i samarbete med King's College London och ZOE Global Ltd

LUNDS UNIVERSITET | ZOE | KCL

FÖLJ OSS

COVID Symptom Study

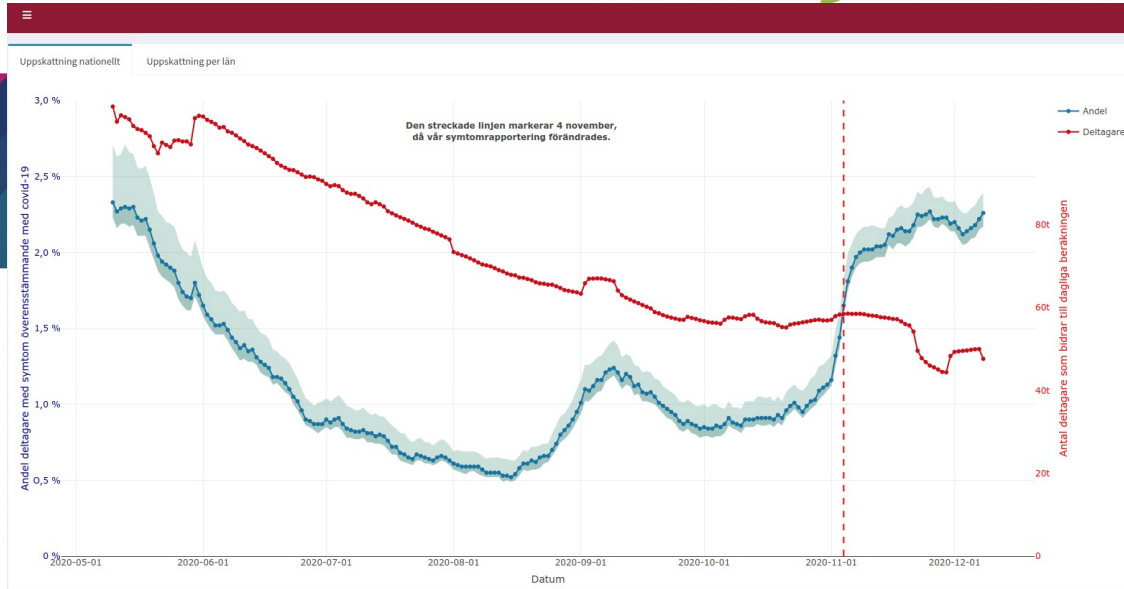
Lägg 1 minut om dagen på att hjälpa oss bekämpa COVID-19.

Läs mer på <https://www.covid19app.lu.se/>

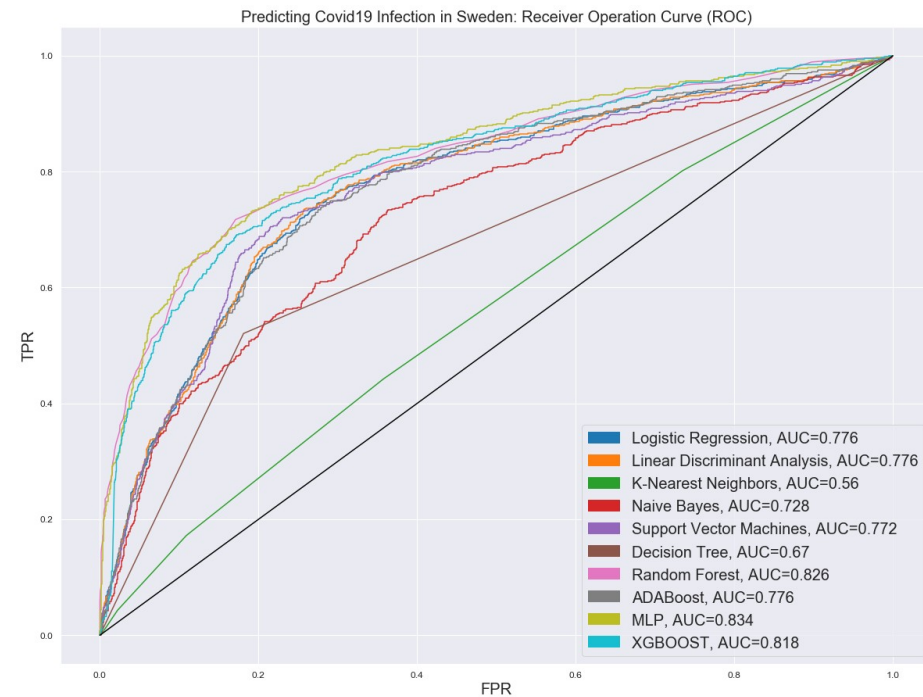
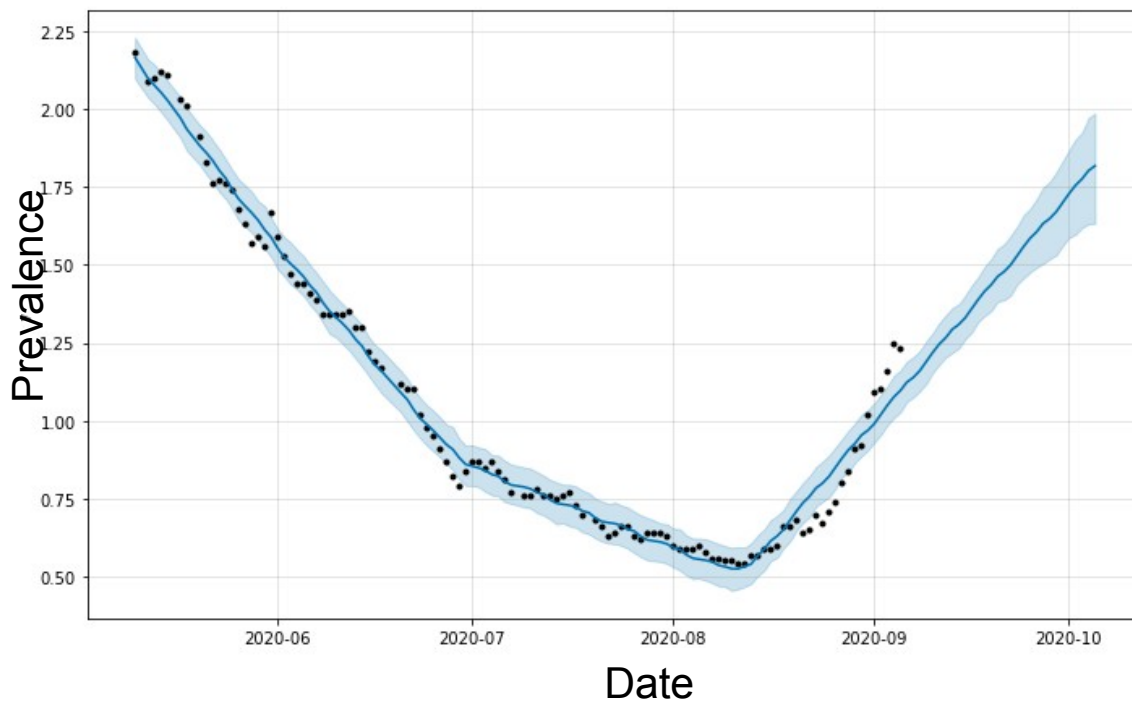
App store | Google Play

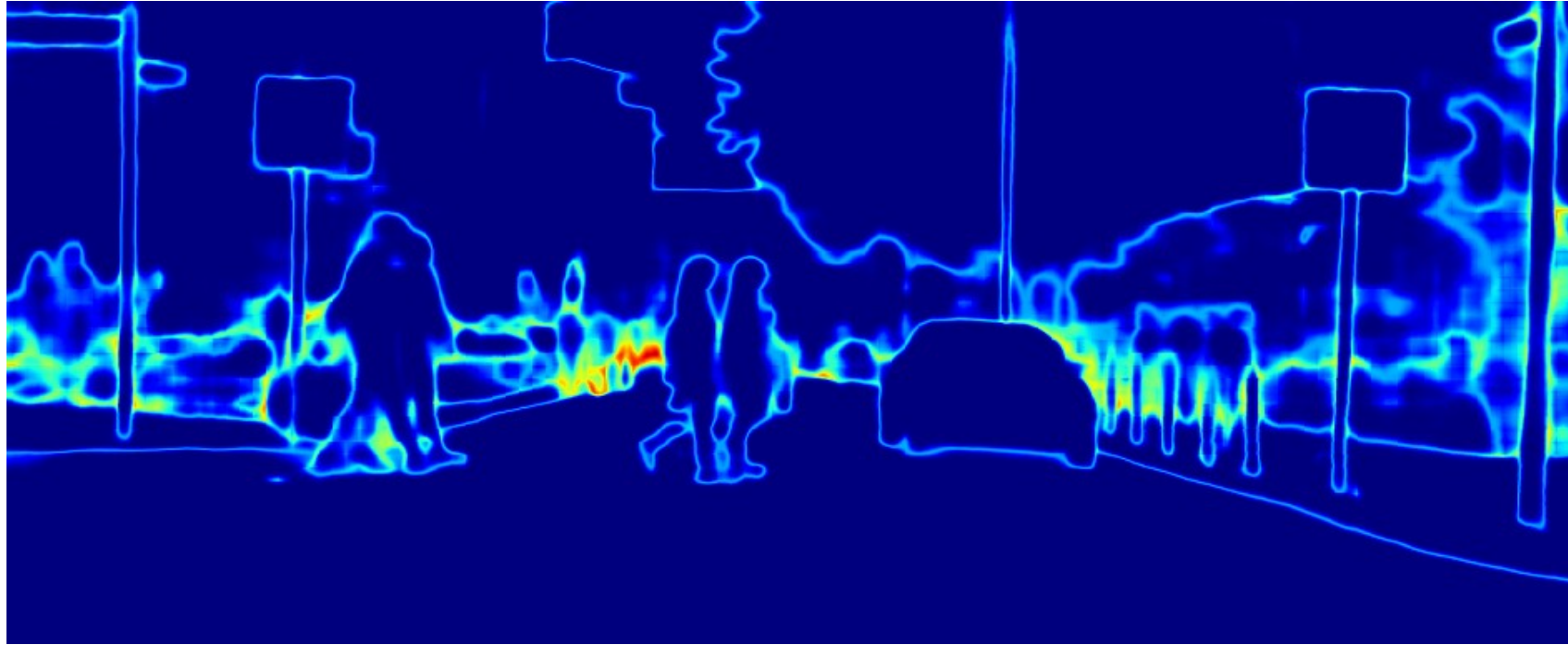
Ett forskningsprojekt vid Lunds universitet i samarbete med King's College London och ZOE Global Ltd

LUNDS UNIVERSITET | ZOE | KCL

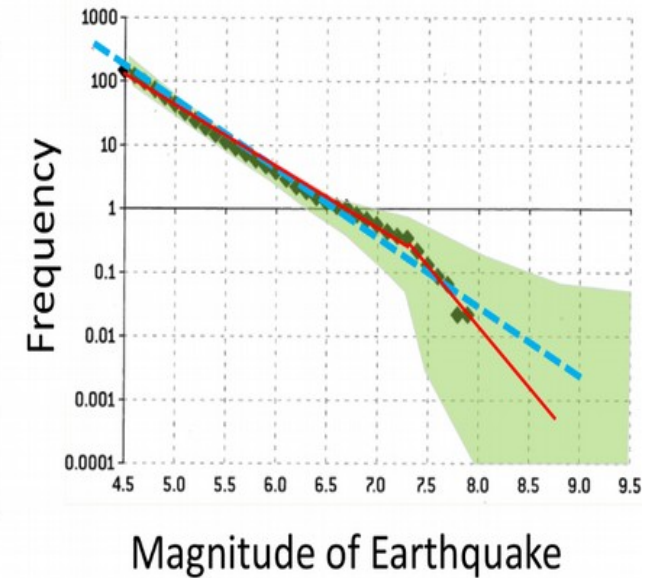
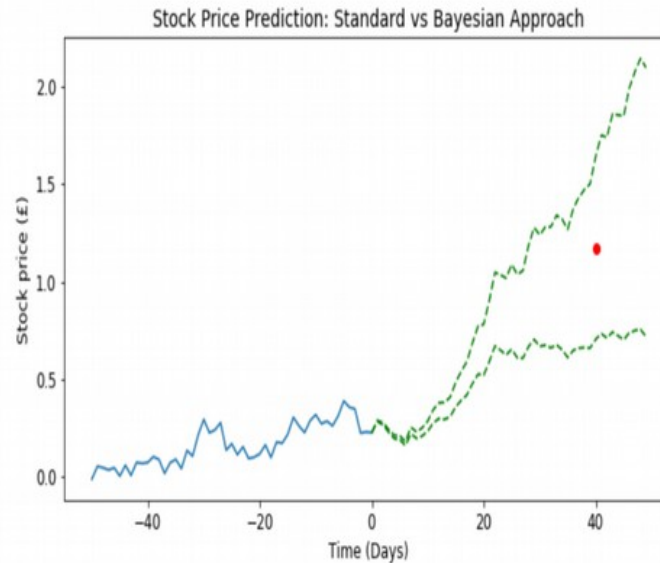


2.2 mln data points from 200 k individuals

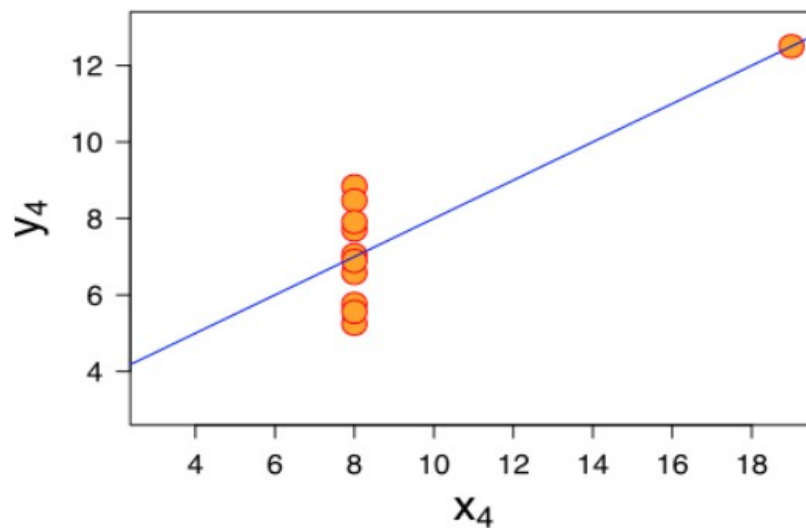
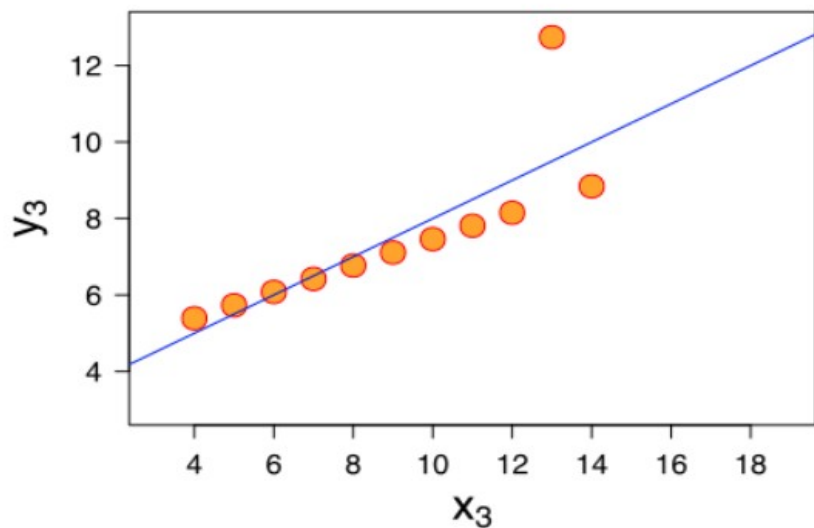
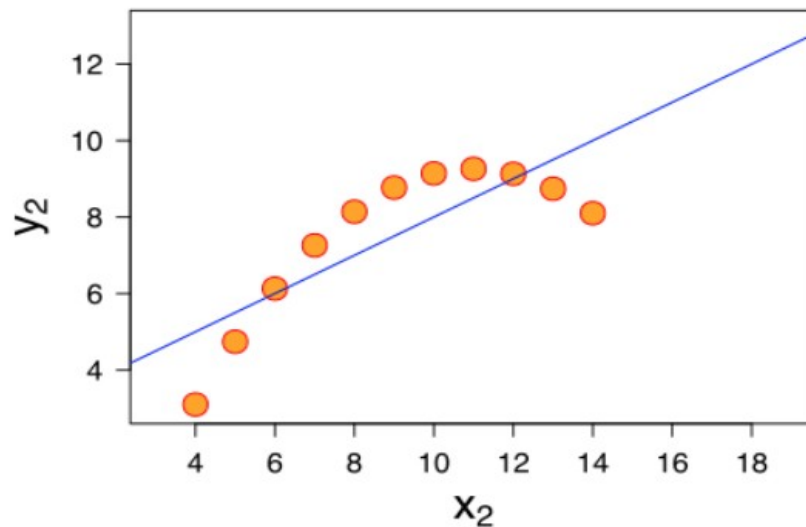
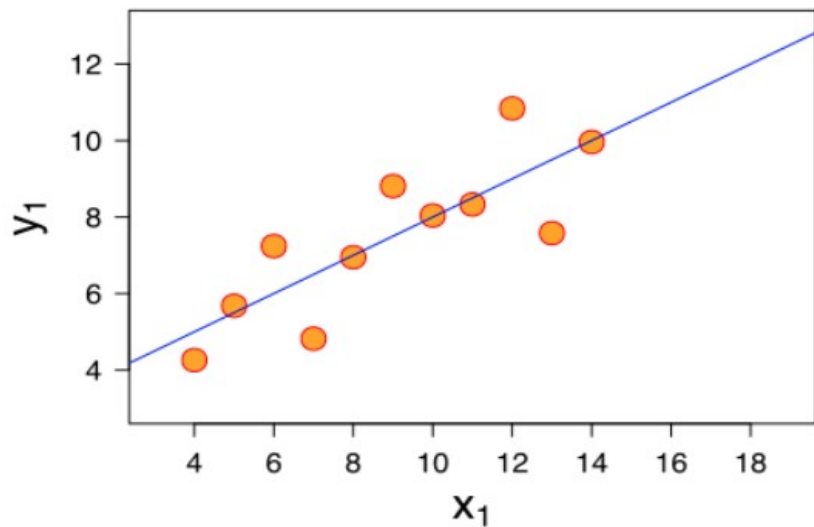




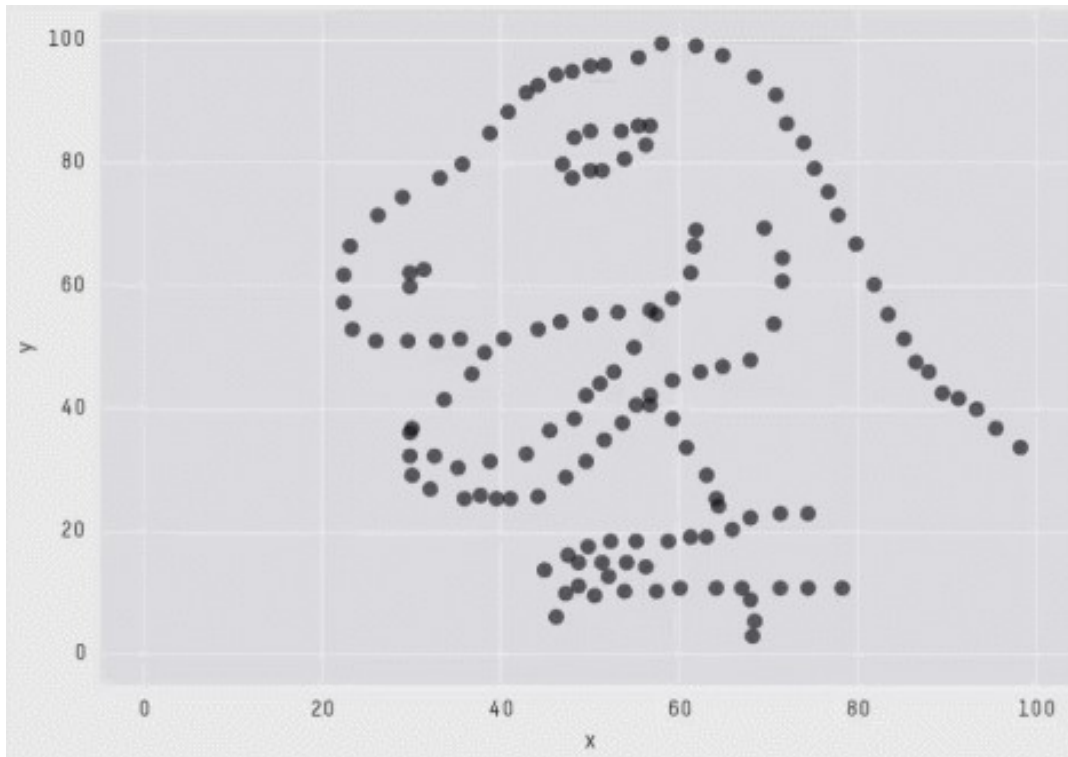
Intelligence is to know how much you do not know



Frequentist Statistics Failure

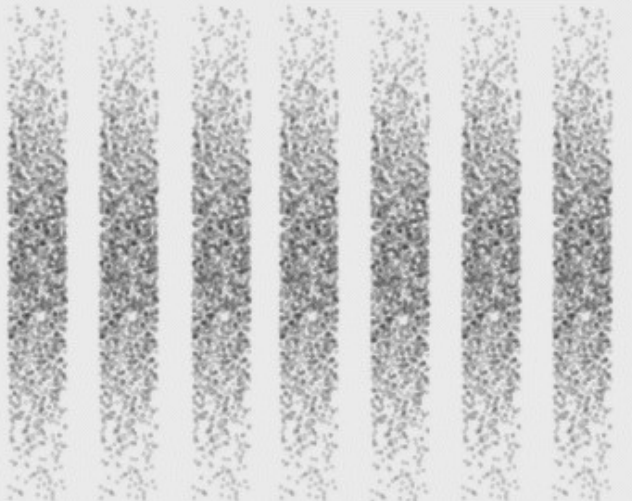


Why Frequentist Statistics is Brain Damaging

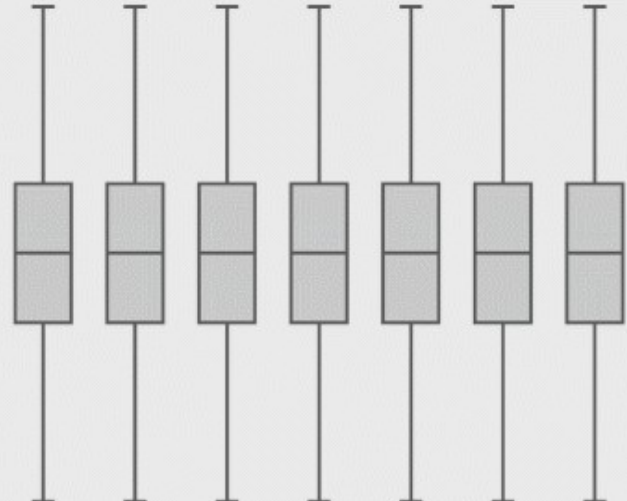


X Mean: 54.2659224
Y Mean: 47.8313999
X SD : 16.7649829
Y SD : 26.9342120
Corr. : -0.0642526

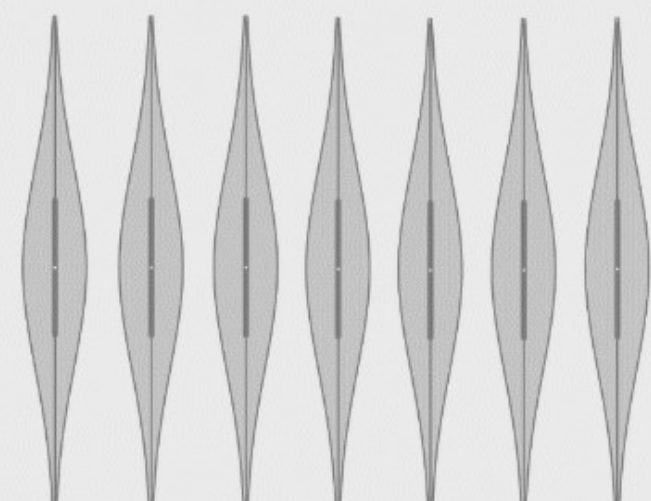
Raw Data



Box-plot of the Data



Violin-plot of the Data



A

B

C

D

E

F

G

A

B

C

D

E

F

G

A

B

C

D

E

F

G

Pvalue is not good for ranking features

nature > comment > article

nature research journal

nature

Subscribe

COMMENT - 20 MARCH 2019

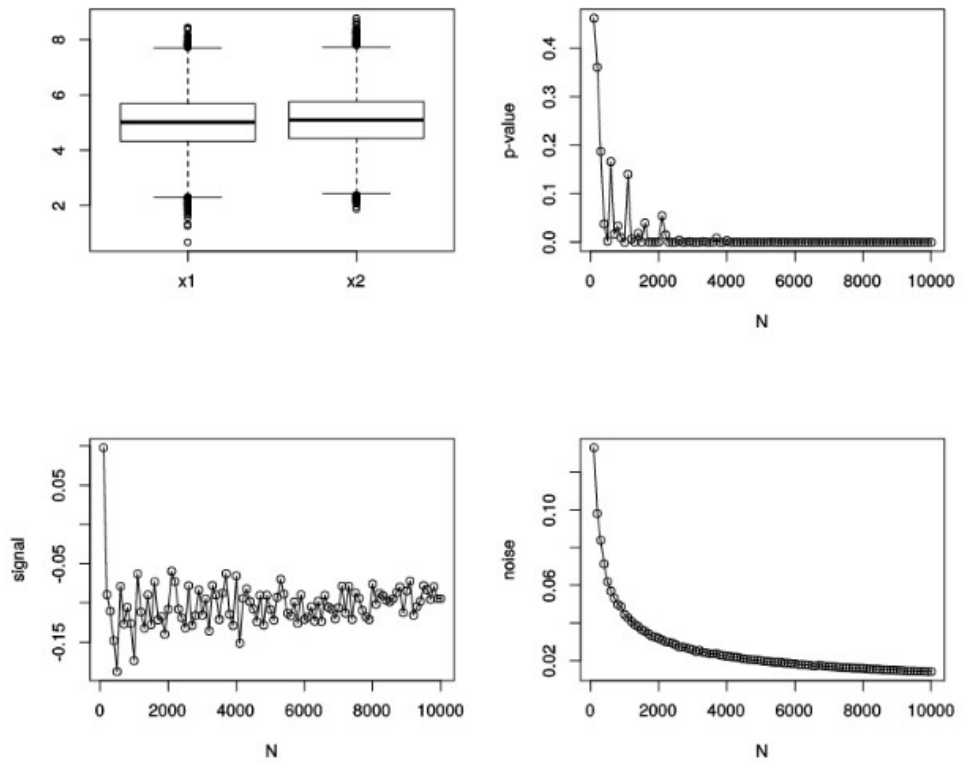
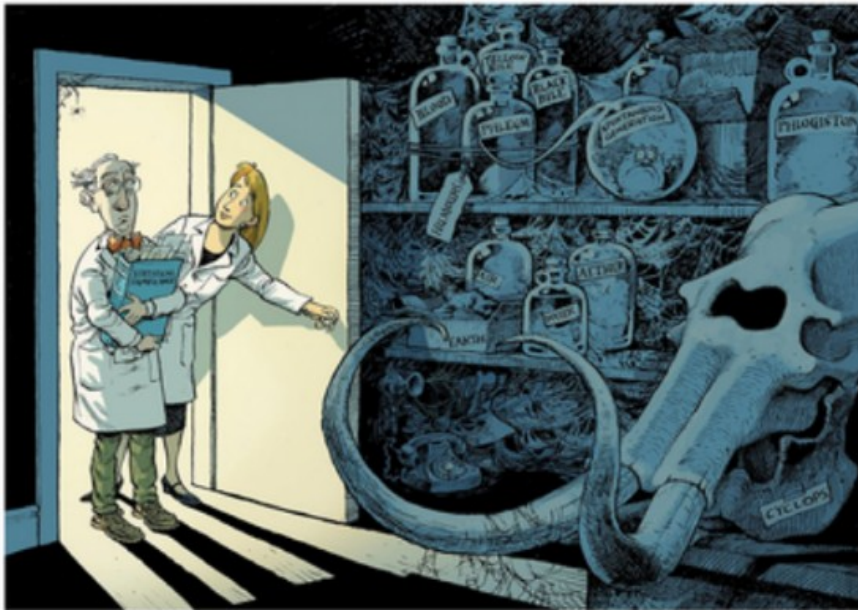
Scientists rise up against statistical significance

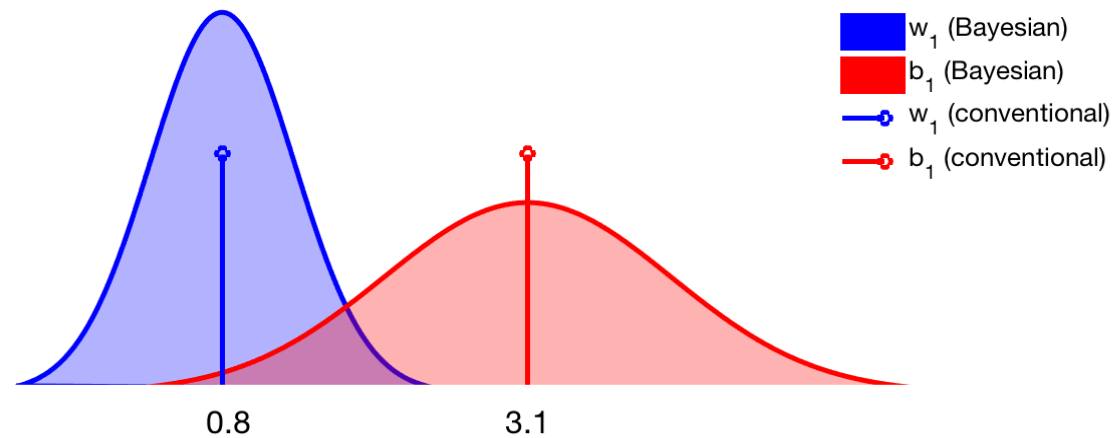
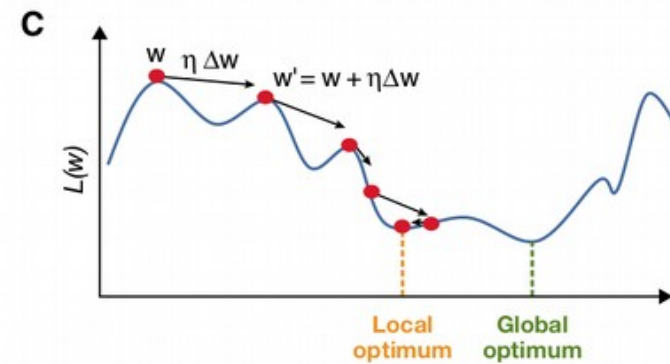
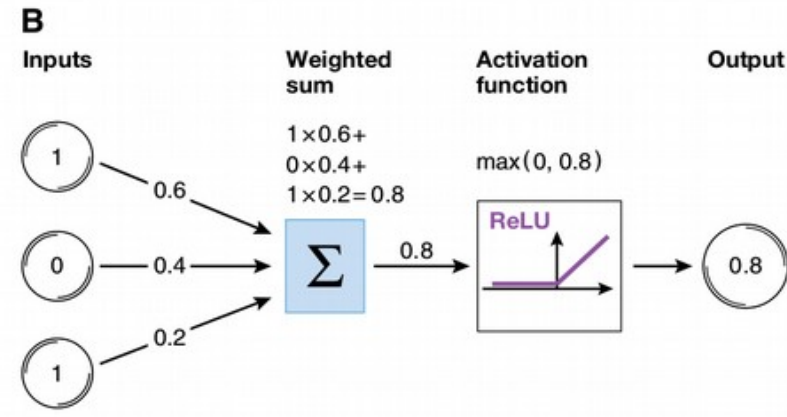
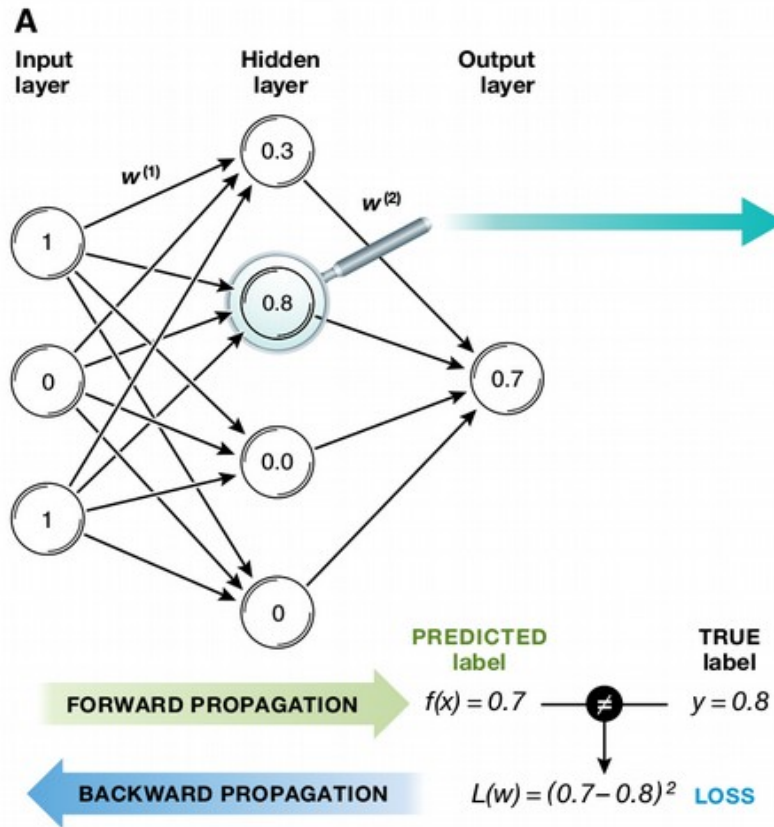
Valentin Amrhein, Sander Greenland, Blake McShane and more than 800 signatories call for an end to hyped claims and the dismissal of possibly crucial effects.

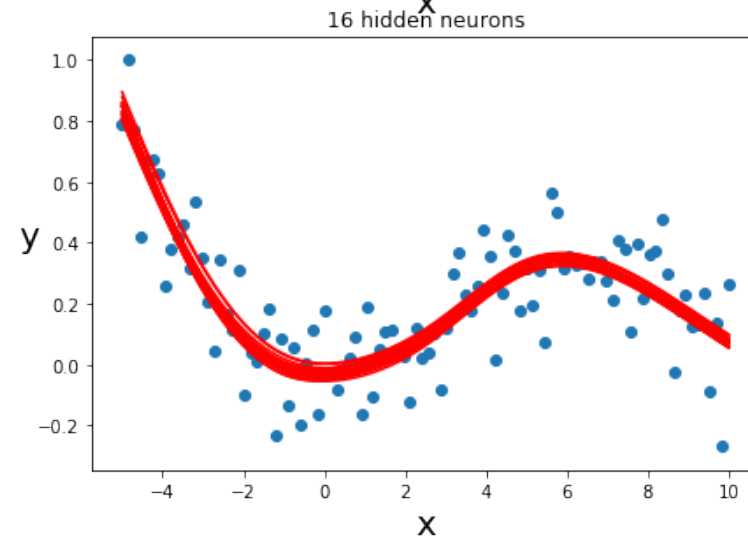
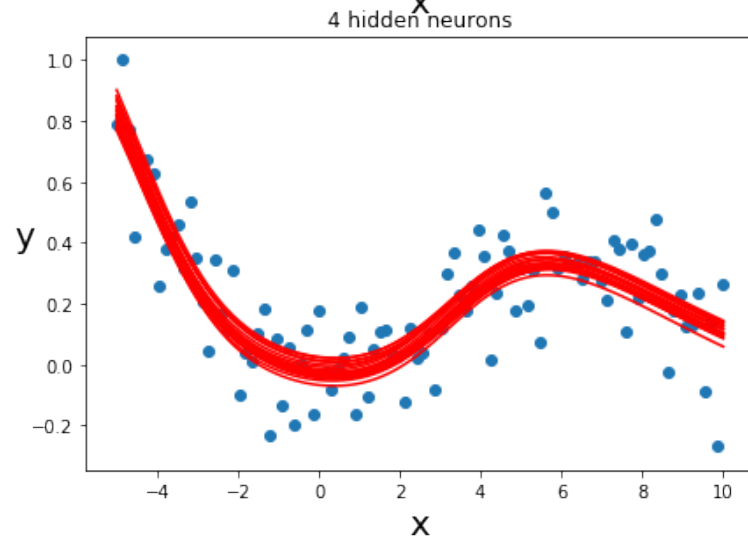
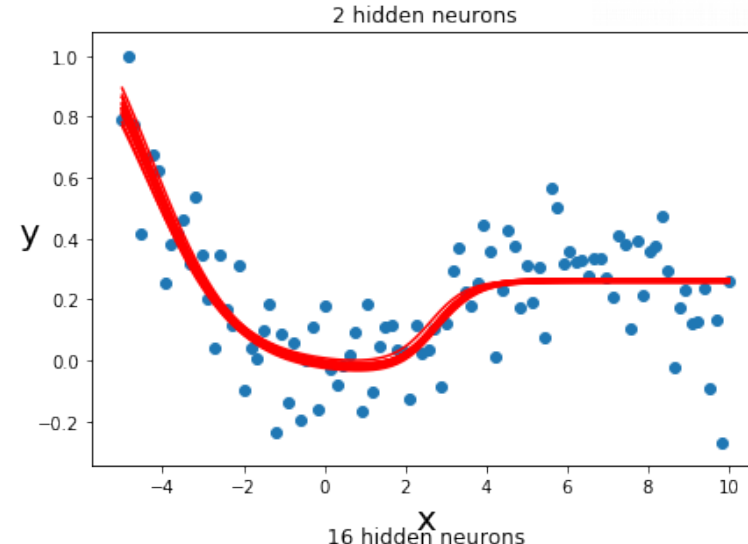
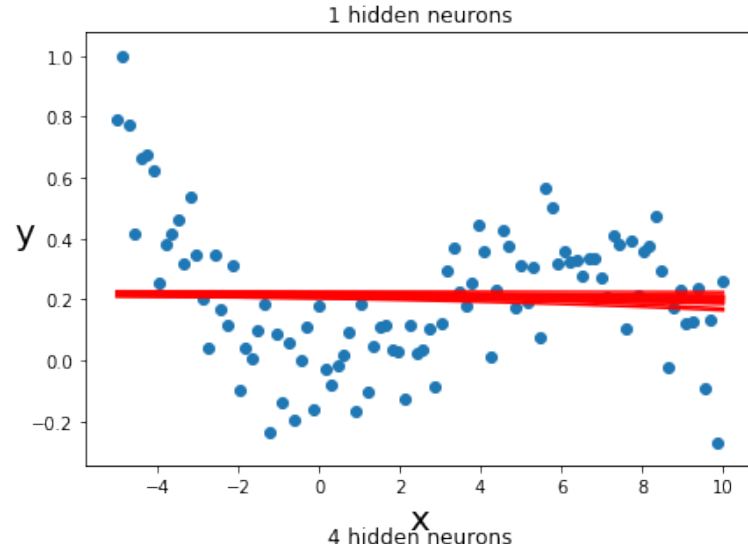
Valentin Amrhein, Sander Greenland & Blake McShane

Twitter Facebook Email

```
FC <- 1.02
x_mean <- 5; x_sd <- 1
N_vector <- seq(from=100, to=10000, by=100)
x1 <- rnorm(N_vector, x_mean, x_sd)
x2 <- rnorm(N_vector, x_mean*FC, x_sd)
```



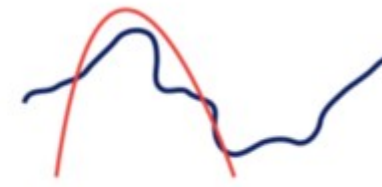




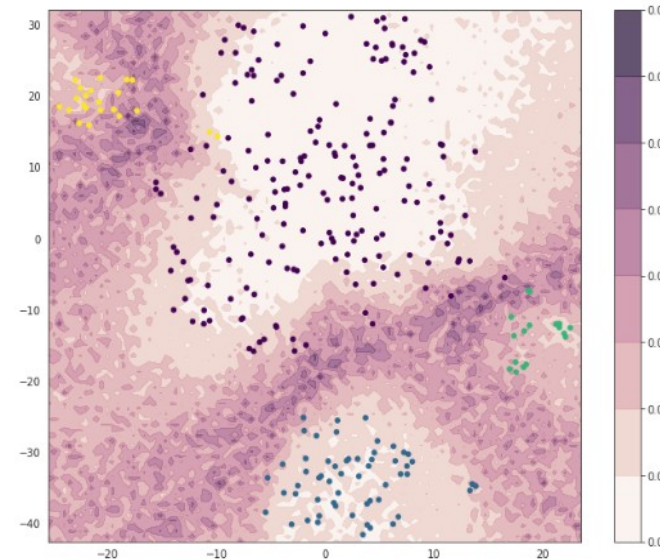
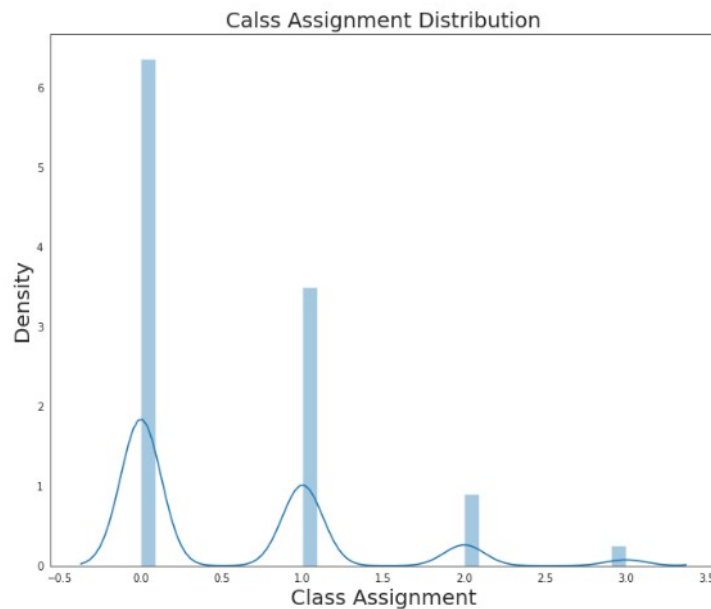
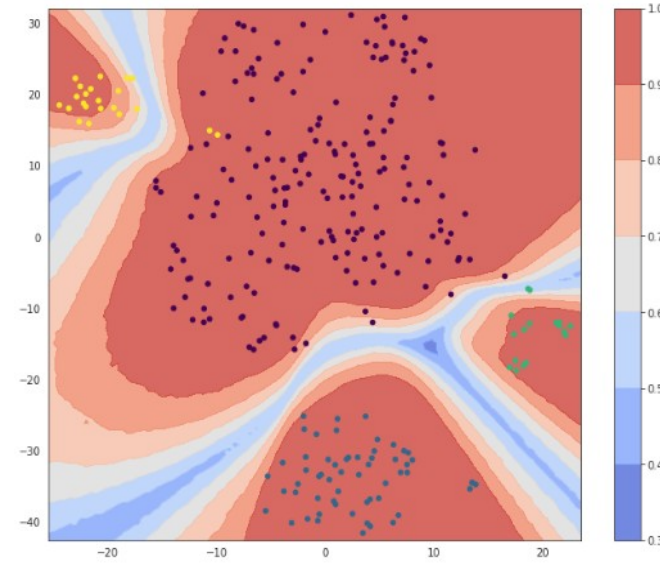
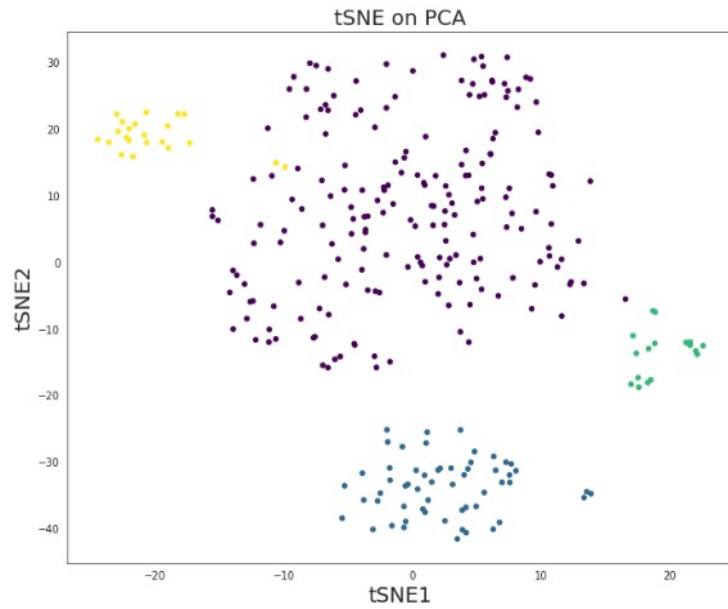
true distribution



Monte Carlo

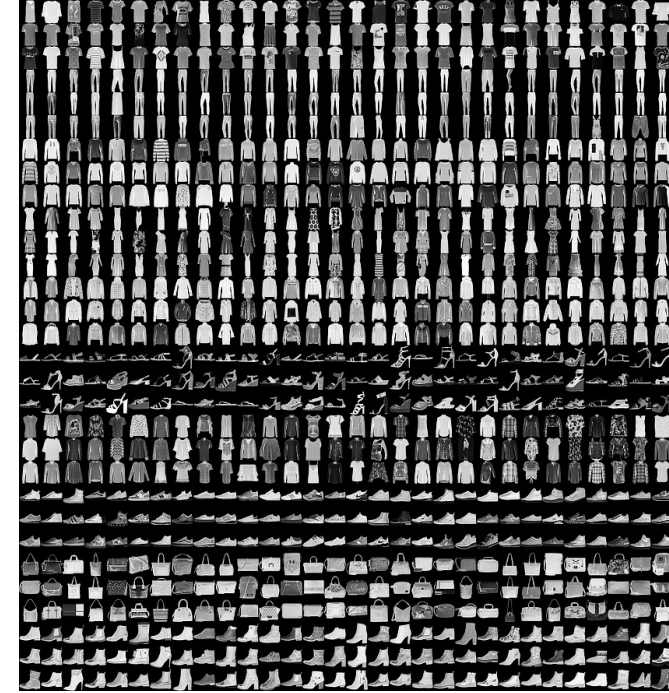


variational distribution



Bartoschek et al. 2018, Nature Communications, 9, 5150

Fashion MNIST



```
In [24]: # normalize inputs from 0-255 to 0.0-1.0
X_train = X_train.reshape(X_train.shape[0], 1, 28, 28).astype('float32')
X_test = X_test.reshape(X_test.shape[0], 1, 28, 28).astype('float32')
X_train = X_train / 255.0
X_test = X_test / 255.0

In [25]: # one hot encode outputs
y_train = np_utils.to_categorical(y_train)
y_test = np_utils.to_categorical(y_test)
num_classes = y_test.shape[1]
print(num_classes)
10

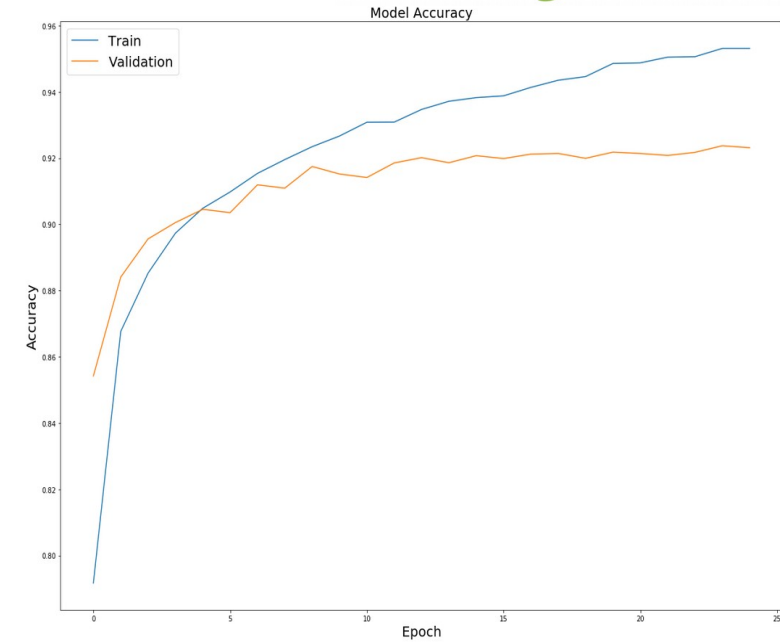
In [27]: # Create the model
model = Sequential()
model.add(Conv2D(32, (3, 3), input_shape=(1, 28, 28), padding='same', activation='relu',
                kernel_constraint=MaxNorm(3)))
model.add(Dropout(0.2))
model.add(Conv2D(32, (3, 3), padding='same', activation='relu',
                kernel_constraint=MaxNorm(3)))
model.add(MaxPooling2D(pool_size=(2, 2)))
model.add(Flatten())
model.add(Dense(512, activation='relu', kernel_constraint=MaxNorm(3)))
model.add(Dropout(0.5))
model.add(Dense(num_classes, activation='softmax'))

# Compile model
epochs = 25
rate = 0.01
decay = rate/epochs
sgd = SGD(lr=rate, momentum=0.9, decay=decay, nesterov=False)
model.compile(loss='categorical_crossentropy', optimizer=sgd, metrics=['accuracy'])
print(model.summary())

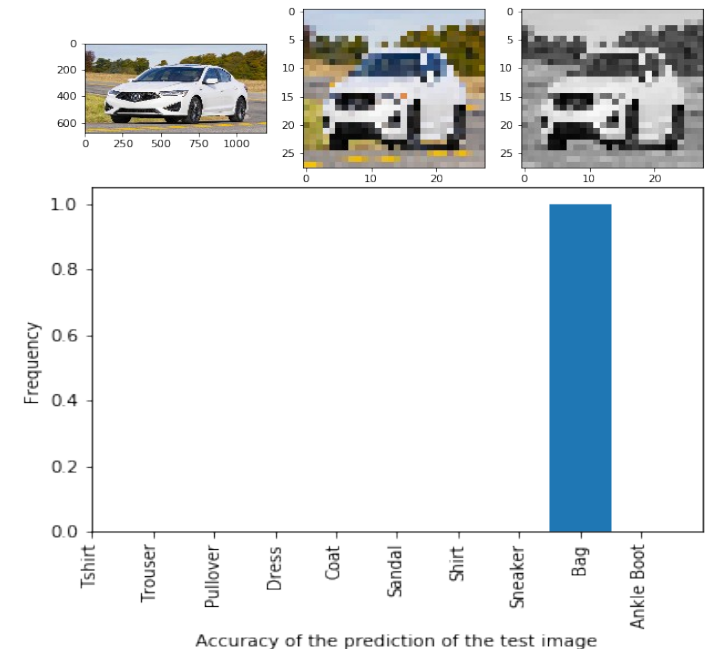
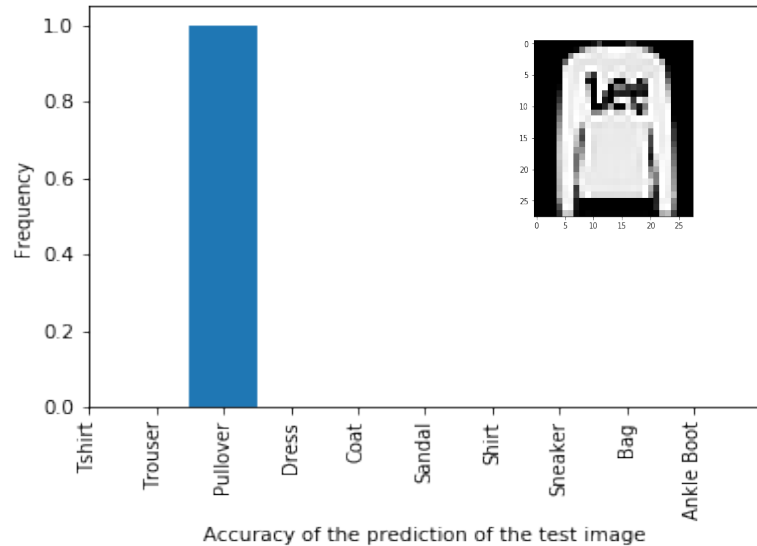
Layer (type)           Output Shape           Param #
-----
conv2d_0 (Conv2D)      (None, 32, 28, 28)     320
dropout_0 (Dropout)    (None, 32, 28, 28)     0
conv2d_1 (Conv2D)      (None, 32, 28, 28)     9248
max_pooling2d_0 (MaxPooling2D) (None, 32, 14, 14)     0
flatten_0 (Flatten)    (None, 6272)           0
dense_0 (Dense)        (None, 512)            3211776
dropout_1 (Dropout)    (None, 512)            0
dense_1 (Dense)        (None, 10)             5130
-----
Total params: 3,226,474
Trainable params: 3,226,474
Non-trainable params: 0
None

In [28]: # Fit the model
model.fit(X_train, y_train, validation_data=(X_test, y_test), epochs=epochs, batch_size=32)
history = model.fit(X_train, y_train, epochs=epochs, verbose=1, validation_split=0.25,
                    batch_size=32, shuffle=True)

Train on 45000 samples, validate on 15000 samples
Epoch 1/25
45000/45000 [=====] - 1158s 26ms/step - loss: 0.5762 - acc: 0.7917 - val_loss: 0.3973 - val_acc: 0.8542
Epoch 2/25
45000/45000 [=====] - 1124s 25ms/step - loss: 0.3643 - acc: 0.8676 - val_loss: 0.3127 - val_acc: 0.8843
Epoch 3/25
45000/45000 [=====] - 1158s 26ms/step - loss: 0.3129 - acc: 0.8653 - val_loss: 0.2825 - val_acc: 0.8956
Epoch 4/25
45000/45000 [=====] - 1609s 36ms/step - loss: 0.2813 - acc: 0.8973 - val_loss: 0.2727 - val_acc: 0.9005
Epoch 5/25
45000/45000 [=====] - 902s 20ms/step - loss: 0.2618 - acc: 0.9048 - val_loss: 0.2588 - val_acc: 0.9045
Epoch 6/25
45000/45000 [=====] - 936s 21ms/step - loss: 0.2451 - acc: 0.9090 - val_loss: 0.2564 - val_acc: 0.9035
Epoch 7/25
```



Prediction



PyMC3, Edward, TensorFlow Probability

Prediction

```
In [8]: x_train = x_train.reshape(x_train.shape[0],D)
x_test = x_test.reshape(x_test.shape[0],D)
print(x_train.shape)
print(x_test.shape)

(60000, 784)
(10000, 784)

In [9]: from keras.utils import to_categorical
y_train = to_categorical(y_train)
y_test = to_categorical(y_test)
print(y_train.shape)
print(y_test.shape)

(60000, 10)
(10000, 10)

In [10]: ed.set_seed(314159)
N = 100 # number of images in a minibatch.
D = D # number of features.
K = 10 # number of classes.

# Create a placeholder to hold the data (in minibatches) in a TensorFlow graph.
x = tf.placeholder(tf.float32, [None, D])
# Normal(0,1) priors for the variables. Note that the syntax assumes TensorFlow 1.1.
w = Normal(loc=tf.zeros([D, K]), scale=tf.ones([D, K]))
b = Normal(loc=tf.zeros(K), scale=tf.ones(K))
# Categorical likelihood for classification.
y = Categorical(tf.matmul(x, w) + b)

In [11]: # Construct the q(w) and q(b). in this case we assume Normal distributions.
qw = Normal(loc=tf.Variable(tf.random_normal([D, K])),
            scale=tf.nn.softplus(tf.Variable(tf.random_normal([D, K]))))
qb = Normal(loc=tf.Variable(tf.random_normal([K])),
            scale=tf.nn.softplus(tf.Variable(tf.random_normal([K]))))

In [12]: def generator(arrays, batch_size = N):
starts = [0] * len(arrays) # pointers to where we are in iteration
while True:
    batches = []
    for i, array in enumerate(arrays):
        start = starts[i]
        stop = start + batch_size
        diff = stop - array.shape[0]
        if diff <= 0:
            batch = array[start:stop]
            starts[i] += batch_size
        else:
            batch = np.concatenate((array[start:], array[:diff]))
            starts[i] = diff
        batches.append(batch)
    yield batches
cifar10 = generator([x_train, y_train], N)

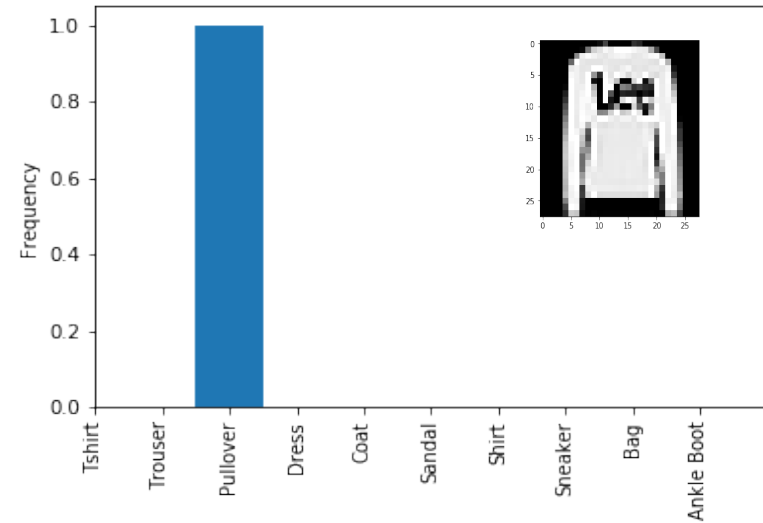
In [13]: # We use a placeholder for the labels in anticipation of the training data.
y_ph = tf.placeholder(tf.int32, [N])
# Define the VI inference technique, ie. minimise the KL divergence between q and p.
inference = ed.KLqp({w: qw, b: qb}, data={y: y_ph})
# Initialise the inference variables
inference.initialize(n_iter=50000, n_print=100, scale={y: float(x_train.shape[0]) / N})
# We will use an interactive session.
sess = tf.InteractiveSession()
# Initialise all the variables in the session.
tf.global_variables_initializer().run()
# Let the training begin. We load the data in minibatches and update the VI inference using each new batch.
for _ in range(inference.n_iter):
    X_batch, Y_batch = next(cifar10)
    x_batch = X_batch.reshape(N, -1)
    # TensorFlow method gives the label data in a one hot vector format. We convert that into a single label.
    Y_batch = np.argmax(Y_batch, axis=1)
    info_dict = inference.update(feed_dict={x: X_batch, y_ph: Y_batch})
    inference.print_progress(info_dict)

50000/50000 [100%] Elapsed: 221s | Loss: 85453.266

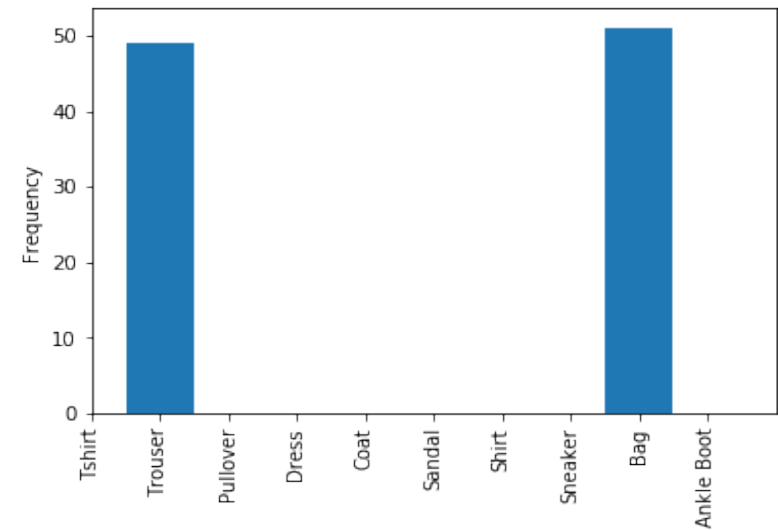
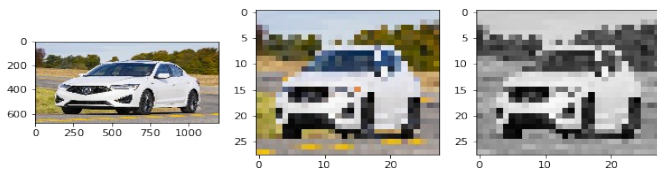
In [14]: # Generate samples the posterior and store them.
n_samples = 100
prob_lst = []
samples = []
w_samples = []
b_samples = []
for _ in range(n_samples):
    w_samp = qw.sample()
    b_samp = qb.sample()
    w_samples.append(w_samp)
    b_samples.append(b_samp)
    # Also compute the probability of each class for each (w,b) sample.
    prob = tf.nn.softmax(tf.matmul(x_test, w_samp) + b_samp)
    prob_lst.append(prob.eval())
    sample = tf.concat([tf.reshape(w_samp,[-1]),b_samp,0])
    samples.append(sample.eval())

In [15]: # Compute the accuracy of the model.
# For each sample we compute the predicted class and compare with the test labels.
# Predicted class is defined as the one which as maximum probability.
# We perform this test for each (w,b) in the posterior giving us a set of accuracies
# Finally we make a histogram of accuracies for the test data.
accy_test = []
for prob in prob_lst:
    y_trn_prd = np.argmax(prob, axis=1).astype(np.float32)
    acc = (y_trn_prd == np.argmax(y_test, axis=1)).mean()*100
    accy_test.append(acc)

plt.hist(accy_test)
plt.title("Histogram of prediction accuracies in the CIFAR10 test data")
plt.xlabel("Accuracy")
plt.ylabel("Frequency")
plt.show()
```



Accuracy of the prediction of the test image



Accuracy of the prediction of the test image



*Knut och Alice
Wallenbergs
Stiftelse*



LUNDS
UNIVERSITET