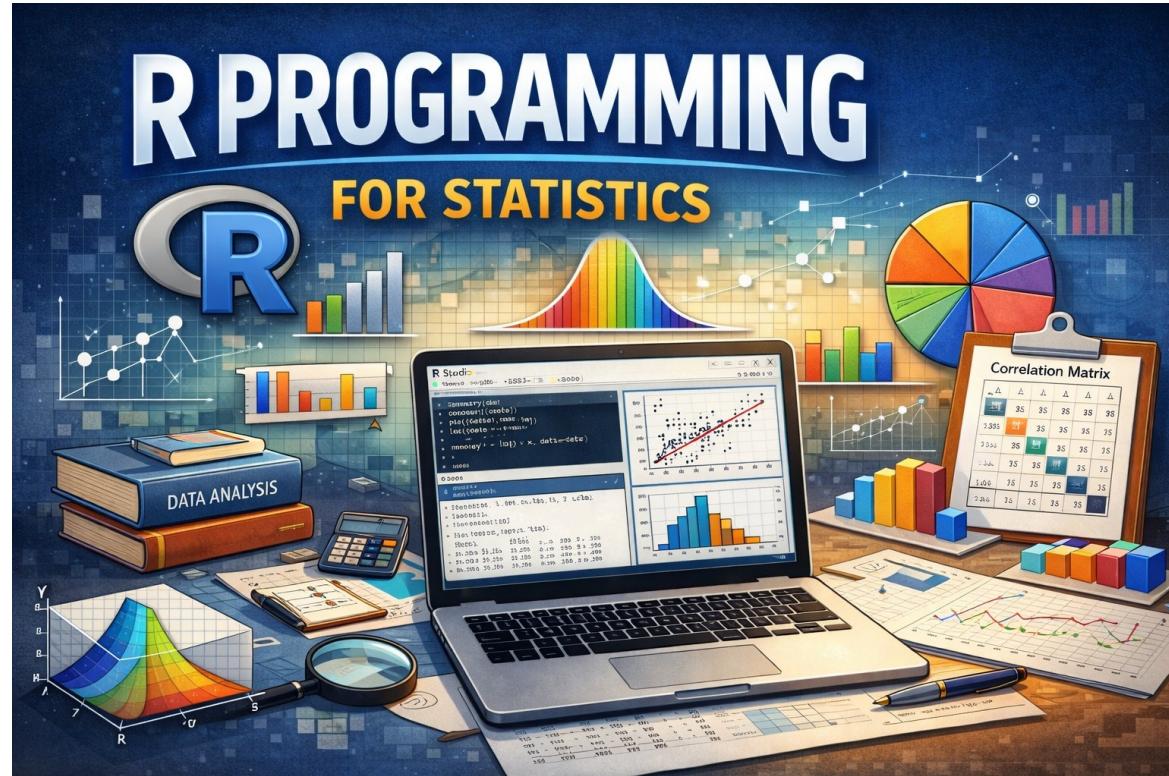




R Programming for Statistics and Machine Learning

Nikolay Oskolkov, MRG Group Leader, LIOS, Riga, Latvia
R course, 05.02.2026



@NikolayOskolkov



@oskolkov.bsky.social

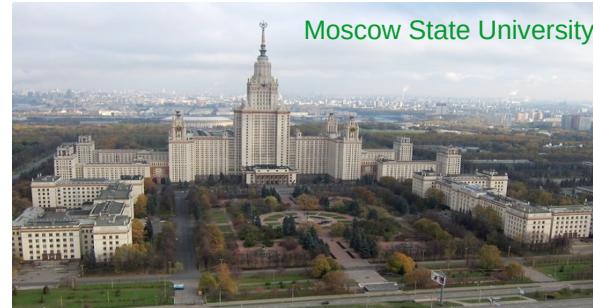
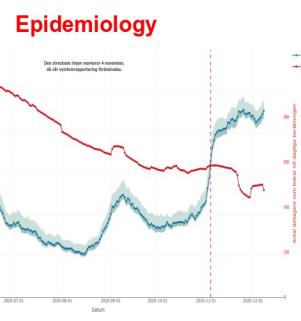
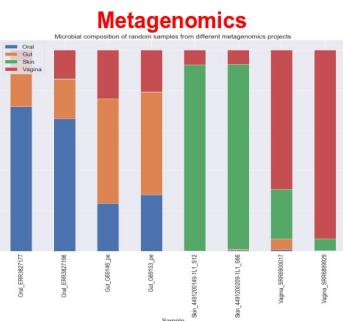
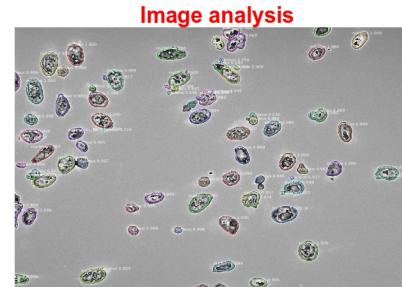
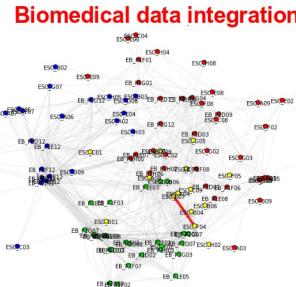
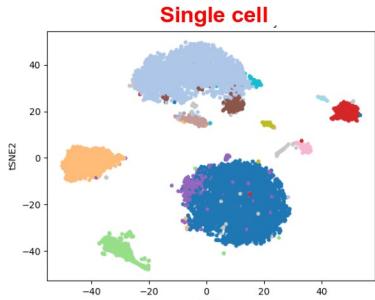


Personal homepage:
<https://nikolay-oskolkov.com>

Brief introduction: who am I

- 2007 PhD in theoretical physics in Moscow, Russia
- 2011 medical genomics at Lund University, Sweden
- 2016 bioinformatician at NBIS SciLifeLab, Sweden
- 2025 Metabolic Research Group leader, LIOS, Latvia

50+ publications; h-index = 25; 4,000+ citations





Metabolic Research Group at LIOS



Publications

Conferences

Metabolic Research Group

Metabolic Research Group

The Metabolic Research Group (MRG) focuses on advancing computational methods to identify and validate novel drug targets for metabolic diseases. Our research profile centers on the development and application of machine learning approaches, combined with statistical modeling, to extract biological knowledge from complex datasets. A key expertise of the group is the integration of diverse multiOmic data—including genomics, transcriptomics, proteomics, metabolomics, and metagenomics—enabling a systems-level understanding of metabolic processes and disease mechanisms. Through this integrative and data-driven approach, we aim to contribute to precision medicine by supporting the discovery of innovative therapeutic strategies within the TARGETWISE project.

1 more postdoctoral fellow and 1 PhD student to be hired

**If you know anyone who might be interested,
please contact me!**



PhD. Nikolay Oskolkov
Group Leader (PI) of the Metabolic
Research Group



Kristina Grausa, PhD student
in Metabolic Research Group



Daniel Rivas, MD, PhD in AI,
postdoctoral fellow
in Metabolic Research Group



About you

- Name
- University / Institute / Company
- Research interest(s)
- Previous experience with computational analysis and bioinformatics
- Motivation to join the course
- Expectations from the course



Housekeeping rules



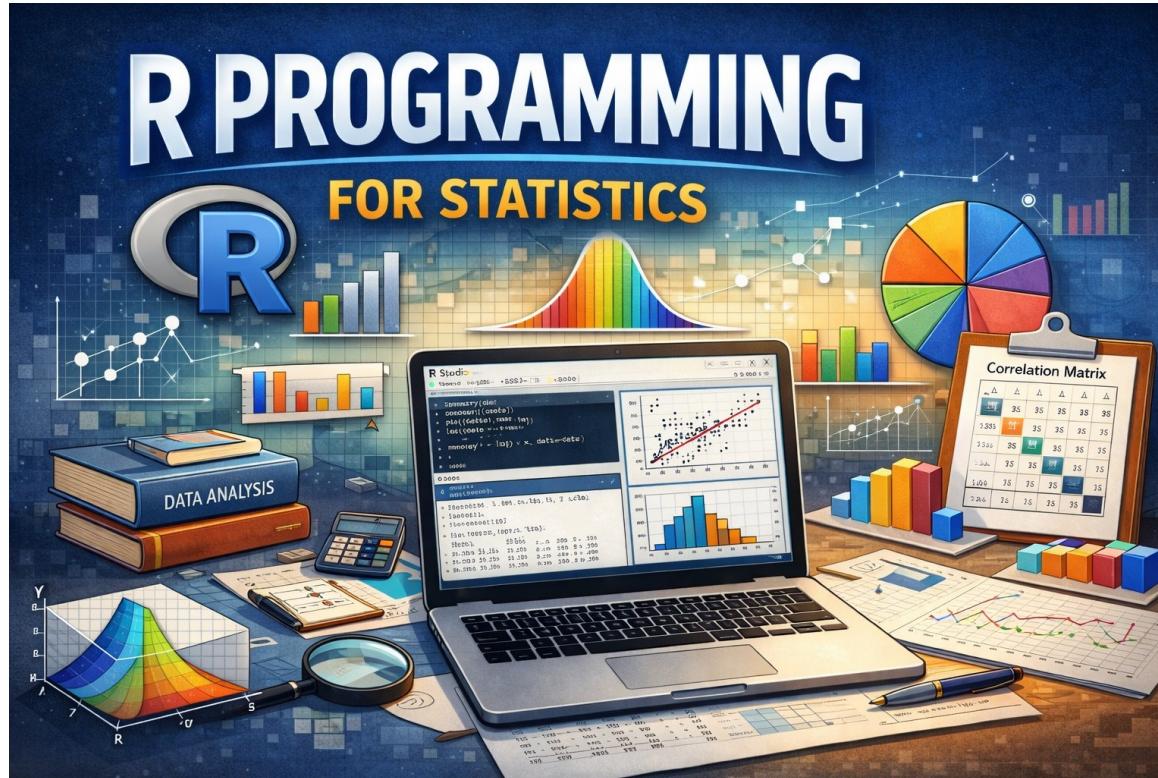
Latvian Institute of
Organic Synthesis

- Please **ask your questions anytime**, it is ok to interrupt me, or you can raise your hand and ask. You can also send me your questions by email.
- The course includes **7 lectures (~1h each)** and **8 practicals (~1h each)**, there are **30 min breaks** after each lecture and each practical.
- During practicals, we will use **Rmarkdown notebooks**. At the end of each practical, I will use rendered html-versions of the notebooks to go through command lines with my explanations.
- The material is based on data and problems from **computational biology**, however the concepts discussed are general and can be applied for other types of data.



Introduction to Statistical Analysis in R

Nikolay Oskolkov, MRG Group Leader, LIOS, Riga, Latvia
R course, 05.02.2026



@oskolkov.bsky.social



Personal homepage:
<https://nikolay-oskolkov.com>



Session content

Topics we'll cover in this session:

- 1) High-dimensional biological data and curse of dimensionality
- 2) Data-driven choice of statistical analysis
- 3) Basics of Frequentist statistics: pros and cons of p-value
- 4) Basics of Bayesian statistics: likelihood, prior and posterior

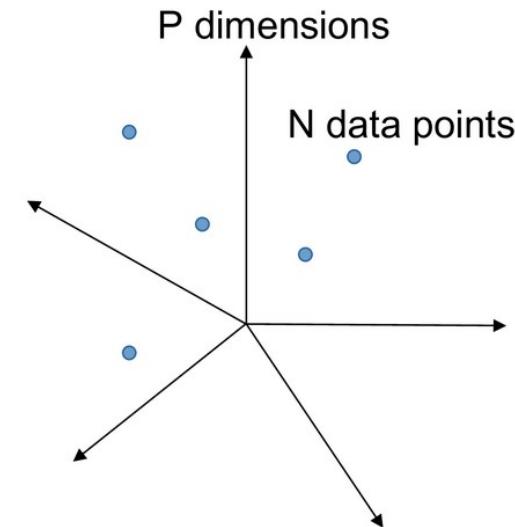
Biological data are high dimensional

Statistical observations:
e.g. samples, cells etc.

Features: genes, proteins,
microbes, metabolites etc.

N

0	3	1	0	2	3	8	1	1	3
1	1	0	0	7	1	2	2	3	3
1	2	2	0	0	6	7	1	2	2
1	2	3	10	0	4	6	1	0	5
3	2	2	1	4	3	2	1	6	0
7	4	4	5	3	9	6	1	6	1
7	1	1	5	2	8	9	1	3	6
5	0	1	6	2	0	0	0	1	5
1	6	3	3	4	6	2	0	1	1
1	2	2	4	1	1	3	0	8	2



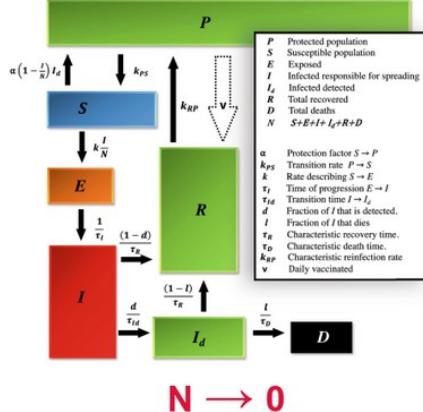
High Dimensional Data:
 $P \gg N$

For a robust statistical analysis, one should properly “sample” the P-dimensional space, hence large sample size is required, $N \gg P$

Types of data analysis

Biology / Biomedicine

Mathematical modeling

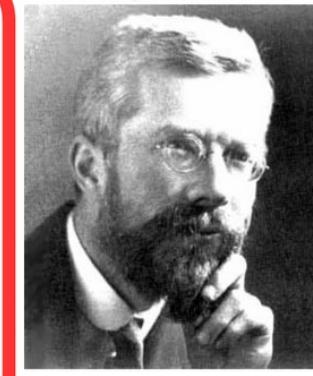


Bayesianism



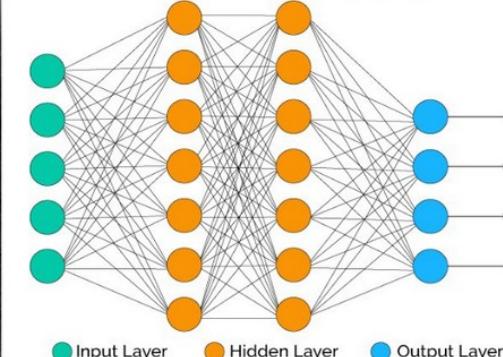
$N \ll P$

Frequentism



$N \approx P$

Machine Learning



$N > P$

Hypothesis-driven

Data-driven

Amount of Data

The Curse of Dimensionality

Ex.1

$$Y = \alpha + \beta X$$

$$\beta = (X^T X)^{-1} X^T Y$$

$$(X^T X)^{-1} \sim \frac{1}{\det(X^T X)} \cdots \rightarrow \infty, \quad n \ll p$$

$$\text{Ex.2} \quad E[\hat{\sigma}^2] = \frac{n-p}{n} \sigma^2$$

Biased Maximum Likelihood variance estimator at $n \ll p$

Some peculiarities of Frequentist statistics

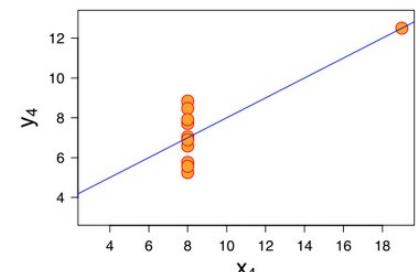
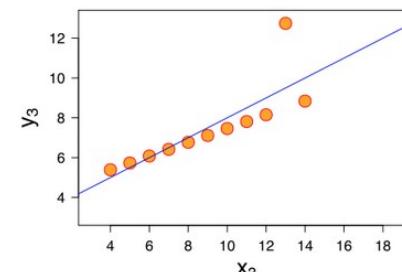
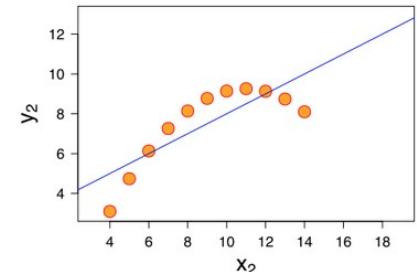
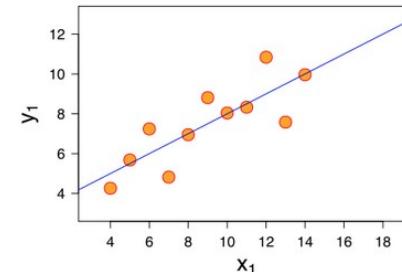
- based on Maximum Likelihood principle
- focus too much on summary statistics

$$L(x_i | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp^{-\frac{\sum_{i=1}^N (x_i - \mu)^2}{2\sigma^2}}$$

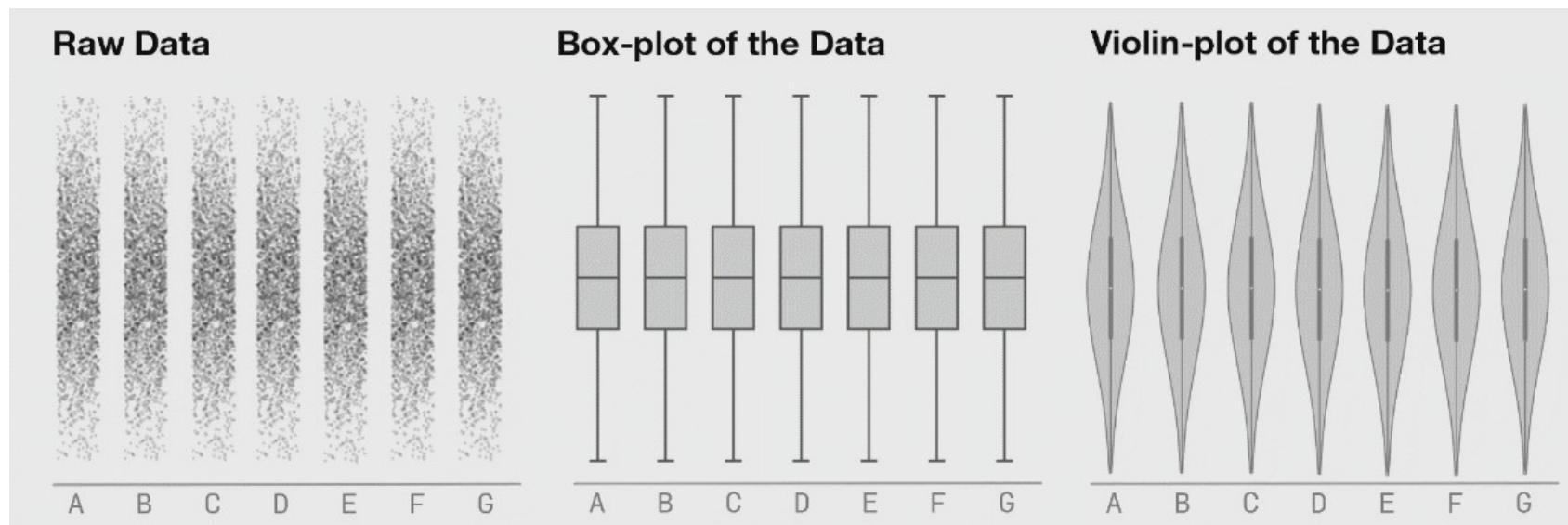
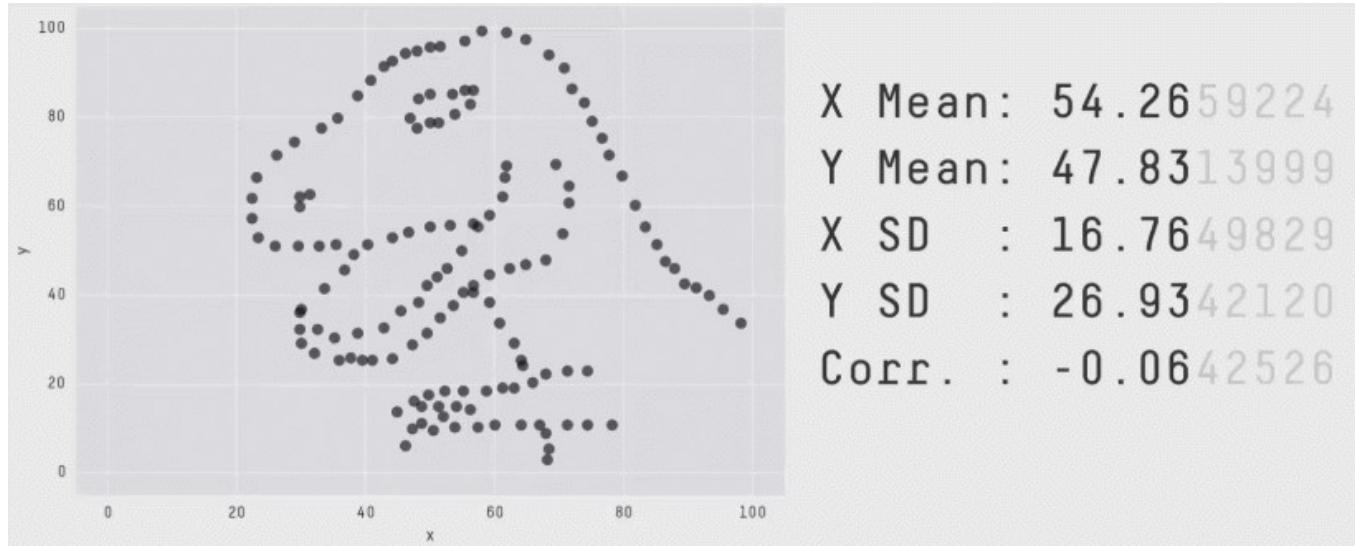
$$\frac{\partial L(x_i | \mu, \sigma^2)}{\partial \mu} = 0; \quad \frac{\partial L(x_i | \mu, \sigma^2)}{\partial \sigma^2} = 0$$

$$\mu = \frac{1}{N} \sum_{i=0}^N x_i - \text{mean estimator}$$

$$\sigma^2 = \frac{1}{N} \sum_{i=0}^N (x_i - \mu)^2 - \text{variance estimator}$$



Summary statistics do not always reasonably describe data (example: Anscombes quartet)

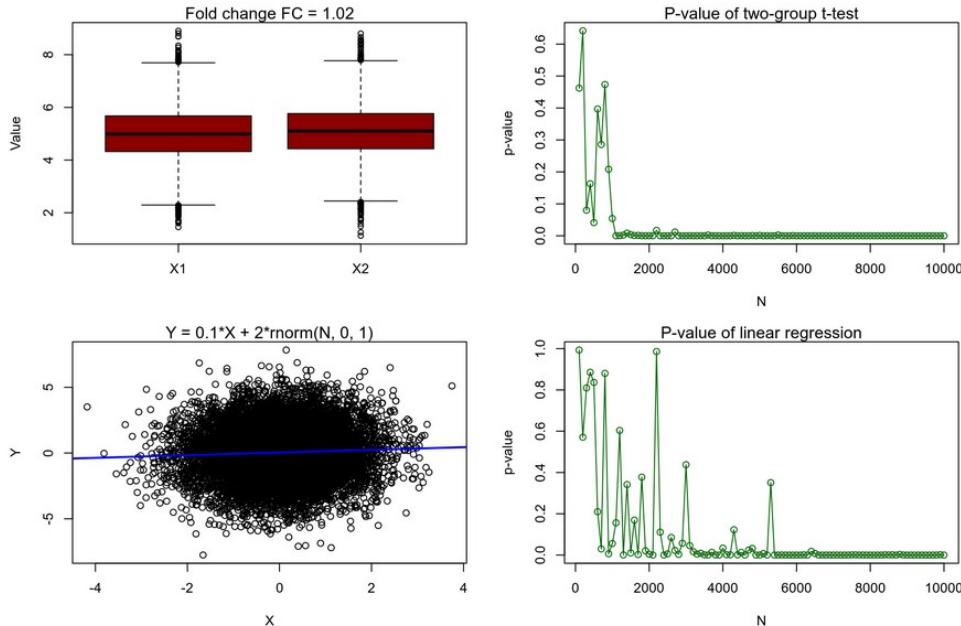


Frequentist statistics: focus too much on p-values

```

1 FC<-1.02; x_mean<-5; x_sd<-1; N_vector<-seq(from=100,to=10000,by=100); pvalue_t<-vector(); pvalue_lm<-vector()
2 for(N in N_vector)
3 {
4   x1 <- rnorm(N, x_mean, x_sd); x2 <- rnorm(N, x_mean+FC, x_sd)
5   t_test_res<-t.test(x1, x2); pvalue_t <- append(pvalue_t, t_test_res$p.value)
6
7   x <- rnorm(N, 0, 1); y <- 0.1*x+2*rnorm(N, 0, 1)
8   lm_res <- summary(lm(y~x)); pvalue_lm <- append(pvalue_lm, lm_res$coefficients[2,4])
9 }
10 par(mfrow=c(2,2)); par(mar = c(5, 5, 1, 1))
11 boxplot(x1, x2, names=c("X1", "X2"), ylab="Value", col="darkred"); mtext("Fold change FC = 1.02")
12 plot(pvalue_t~N_vector,type='o',xlab="N",ylab="p-value",col="darkgreen"); mtext("P-value of two-group t-test")
13 plot(y~x, xlab="X", ylab="Y"); abline(lm(y~x), col="blue", lwd=2); mtext("Y = 0.1*X + 2*rnorm(N, 0, 1)")
14 plot(pvalue_lm~N_vector,type='o',xlab="N",ylab="p-value",col="darkgreen"); mtext("P-value of linear regression")

```

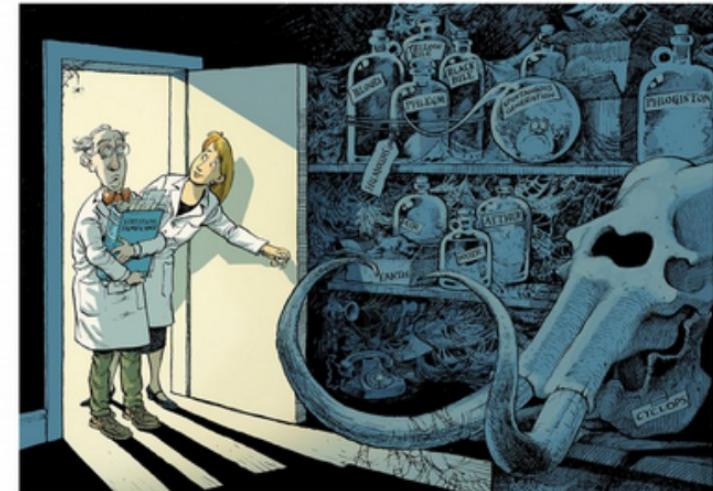


COMMENT · 20 MARCH 2019

Scientists rise up against statistical significance

Vulnent in 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



Questionable whether p-value is a best metric for ranking features (biomarkers)

Frequentist statistics struggles with high-dimensional data

```

1 n <- 20 # number of samples
2 p <- 2 # number of features / dimensions
3 Y <- rnorm(n)
4 X <- matrix(rnorm(n * p), n, p)
5 summary(lm(Y ~ X))

```

Call:
`lm(formula = Y ~ X)`

Residuals:

	Min	1Q	Median	3Q	Max
	-2.0522	-0.6380	0.1451	0.3911	1.8829

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.14950	0.22949	0.651	0.523
X1	-0.09405	0.28245	-0.333	0.743
X2	-0.11919	0.24486	-0.487	0.633

Residual standard error: 1.017 on 17 degrees of freedom
 Multiple R-squared: 0.02204, Adjusted R-squared: -0.09301
 F-statistic: 0.1916 on 2 and 17 DF, p-value: 0.8274

Going to higher dimensions →

```

1 n <- 20 # number of samples
2 p <- 10 # number of features / dimensions
3 Y <- rnorm(n)
4 X <- matrix(rnorm(n * p), n, p)
5 summary(lm(Y ~ X))

```

Call:
`lm(formula = Y ~ X)`

Residuals:

	Min	1Q	Median	3Q	Max
	-1.0255	-0.4320	0.1056	0.4493	1.0617

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.54916	0.26472	2.075	0.0679 .
X1	0.38013	0.21698	1.384	0.1998
X2	0.68053	0.27693	2.457	0.0363 *
X3	-0.10675	0.26018	-0.410	0.6911
X4	-0.21367	0.33690	-0.634	0.5417
X5	-0.19123	0.31881	-0.600	0.5634
X6	0.81074	0.25221	3.214	0.0106 *
X7	0.09634	0.24143	0.399	0.6992
X8	-0.29864	0.19084	-1.571	0.1505
X9	-0.78175	0.35408	-2.208	0.0546 .
X10	0.83736	0.36936	2.267	0.0496 *

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.8692 on 9 degrees of freedom
 Multiple R-squared: 0.6592, Adjusted R-squared: 0.2805
 F-statistic: 1.741 on 10 and 9 DF, p-value: 0.2089

Going to even higher dimensions →

```

1 n <- 20 # number of samples
2 p <- 20 # number of features / dimensions
3 Y <- rnorm(n)
4 X <- matrix(rnorm(n * p), n, p)
5 summary(lm(Y ~ X))

```

Call:
`lm(formula = Y ~ X)`

Residuals:
 ALL 20 residuals are 0: no residual degrees of freedom!

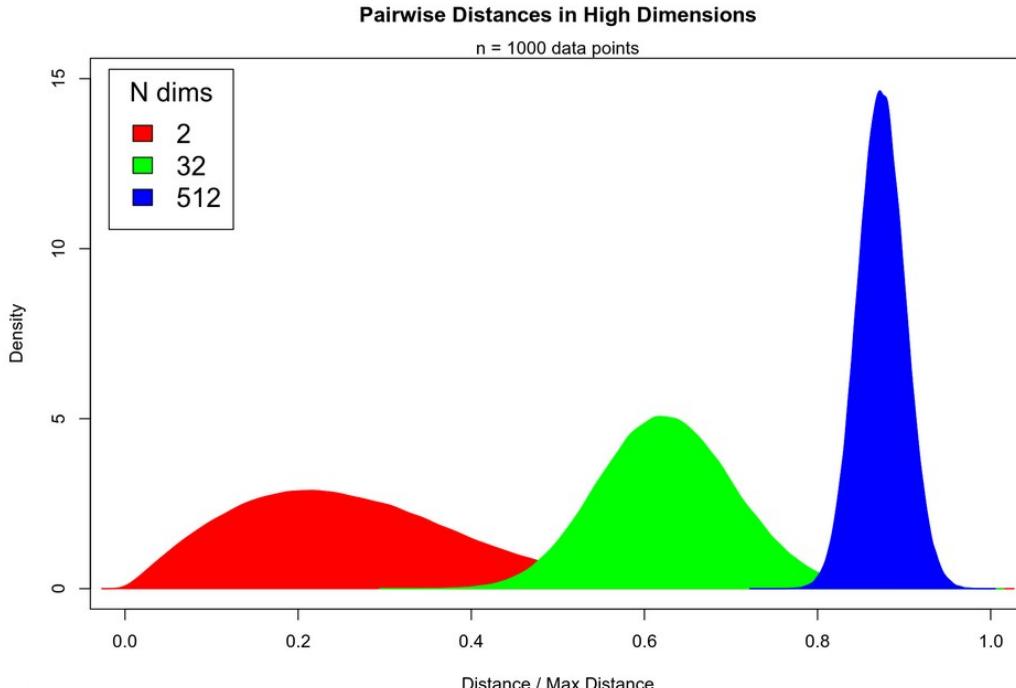
Coefficients: (1 not defined because of singularities)

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	1.34889	NaN	NaN	NaN
X1	0.66218	NaN	NaN	NaN
X2	0.76212	NaN	NaN	NaN
X3	-1.35033	NaN	NaN	NaN
X4	-0.57487	NaN	NaN	NaN
X5	0.02142	NaN	NaN	NaN
X6	0.40290	NaN	NaN	NaN
X7	0.03313	NaN	NaN	NaN
X8	-0.31983	NaN	NaN	NaN
X9	-0.92833	NaN	NaN	NaN
X10	0.18091	NaN	NaN	NaN
X11	-1.37618	NaN	NaN	NaN
X12	2.11438	NaN	NaN	NaN
X13	-1.75103	NaN	NaN	NaN
X14	-1.55073	NaN	NaN	NaN
X15	0.01112	NaN	NaN	NaN
X16	-0.50943	NaN	NaN	NaN
X17	-0.47576	NaN	NaN	NaN
X18	0.31793	NaN	NaN	NaN
X19	1.43615	NaN	NaN	NaN
X20	NA	NA	NA	NA

Residual standard error: NaN on 0 degrees of freedom
 Multiple R-squared: 1, Adjusted R-squared: NaN
 F-statistic: NaN on 19 and 0 DF, p-value: NA

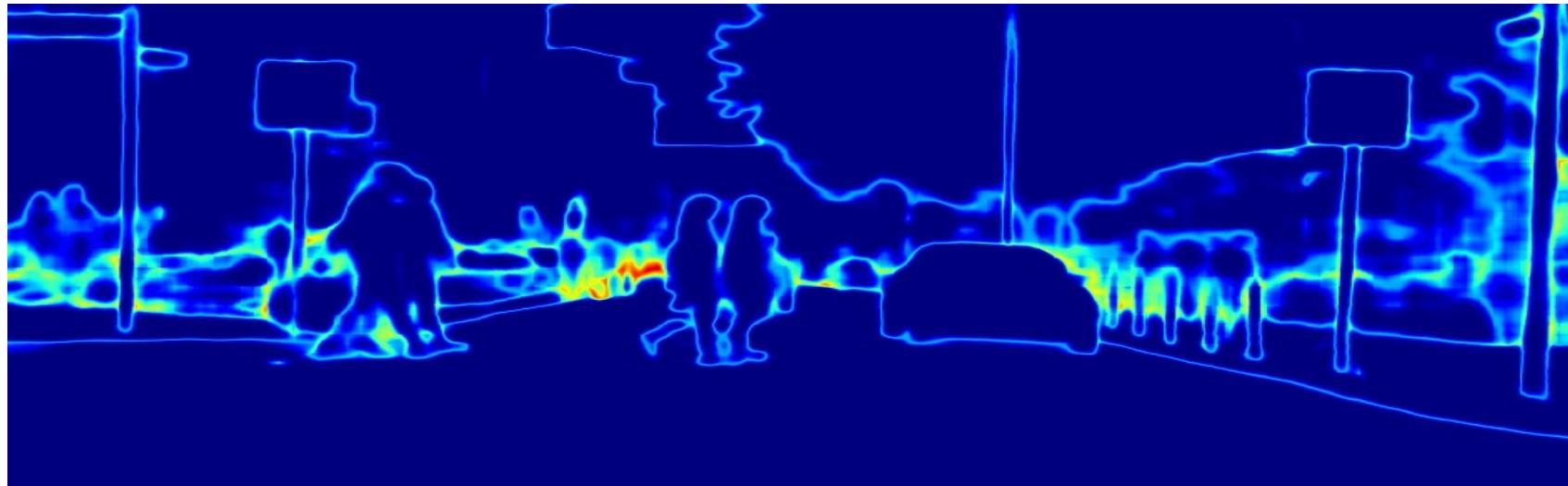
Equidistant points in high dimensions

```
1 n <- 1000; p <- c(2, 32, 512); pair_dist <- list()
2 for(i in 1:length(p)) {
3   X <- matrix(rnorm(n * p[i]), n, p[i])
4   pair_dist[[i]] <- as.vector(dist(X));
5   pair_dist[[i]] <- pair_dist[[i]] / max(pair_dist[[i]])
6 }
```

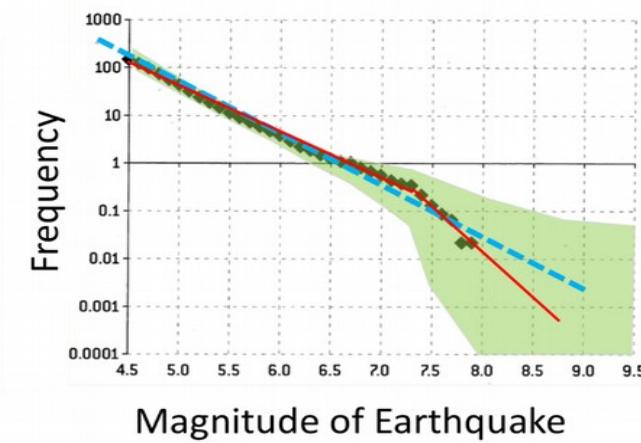
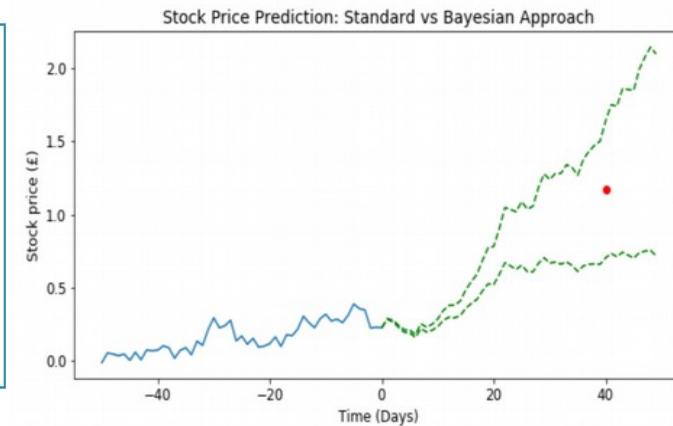


- Data points in high dimensions:
 - move away from each other
 - become **equidistant** and similar
- Impossible to see differences between cases and controls

When Frequentist stats are not good enough



Intelligence is to know how much you do not know



Regularizations are priors in Bayesian statistics

$$Y = \beta_1 X_1 + \beta_2 X_2 + \epsilon; \quad Y \sim N(\beta_1 X_1 + \beta_2 X_2, \sigma^2) \equiv L(Y | \beta_1, \beta_2)$$

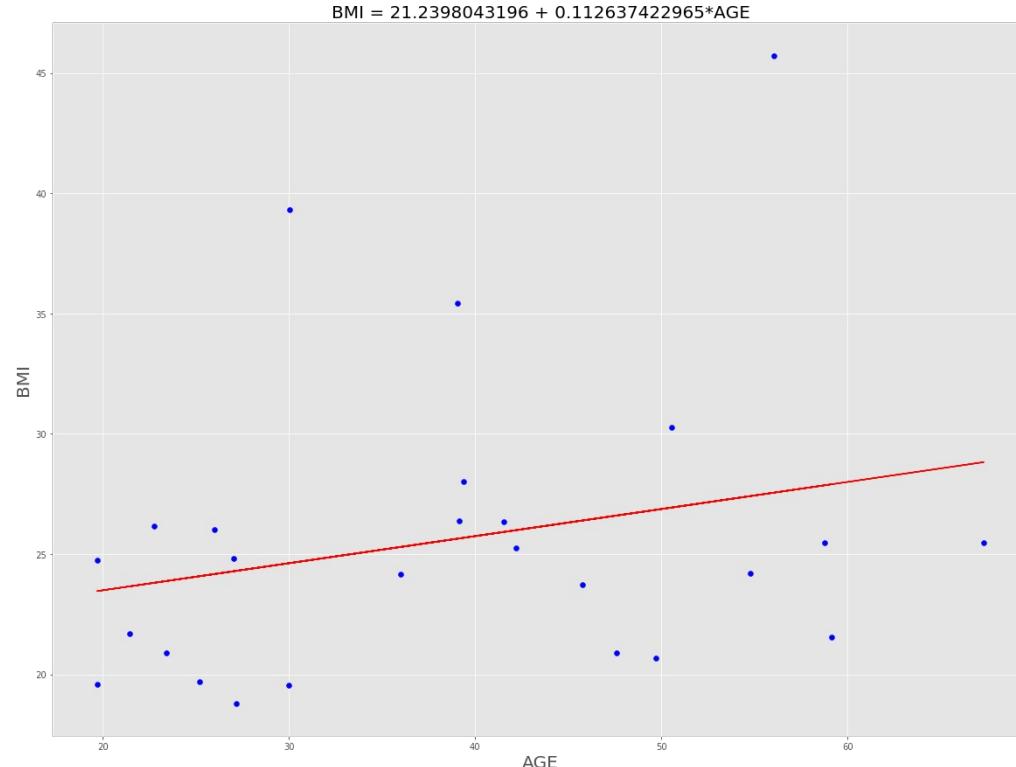
- Maximum Likelihood principle: maximize probability to observe data given parameters:

$$L(Y | \beta_1, \beta_2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(-\frac{(Y - \beta_1 X_1 - \beta_2 X_2)^2}{2\sigma^2} \right)$$

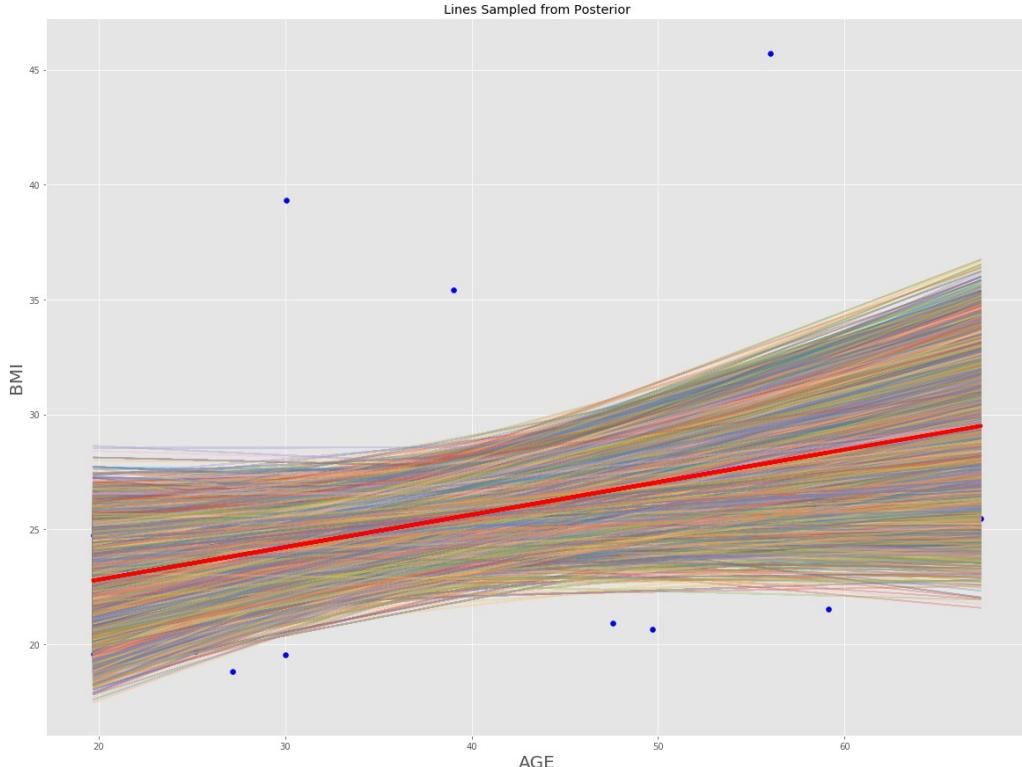
- Bayes theorem: maximize posterior probability of observing parameters given data:

$$\text{Posterior(params | data)} = \frac{L(\text{data} | \text{params}) * \text{Prior}(\text{params})}{\int L(\text{data} | \text{params}) * \text{Prior}(\text{params}) d(\text{params})}$$

Frequentist vs. Bayesian linear fitting



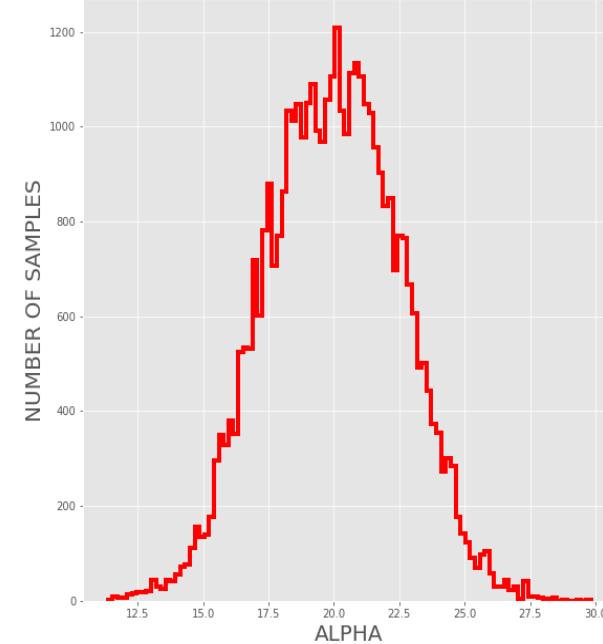
Fitting Frequentist linear model



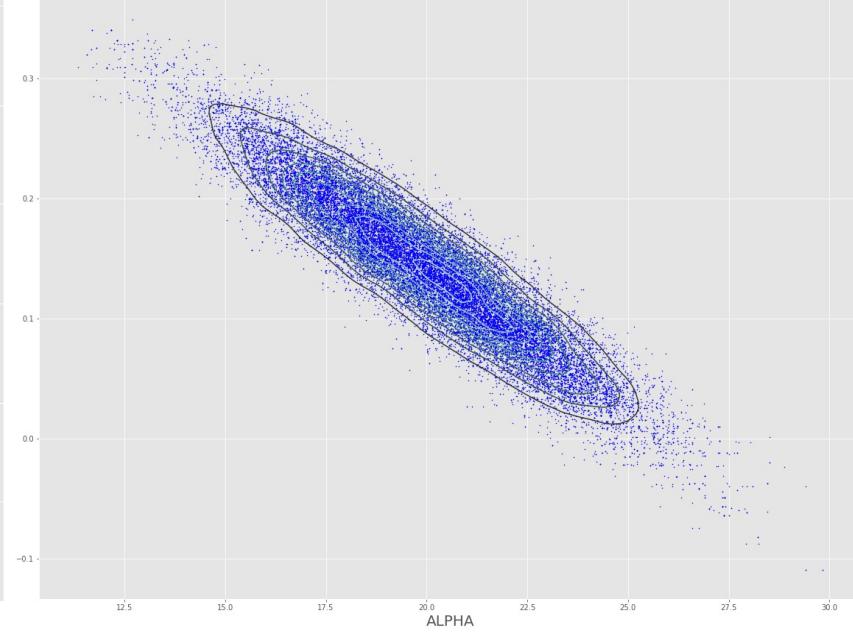
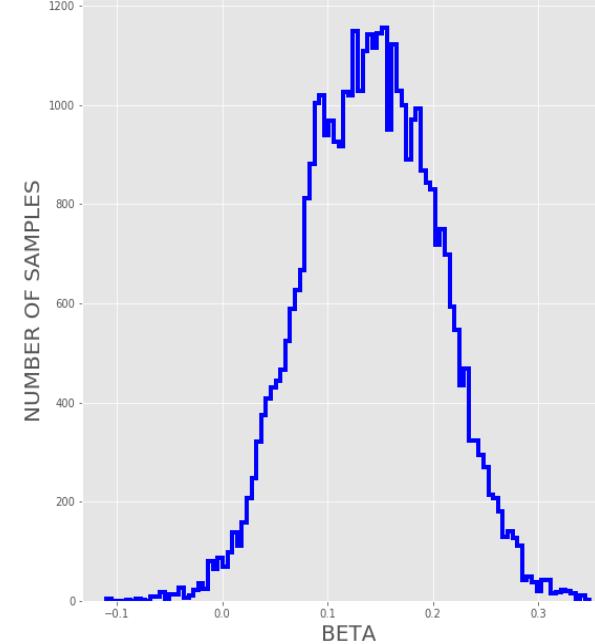
Fitting Bayesian linear model

Drawing possible fits from Posterior

POSTERIOR OF ALPHA



POSTERIOR OF BETA



Slope and intersect are not fixed values any more

MCMC can draw infinite number of fits from proximity of optimal slope and intercept parameters

Fashion MNIST



```
In [24]: # normalize inputs from 0-255 to 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

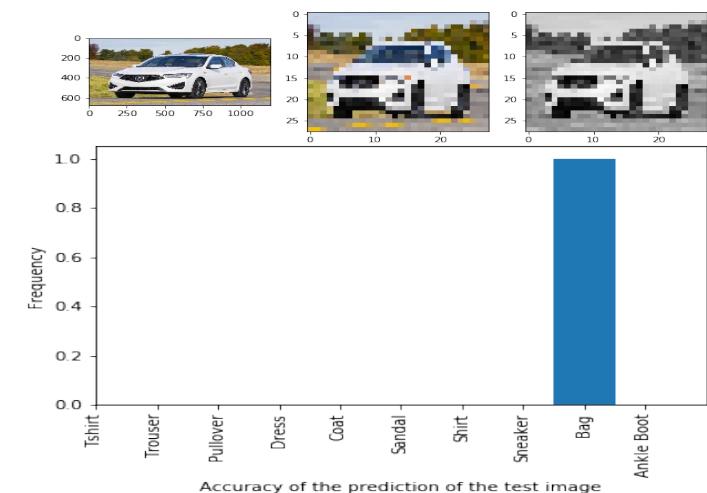
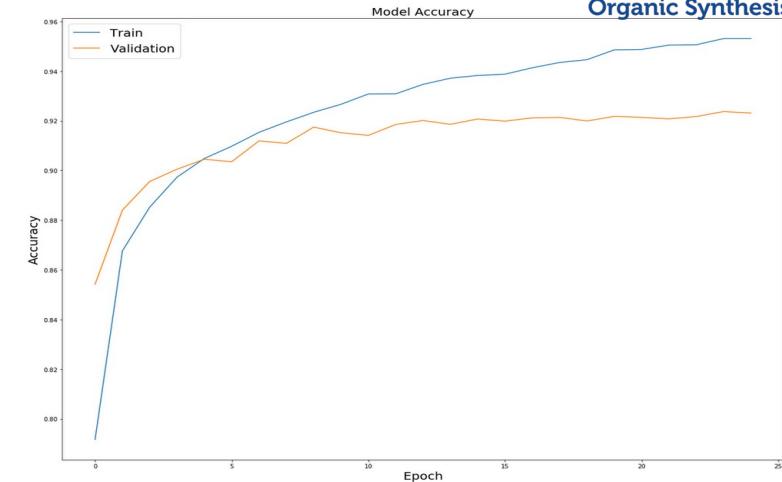
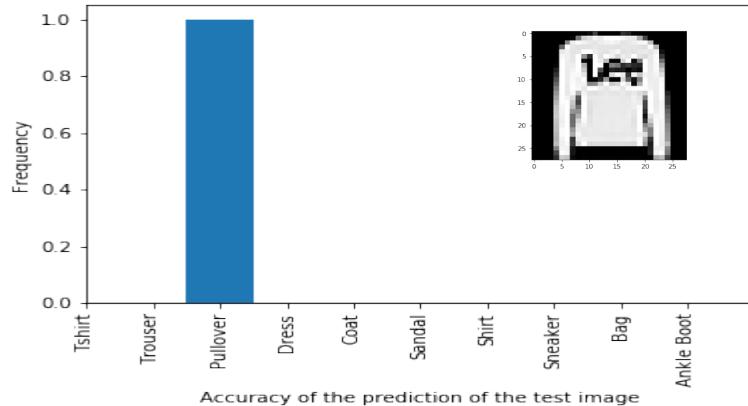
In [25]: # one hot encode outputs
y_train = np_utils.to_categorical(y_train)
y_test = np_utils.to_categorical(y_test)
num_classes = y_train.shape[1]
print(num_classes)
10

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(128, activation='relu', kernel_constraint=maxnorm(3)))
model.add(Dense(10, activation='softmax'))

# Compile the model
model.compile(optimizer='sgd', loss='categorical_crossentropy', metrics=['accuracy'])

Layer (type):          Output Shape:           Param #
conv2d_8 (Conv2D):      (None, 32, 28, 28)    320
dropout_7 (Dropout):   (None, 32, 28, 28)    0
conv2d_9 (Conv2D):      (None, 32, 28, 28)    9248
max_pooling2d_4 (MaxPooling2D): (None, 32, 14, 14) 0
flatten_4 (Flatten):   (None, 6272)            0
dense_7 (Dense):       (None, 512)             3211776
dropout_8 (Dropout):   (None, 512)             0
dense_8 (Dense):       (None, 10)              5130
=====
Total params: 3,226,474
Trainable params: 3,226,474
Non-trainable params: 0
None

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



Bayesian image recognition

PyMC3, Edward, TensorFlow Probability

```
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)
print(w.shape)
print(b.shape)
(60000, 784)
(10000, 784)
(10000, 10)

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)
print(w.shape)
print(b.shape)
(60000, 10)
(10000, 10)

In [10]: ed.set_seed(314159)
N = 100 # number of images in a minibatch.
D = 784 # 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, shape=[None, D])
# Normal(0,1) priors for the variables. Note that the syntax assumes TensorFlow 1.1.
w = Normal(loc=0, scale=tf.zeros([D, K]), scale_fn=lambda x: x)
b = Normal(loc=0, scale=tf.ones([K]), scale_fn=lambda x: x)
# 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 arr 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, i.e. minimise the KL divergence between q and p.
inf_inference = ed.inference.vi.QKL(qw, qb, data=(y: y_ph))
# Initialise the inference variable.
inference.initialize(n_iter=50000, n_print=100, scale=(y: float(x_train.shape[0]) / N))
# Set up the session.
sess = tf.InteractiveSession()
# Initialise the variable in the session.
tf.global_variables_initializer().run()

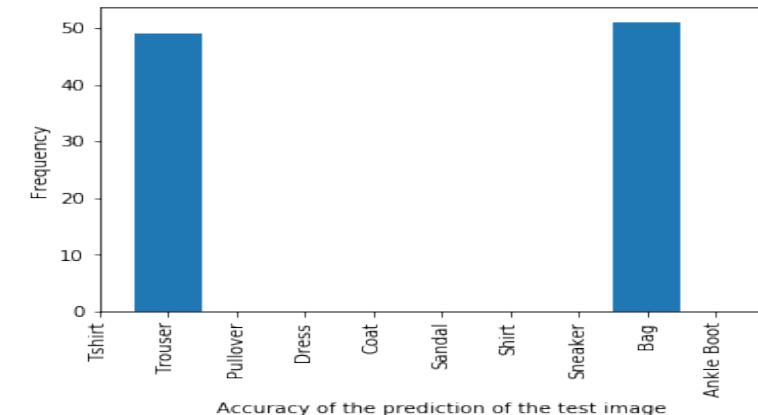
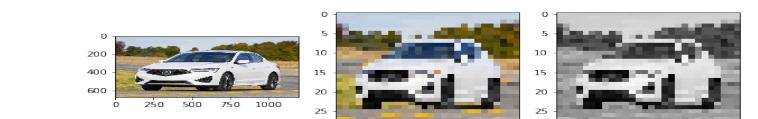
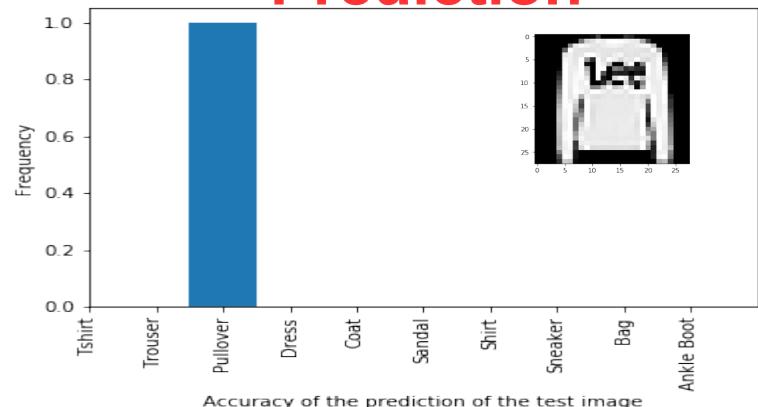
# Let the training begin. We load the data in minibatches and update the VI inference using each new batch.
for i in range(50000):
    X_batch, Y_batch = next(cifar10)
    X_batch = X_batch.reshape(N, -1)
    # We convert the raw pixel data into one hot vector format. We convert that into a single label.
    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)
    print("50000/50000 [100%] Elapsed: 221s | Loss: 85453.266")

In [14]: n_samples = 100
prob = []
samples = []
w_samples = []
b_samples = []
for i in range(n_samples):
    w_samp = ed.sample(w)
    b_samp = ed.sample(b)
    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 = tf.reshape(prob, [-1])
    sample = tf.concat((tf.reshape(w_samp, [-1]), b_samp), 0)
    samples.append(sample.eval())
    prob_hist.append(prob)

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 with maximum probability.
# We store the accuracy 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_hist:
    y_trn_prd = np.argmax(prob, axis=1).astype(np.float32)
    y_trn_label = np.argmax(y_test, axis=1).mean() * 100
    accy_test.append(accy)

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

Prediction





Session summary

Take home messages of the session:

- 1) Biological data are high dimensional and this should be taken into account when performing statistical analyses
- 2) The choice of analysis (e.g. Frequentist, Bayesian or machine learning) is driven by particular data type
- 3) Frequentist stats are focused on summary statistics and p-values which can be misleading
- 4) Bayesian stats deliver uncertainty which facilitates safer predictive analyses



Acknowledgments: LIOS + TARGETWISE



2027
National
Development Plan



Funded by
the European Union