

In this lecture, we deal with the issue of reconstructing missing values in our data set and with the problem of batch effects in the data set.

We will discuss the principles of two tools, ComBat and FunNorm that are widely used for removing batch effects.

Then we will also look at the tool BEclear from our group.

At the end I have summarized some basics from probability theory that are worth browsing over.

staphyType]	Fest Repor	t	MPSA (moch)	0	0	0	0	
		ā	WINSA (ITIECA)	0	0	0	0	0
Operator			PVL	0	0	0	0	0
Sample ID	2192119	201 A DOC 7D40 (ED0 2005 E 500000ED(206)	23S-rRNA	1	1	1	1	1
Date of Recult	2192119 - {+	10-16-01 2011						
Assay Name	StanhyType	10.40.01 2011	вард	1	1	1	1	1
Assav ID	10248		katA	1	1	1	1	1
Well Position	01 (01-A)					1	1	1
Software Version	2009-07-09		COA	1	0	1	1	1
Device	04a0022		Protein A	1	1	1	1	1
			sbi	1	1	1	1	1
nternal Controls			nuc	1	1	1	1	1
Data Quality	passed		fnbA	1	1	1	1	1
			vraS	1	1	1	1	1
enetic markers for	S. aureus / MR	SA / PVL	sarA	1	1	1	1	1
	0	to difference a second s	eno	1	1	1	1	1
atonomy	Species Mar	ker (S. auveta) positive	caos	1	1	1	1	1
MRSA (mecA)	positive		saes	0	1	1	1	1
	in Bauto		meca	0	-	0	0	0
esistance Conotan			DIAZ	0		0	0	0
cesistance Genotyp	e .		Diai	0	1	0	0	0
Iybridisation (Gene)	Result	Expected Resistance	blaR	0		0	0	0
necA	positive	Methicillin, Oxacillin and all Beta-Lactams, defining MRSA	ermA	0	0	0	0	0
olaZ	negative	Beta-Laktamase						
ma.	positive	Macrolide, Lincosamide, Streptogramm	ermB	0	0	0	0	0
mic .	negative	Macrolide, Efficosamide, Streptogramm	ermC	0	0	0	0	0
inA	negative	Lincosamides	linA	0	0	0	0	0
		Compute Euclid $ a - b _2 = \sqrt{\sum_{i=1}^{n}}$	dian distance $(a_i - b_i)^2$	betwee	en sar	nples		

First, we will look again at the microarray data set that we discussed in the first lecture.



The image-reader device generates 3 sorts of output "positive" (dark circle), "negative" (white field), and "ambiguous" for fields that cannot be precisely determined.

There are various possible reasons why about 5% of the microarray tests yielded ambiguous densities.



In the large scale project discussed in the first lecture, ambiguous values are disturbing the process of data analysis. They need to be cleaned up and replaced by either "positive" or "negative" values.

A simple approach would be to replace them either by the average signal of the data points for this particular gene probe ("gene average"), or by the average of the data points in this particular sample ("sample average") or by the average value of the full data matrix.

One could even compute the average of these 3 averages -> b_prediction.

Because we can only deal with 0 or 1 entries, the computed averages need to be thresholded by a suitable value, e.g. 0.5. Averages below 0.5 would be set to 0, those above 0.5 to 1.

We tested how well this works for some randomly selected data points. If we regenerate their entries and compare them to the correct values, this gives an agreement of 85% which is much better than random (50%).

We will now introduce a method that uses latent factor models that even generates predictions that are about 95% correct.



Latent Factor Models are very successful in image reconstruction.

If we delete 90% of the data points, the upper row shows that SVD is not useful for reconstructing the missing values.

However, a LFM can recover enough contrast so that we can recognize the face in this picture.



This slide illustrates the principles of LFM.

The idea is to represent the data matrix D_{ij} as the product of two matrices L and R. Once L and R are found, they can be used to compute all missing data points.

The algorithm iteratively refines guesses for L and R so that the squared difference of their product from the known data points is minimal.

Since this problem is usually underdetermined, there would be many different equally good solutions.

Therefore, one also applies the principle of regularization meaning that the algorithm constructs L and R in a way so that their norm is minimal.

A parameter lambda controls the balance between the two terms.



BEclear implements a stochastic gradient descent algorithm following a classic paper by Koren et al. This paper has been cited more than 9000 times.

It mentions that there are two popular approaches to solve the minimization task, stochastic gradient descent and alternating least squares (ALS).

Gradient descent is a well-known algorithm for optimization. An initial guess is iteratively refined by taking small steps along the direction of steepest descent which is the direction where the negative of the first derivative of the objective function is largest.

In **stochastic gradient descent**, the actual gradient that is calculated from the entire data set is replaced by an estimate of the gradient that is calculated from a randomly selected subset of the data.

MA assignment to clonal complexes + LFM predictions confirmed by WGS

154 *S. aureus* isolates (182 target genes) from Germany-vs-Africa study Table 1A

				Fu	nctional Cate	gory of genes			8
Result	Category	Res	ult caused by	Identification Regulation Resistance Virulen		Virulence	Total 11,374	% Total	
Concordant Positive		Microarray and WGS (de novo)		829	990	1,060		8,495	40.6%
n=27,119	Negative	Microarray a	and WGS (de novo)	0	1,159	8,100	6,486	15,745	56.2%
(96.8 %)									
Discrepant	False Positive	Microarray	Mishybridizations	0	78	21	103	202	0.7%
n=909 (3.2 %)		LFM	Misprediction	0	17	2	9	28	0.1%
	False Negative	Microarray	Polymorphisms	0	3	14	140	157	0.6%
		LFM	Misprediction	Ō	0	0	5	5	< 0.1%
		WGS	Assembly error	88	42	16	164	310	1.1%
			Cropped contig	1	12	15	28	56	0.2%
			Not sequenced or	6	9	8	100	123	0.4%
	Unknown		aberrant anere	0	0	4	24	28	0.1%
		Total num	ber of typing results	924	2,310	9,235	15,554	28,028	100%
Strauga	tal I Clin I	liorobial	(2016)	Very	few erro	rs due t	o LFM r	nis-pr	ediction
V3	tai. J Clin i		(2010) Processir	ng of Biologic 2021/22	al Data W	S			

In this comparison that was already shown in the first lecture, we were using data for 334 probe IDs from 154 isolates.

Out of this data, n = 2,788 or 5.4% of the hybridization signals were assigned as ambiguous value.

As just described, ambiguous were replaced by 1 or 0 values according to an LFM prediction based on the entries in neighboring fields of the involved columns and rows.

First, the accuracy of this approach was tested by a bootstrap approach as follows: 5% of randomly selected entries that were known to be positive or negative were removed from the dataset. This fraction corresponds to the typical number of targets typed as ambiguous in the microarray experiments. Then, these missing entries were predicted using LFM and were compared to the original values. As a result, LFM yielded an accuracy of 97% against the original values. Thus, the error rate of predicted values can be estimated as about 3%.

By comparison to the WGS data, LFM predictions were actually only wrong in 33 or about 0.1 % of the cases.



Now we come to the detection of batch effects.

A batch effect describes a case when a subgroup of measurements in the data set shows a qualitatively different behavior from the rest of the data.

Listed here are possible reasons why batch effects may occur.

This is a link to this paper: https://www.nature.com/articles/tpj201057

Another typical **error source** is mislabeling of samples either as healthy or as tumor. But this introduces random noise, not systematic deviations.



In these examples taken from the literature, significant batch effects can be seen by the perfect separation of different batches on the PCA score plots.

For the Hamner data set (panel B), batch effects exist with overlaps between several batches.



These two plots show batch effects due to using different fluorescent dyes and due to using different microarray platforms.



The left plot shows a box plot of microarray data. Each line represents the expression of all genes in one sample.

Obviously, the medians are very different. The left sample is highest.

The right plot shows the same data after RMA normalization. This algorithm uses quantile normalization.

Now, the distributions are very similar to each other.

Q	uanti	le no	rmal	isation: adju	sts	multi	ple d	listri	butions	
Give	<u>n</u> : 3 me	asurer	nents o	f 4 variables A – D) .					
<u>Aim</u> :	all mea	surem	ents sh	ould get identical	distril	outions	of valu	les		
Origin	al data				Determ	nine in ea	ch colun	nn the ra	ank of each va	ue
А	5	4	3		А	iv	iii	i		
в	2	1	4	\rightarrow	В	i	i	ii		
С	3	4	6		С	ii	iii	iii		
D	4	2	8		D	iii	ii	iv		
Sort c	olumns by	/ magnit	ude		Compu	ute mean	of each	row		
А	2	1	3	ς.	А	2	Rank	i		
В	3	2	4	—	В	3	Rank	ii		
С	4	4	6		С	4.67	Rank	iii		
D	5	4	8		D	5.67	Rank	iv		
	F 07	4.07	0	Deplese stisted valu						
A	5.67	4.67	2	according to the rank	of the	data fielo	les I.			
В	2	2	3		tain tha		luce			
С	3	4.67	4.67	(except of duplicates) so tha	at they ca	n be			
D	4.67	3	5.67	easily compared.		-				
V3				Processing of Biolo 2021/2	ogical D 22	ata WS				13

This slide reviews the quantile normalization method.

All data points are replaced by row averages so that the distributions become identical (except of duplicates).

Methods to correct batch effects
Available batch effect removal methods can be classified in 2 main approaches:
location-scale methods and matrix factorization methods.
The location-scale methods assume a model for the data distribution within
batches, and adjust the data within each batch to fit this model.
This approach is the most straight-forward one and many methods have been
proposed: ratio-based methods, ComBat, quantile based methods, mean or median
centering etc.
The matrix factorization based methods assume that the gene-by-sample
expression matrix can be represented by a small set of rank-one components which
can be estimated by means of matrix factorization.
The components that correlate with the batch number are then removed to obtain the
normalized dataset
Emilie Renard, PA. Absil 2017 IEEE International Conference on Bioinformatics and Biomedicine (BIBM), pp. 1511-1518, 2017
V3 Processing of Biological Data WS 14
2021/22

This is an overview over existing methods for removing batch effects.

Global methods to correct batch effects Mean-centering : after the transformation, the mean of each feature across all the samples within each batch is set to zero.	
Standardization: Beyond mean-centering, this approach normalizes the standard deviation of all features across samples within each batch to unity.	
The combination of these 2 yields the z-score of feature <i>i</i> : $\frac{x_{ij}-\mu_i}{\sigma_i}$	
Ratio-based: All samples are scaled by a reference array. This can be the average of multiple reference arrays, such as the measurement of universal human reference RNA samples for clinical data and vehicle control sample for toxicogenomics data.	s
Such global normalization methods do not remove batch effects if these affect specific subsets of genes so that different genes are affected in different ways.	
Luo et al. Pharmacogenomics J. (2010) 10: 278–291. V3 Processing of Biological Data WS 1: 2021/22	5

Listed here are three global methods that correct all data entries.



This is again the microarray data set that was normalized by RMA.

Although the overall distributions of the samples have been homogenized, there are hundreds of genes left that show clear batch effects.

Note that this plot shows the expression of individual genes.

If one clusters this normalized data (see right plot), the samples cluster according to processing date (green and orange represent two different dates).

This indicates that **RMA did not manage to remove the batch effect** for these genes.



This is another example for a large-scale batch effect in a famous genomic project.

For some reason, sequencing in the 1000 genome project generated higher read coverage during days 243 and 251.

ComBat

A widely used location-scale method is ComBat.

Here, the expression value of gene *i* for sample *j* in batch *b* is modeled as

 $X_{bij} = \alpha_i + \beta_i C_j + \gamma_{bi} + \delta_{bi} \varepsilon_{bij}$

where α_i is the overall gene expression, and C_j is the vector of known covariates representing the sample conditions (such as batch membership).

The error term ε_{bij} is assumed to follow a normal distribution N(0, σ_i^2). Additive and multiplicative batch effects are represented by parameters γ_{bi} and δ_{bi} .

ComBat uses a Bayesian approach to model the different parameters, and then removes the batch effects from the data to obtain the clean data:

$$X_{_{bij}}^* = \hat{\alpha}_i + \hat{\beta}_i C_j + \hat{\varepsilon}_{bij}$$

Emilie Renard, P.-A. Absil 2017 IEEE International Conference on Bioinformatics and Biomedicine (BIBM), pp. 1511-1518, 2017 See also discussion of Combat in Y.Zhang et al. BMC Bioinformatics 19, 262 (2018) V3 Processing of Biological Data WS 2021/22

A widely used tool for removing batch effects is ComBat. It is a location-scale method. The slide explains the basic principles of ComBat.

18

Although it is widely used, experience has shown that ComBat also has caveats. The listed Zhang paper

(https://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-018-2263-6) discusses some of them. For example, ComBat removes batch effects impacting both the means and variances of each gene across the batches. However, in some cases, the data might require a less (or more) extreme batch adjustment. Also, ComBat suffers from sample 'set bias', meaning that if samples or batches are added to or removed from the set of samples on hand, the batch adjustment must be reapplied, and the adjusted values will be different–even for the samples that remained in the dataset in all scenarios.



After a high-throughput study has been performed, the statistical approach for dealing with batch effects consists of two key steps.

Exploratory analyses must be carried out to identify the existence of batch effects and quantify their effect, as well as the effect of other technical artefacts in the data.

Downstream statistical analyses must then be adjusted to account for these unwanted effects.



In late 2011, we started working with DNA methylation data from the TCGA breast cancer study.

Soon we detected a severe batch effect that only affected segments on the Illumina chips.

In our 2013 publication, we omitted all affected genes (ca. 25% of the data).

Later, we developed a method termed BEclear (stands for clearing of batch effects) and published that tool in 2016.



This is an example of exploratory analysis. The top left panel shows a boxplot of the DNA methylation data in different batches of the TCGA data set.

The top right panel shows hierarchical clustering of the same data. The middle right panel shows a PCA of the same data. The bottom right panel shows a density distribution plot.

All plots illustrate clearly that, in batch 136, the distribution of β -values of genes is shifted to larger values than in the other batches.

The per sample plot (<u>top left</u>) shows that the difference in batch 136 is not due to only one sample but exists in all but two samples from this batch.

Beclear: Identify batch effected genes	
(1) Compare the distribution of every gene in one batch to its distribut other batches using the nonparametric Kolmogorov-Smirnov (KS) test	ion in all t.
P-values are corrected by False Discovery Rate.	
(2) To consider only biologically relevant differences in methylation levels the absolute difference between the median of all β-values within a bas specific gene and the respective median of the same gene in batches.	s, identify atch for a all other
Beta-values range between 0 and 1. The exp. error was estimated as 5%).
-> Smaller variations are not considered meaningful.	
Therefore, only those genes that have a FDR-corrected significance below 0.01 (KS-test) AND a median difference larger than 0.05 are consi batch effected (BE) genes in a specific batch or sample.	e p-value dered as
V3 Processing of Biological Data WS 2021/22	22

We suspected that the batch effect of the analyzed data affected various genes on the chip in different ways.

Therefore, we first had to identify which genes contain data points that differ largely from the remaining data points.



Each batch is assigned a BEscore value that considers the number of BE genes in that batch and the magnitude of their batch effects.

The question was now which values should be replaced: only the individual data points of BE genes in this batch or all of them?

We reasoned that if a sample (or batch) has a BEscore that is significantly larger than the other BEscore values, all values of that sample (or batch) should be replaced by LFM predictions.

Comparison of BEscores is done using the tabulated Dixon test. This test considers the absolute difference (gap) between the outlier in question and the closest value to it relative to the range of values (max - min).



This figure shows the outcome of BEclear for the tumor data.



This figure shows the normalization result by the tool FunNorm, another tool.



FunNorm builds on the idea of quantile normalization. This is a link to the paper: https://genomebiology.biomedcentral.com/articles/10.1186/s13059-014-0503-2

It is particularly tailored to the Illumina 450k chip that detects methylation levels for 450.000 CpG sites in the human genome.

It also contains close to 1000 control probes that do not measure CpG methylation of the sample, but are used to test the correctness of the biochemical processing steps carried out.

Functional Normalization
For each high-dimensional observation \mathbf{Y}_i , we form the empirical quantile function $r \in [0,1]$ for its marginal distribution, and denote it by q_i^{emp} .
What is a quantile function ?
 The k-th percentile of a set of values divides them so that k % of the values lie below and (100-k)% of the values lie above. The 25th percentile is known as the lower quartile. The 50th percentile is known as the median. The 75th percentile is known as the upper quartile.
It is more common in statistics to refer to quantiles . These are the same as percentiles, but are indexed by sample fractions rather than by sample percentages.
V3 Processing of Biological Data WS 27 2021/22

The data sets Y to be analyzed are transformed into their quantile functions. Here, we review what quantiles of a data set are.



Let us look at the quantile function of the standard normal distribution (blue curve in the upper plot).

Its quantile function is shown below.

For p=0.5, the variable has a 50% chance to be smaller than 0 in the normal distribution. Thus 0 is plotted on the y-axis for p = 0.5.

For p=0.1, the variable has a 10% chance to be smaller than (about) -1.3 in the normal distribution. Thus -1.3 is plotted on the y-axis for p = 0.1

For p=0.05 (the normal significance threshold), the value is -1.7. It is not -2 as we are used to (deviations of at least two standard deviations yield a p-value of 0.05) because we are only looking at one tail of the distribution.



FunNorm considers the quantile functions of the methylation value in all samples and takes its mean. This is termed alpha.

Then FunNorm assumes that the quantile function of a particlular sample i shows variation due to the covariates and some error term.

Functional Normalization

 $\hat{\beta}_{j} \text{ for } j = 1, ..., m$ are estimated using regression from the values observed for the **control probes**. Assuming we have obtained estimates $\hat{\beta}_{j}$ for j = 1, ..., m, we form the functional normalized quantiles by $q_{i}^{\text{Funnorm}}(r) = q_{i}^{\text{emp}}(r) - \sum_{j=1}^{m} Z_{i,j} \hat{\beta}_{j}(r)$ We then transform \mathbf{Y}_{i} into the functional normalized quantity $\tilde{\mathbf{Y}}_{i}$ using the formula $\tilde{\mathbf{Y}}_{i} = q_{i}^{\text{Funnorm}}\left(\left(q_{i}^{\text{emp}}\right)^{-1}(\mathbf{Y}_{i})\right)$ This ensures that the marginal distribution of $\tilde{\mathbf{Y}}_{i}$ has q_{i}^{Funnorm} as its quantile function. 23 Processing of Biological Data WS

The aim is to subtract the variation due to the covariates Z. The coefficients are estimated based on the values observed for the control probes.

The control probes are explained in the supplementary material of the FunNorm paper: "For "Bisulfite Conversion I" probes, 3 probes (C1,C2,C3) are expected to have high signal in the green channel in case the bisulfite conversion reaction was successful, and similarly 3 additional probes (C4,C5,C6) are expected to have high signal in the red channel. We therefore consider these 6 intensities and take the mean as a single summary value."

So these probes do not detect methylation levels of CpG sites in the sample, but rather are a quality measure for the performed experiments.



Here, we generated synthetic data sets with "known" batch effects.

First, we determined the standard deviation of the methylation value of each promoter probe in level 1 adjacent normal samples (samples belonging to batch 136 were excluded due to the existing batch effect).

Then we randomly selected 8000 promoter probes (approximately 10% of all promoter probes present on the chip) and increased the methylation values of 4000 of these promoter probes by a specified multiple of their specific standard deviation plus a noise term. The original probe values before introducing the synthetic batch effect were considered as our gold standard.

Because the methods Funnorm, ComBat and SVA adjust all values, the summed deviation of the corrected values from the original values (y-axis) is quite large.

In contrast, BEclear modifies only the values that are affected by batch effects. Therefore, the summed deviations are much smaller.



Maybe the previous analysis was a bit unfair to the other methods.

Therefore, we now only inspect the deviation of the batch effected data points.

For small batch effects of 2 standard deviations or less (which is a typical magnitude), BEclear still produces the smallest deviations.

Only for larger deviations, BEclear-adjusted values differ more strongly from the original data that with the other methods.



Then, we considered the identities of differentially methylated genes in breast tumor samples vs. normal samples.

As gold standard reference, we used the list of differentially methylated probes identified in the unaffected data using the limma package.

Then, we designed a synthetic batch effect in a similar fashion as before and applied BEclear, RUVm, FunNorm, ComBat, and SVA to this data.

Then, again we identified differentially methylated genes in this BE-adjusted data with limma and compared the results to the original data.

Shown here is the accuracy defined as (TP + TN) / (TP + TN + FP + FN) for the different BE-adjustment methods.

BEclear yielded a similar accuracy as the RUVm method that is not explained.

Both methods were more accurate compared to all other methods.



Today, we started the lecture by discussing various approaches to reconstruct missing data points.

Then, we met the important problem of batch effects in the raw data. If one does not care about batch effects, the downstream analysis may be heavily corrupted.

Therefore, as a bioinformatician, it is your job to check for possible batch effects.

We discussed different approaches that are implemented in software tools for removing unwanted batch effects. In our view, there is no "best" tool.

Certain approaches will offer advantages in certain situations and will give mediocre results in other cases. Identifying the best suited tool depends on the data to be analyzed.



Here, I have compiled some basics from probability theory. Some of this will be considered as known to you in the following lectures.

Probably you know most of this already.

Quickly browsing over these slides will help you fresh up these things.



These are the 3 basic properties that every event space needs to fulfill.



These are the 3 basic conditions that any probability distribution must obey.

Conditional probability

The **conditional probability** of β given α is defined as

$$P(\beta|\alpha) = \frac{P(\alpha \cap \beta)}{P(\alpha)}$$

In other words: The probability that β is true given that we know α is the relative proportion of outcomes satisfying β among those that satisfy α .

From this we see that

$$P(\alpha \cap \beta) = P(\alpha)P(\beta|\alpha)$$

This equality is know as the **chain rule** of conditional probabilities.

More generally, if $\alpha_1,\,\alpha_2,\,\ldots\,\alpha_k$ are events, we can write

$$P(\alpha_1 \cap \alpha_2 \cap \dots \cap \alpha_k) = P(\alpha_1)P(\alpha_2 | \alpha_1) \dots P(\alpha_k | \alpha_1 \cap \dots \cap \alpha_{k-1})$$

V3

Processing of Biological Data WS 2021/22

Bayes rule

Another immediate consequence of the definition of conditional probability is **Bayes' rule**.

Due to symmetry, we can swap the 2 variables α and β in the definition

 $P(\beta|\alpha) = \frac{P(\alpha \cap \beta)}{P(\alpha)}$ and get the equivalent expression $P(\alpha|\beta) = \frac{P(\beta \cap \alpha)}{P(\beta)}$

If we rearrange, we get Bayes' rule $P(\beta|\alpha)P(\alpha) = P(\alpha|\beta)P(\beta)$ or

$$P(\alpha|\beta) = \frac{P(\beta|\alpha)P(\alpha)}{P(\beta)}$$

A more general conditional version of Bayes' rule where all probabilities are conditioned on some background event γ also holds:

$$P(\alpha|\beta \cap \gamma) = \frac{P(\beta|\alpha \cap \gamma)P(\alpha|\gamma)}{P(\beta|\gamma)}$$

V3

Processing of Biological Data WS 2021/22

Example 1 for Bayes rule

Consider a student population.

Let Smart denote smart students and GradeA denote students who got grade A.

Assume we believe that P(GradeA|Smart) = 0.6, and that we get to know that a particular student received grade A.

Suppose that P (Smart) = 0.3 and P (GradeA) = 0.2

Then we have $P(\text{Smart} | \text{GradeA}) = 0.6 \times 0.3 / 0.2 = 0.9$

In this model, an A grade strongly suggests that the student is smart.

On the other hand, if the test was easier and high grades were more common, e.g. P(GradeA) = 0.4, then we would get

 $P(\text{Smart} | \text{GradeA}) = 0.6 \times 0.3 / 0.4 = 0.45$ which is much less conclusive.

V3

Processing of Biological Data WS 2021/22

Example 2 for Bayes rule

Suppose that a tuberculosis skin test is 95% percent accurate.

That is, if the patient is TB-infected, then the test will be positive with probability 0.95 and if the patient is not infected, the test will be negative with probability 0.95.

Now suppose that a person gets a positive test result.

What is the probability that the person is infected?

Naive reasoning suggests that if the test result is wrong 5% of the time, then the probability that the subject is infected is 0.95.

That would mean that 95% of subjects with positive results have TB.

V3

Processing of Biological Data WS 2021/22

Example 2 for Bayes rule

If we consider the problem by applying Bayes' rule, we need to consider the prior probability of TB infection, and the probability of getting a positive test result.

Suppose that 1 in 1000 of the subjects who get tested is infected \rightarrow P(TB) = 0.001

We see that 0.001×0.95 infected subjects get a positive result and 0.999×0.05 uninfected subjects get a positive result.

Thus P(Positive) = $0.001 \times 0.95 + 0.999 \times 0.05 = 0.0509$

Applying Bayes' rule, we get $P(TB|Positive) = P(TB) \times P(Positive|TB) / P(Positive)$ = 0.001 × 0.95 / 0.0509 \cong 0.0187

Thus, although a subject with a positive test is much more probable to be TB-infected than a random subject, fewer than 2% of these subjects are TB-infected.

V3

Processing of Biological Data WS 2021/22

Random Variables	
A random variable is defined by a function that associates a value with each outcome in Ω .	
For students in a class, this could be a function f_{grade} that maps each student in the class (in Ω) to his or her grade (1,, 5).	
The event grade = A is a shorthand for the event $\{\omega \in \Omega: f_{grade}(\omega) = A\}$.	
There exist categorical (or discrete) random values that take on one of a few values, e.g. intelligence could be "high" or "low".	
There also exist integer or real random variable that can take on an infinite number of continuous values, e.g. the height of students.	
By Val(X) we denote the set of values that a random variable X can take.	
V3 Processing of Biological Data WS 2021/22	43



Marginal Distributions

Once we define a random variable X, we can consider the **marginal distribution** P(X) over events that can be described using X.

E.g. let us take the two random variables Intelligence and Grade and their marginal distributions P(Intelligence) and P(Grade)

Let us suppose that

P(Intelligence=high) = 0.3P(Intelligence=low) = 0.7

P(Grade=A) = 0.25P(Grade=B) = 0.37P(Grade=C) = 0.38

These marginal distributions are probability distributions satisfying the 3 properties.

V3

Processing of Biological Data WS 2021/22

Joint Distributions Often we are interested in questions that involve the values of several random variables.							
E.g. we might be interested in the event "Intelligence = high and Grade = A".							A".
In that case we need to consider the joint distribution $P(X_1,, X_n)$ over these two random variables. The joint distribution of 2 random variables has to be consistent with the marginal distribution in that $P(x) = \sum_{y} P(x, y)$.							
	Intelligence						
			low	high	Σ		
		A	0.07	0.18	0.25		
	Grade	В	0.28	0.09	0.37		
		С	0.35	0.03	0.38		
		Σ	0.7	0.3	1		
V3		Pro	cessing of Biological	Data WS 2021/22			46

Conditional Probability

The notion of conditional probability extends to induced distributions over random variables.

P(Intelligence|Grade=A) denotes the conditional distribution over the events describable by Intelligence given the knowledge that the student's grade is A.

Note that the conditional probability $P(\text{Intelligence}=\text{high}|\text{Grade}=\text{A}) = \frac{0.18}{0.25} = 0.72$ is quite different from the marginal distribution P(Intelligence=high) = 0.3.

We will use the notation P(X|Y) to present a set of conditional probability distributions.

Bayes' rule in terms of conditional probability distributions reads

$$P(X|Y) = \frac{P(X)P(Y|X)}{P(Y)}$$

V3

Processing of Biological Data WS 2021/22

Probability Density Functions

A function $p: \mathbb{R} \to \mathbb{R}$

is a probability density function (PDF) for \boldsymbol{X}

if it is a nonnegative integrable function so that $\int_{Va} \int_{Va} p(x) dx = 1$

The function $P(X \le a) = \int_{-\infty}^{a} p(x) dx$ is the **cumulative distribution** for X.

By using the density function we can evaluate the probability of other events. E.g.

$$P(a \le X \le b) = \int_{a}^{b} p(x) dx$$

V3

Processing of Biological Data WS 2021/22

Uniform distribution

The simplest PDF is the uniform distribution

<u>Definition</u>: A variable X has a uniform distribution over [a,b] denoted X \sim Unif[a,b] if it has the PDF

$$p(x) = \begin{cases} \frac{1}{b-a} & b \ge x \ge a\\ 0 & \text{otherwise} \end{cases}$$

Thus the probability of any subinterval of [a,b] is proportional to its size relative to the size of [a,b].

If b - a < 1, the density can be greater than 1.

We only have to satisfy the constraint that the total area under the PDF is 1.

V3

Processing of Biological Data WS 2021/22

Gaussian distribution

A random variable X has a Gaussian distribution with mean μ and variance σ^2 , denoted X $\sim N(\mu;\sigma^2)$ if it has the PDF

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

A standard Gaussian has mean 0 and variance 1.



Expectation

Let X be a discrete random variable that takes numerical values.

Then, the expectation of X under the distribution P is

$$\mathbf{E}_P[X] = \sum_x x \cdot P(x)$$

If X is a continuous variable, then we use the density function

$$\mathbf{E}_P[X] = \int x \cdot p(x) dx$$

E.g. if we consider X to be the outcome of rolling a good die with probability 1/6 for each outcome, then $E[X] = 1 \cdot 1/6 + 2 \cdot 1/6 + ... + 6 \cdot 1/6 = 3.5$

V3

Processing of Biological Data WS 2021/22



Properties of the expectation of a random variable

If X and Y are independent then $\textbf{E}[X \cdot Y] = \textbf{E}[X] \cdot \textbf{E}[Y]$

The conditional expectation of X given y is

$$E_P[X|y] = \sum_x x \cdot P(x|y)$$

V3

Processing of Biological Data WS 2021/22

Variance

The expectation of X tells us the mean value of X. However, it does not indicate how much X deviates from this value. A measure of this deviation is the **variance** of X:

 $Var_P[X] = \mathbf{E}_P[(X - \mathbf{E}_P[X])^2]$

The variance is the **expectation** of the **squared difference** between X and its expected value. An alternative formulation of the variance is $Var[X] = \mathbf{E}[X^2] - (\mathbf{E}[X])^2$

If X and Y are independent, then Var[X + Y] = Var[X] + Var[Y]

$$Var[a \cdot X + b] = a^2 Var[X]$$

For this reason, we are often interested in the square root of the variance, which is called the **standard deviation** of the random variable. We define

 $\sigma_X = \sqrt{Var[X]}$

V3

Processing of Biological Data WS 2021/22

Variance

Let X be a random variable with Gaussian distribution $N(\mu;\sigma^2)$.

Then **E**[X] = μ and Var[X] = σ^2 .

Thus, the parameters of the Gaussian distribution specify the expectation and the variance of the distribution.

The form of the Gaussian distribution implies that the density of values of X drops exponentially fast in the distance (x - μ) / σ .

Not all distributions show such a rapid decline in the probability of outcomes that are distant from the expectation.

However, even for arbitrary distributions, one can show that there is a decline.

The **Chebyshev** inequality states $P(|X - \mathbf{E}_P[X]| \ge t) \le \frac{Var_P[X]}{t^2}$

or in terms of $\sigma_{_{V3}}$

$P(|X - \mathbf{E}_P[X]| \ge k\sigma_X) \le \frac{1}{k^2}$ Processing of Biological Data WS 2021/22

Resources
Nice online resources on statistics:
https://www.khanacademy.org/math/statistics-probability
http://tutorials.istudy.psu.edu/basicstatistics/
https://stattrek.com/statistics/problems.aspx
V3 Processing of Biological Data WS 2021/22