

Today, in lecture #5, we will discuss the issue of identifying peaks in a series of data points.

This is a typical problem in diverse areas of bioinformatics and in data analysis in general.

Of course, there exist many different solutions.

Which one is most suitable for a particular problem depends a lot on the kind of data.



In computer science, one typically deals with very accurate data.

In the 1D example shown on the right, one can easily see that the red-circled entries in fields 2 and 5 are local peaks.

They fulfil the simple requirement that they shouldn't be smaller than their left and right neighbors.

Algorithms for finding peaks in such perfect data are described e.g. in the classic book by Cormen et al. with the title "Introduction to Algorithms".

In contrast, bioinformaticians must detect peaks in inherently "noisy" data = data that is subject to sizeable fluctuations due to biological and technical variation.



As first example, we will discuss the case of histone modifications.

These are an important type of epigenetic marks and consist of posttranslational modifications (methylation, acetylation, phosphorylation ...) of lysine and other amino acids in the N-terminal flexible tails of histone proteins.

Shown in the figure are the two marks H3K36 me3 (tri-methylation of lysine36 of histone #3) and H3K27me3 along the genome sequence.

Also marked are the exons of two genes, FBXO7 and SYN3. The vertical lines or bars indicate the position of exons.

H3K36me3 is typically enriched in the gene body region (inside the gene, not in its promoter or enhancer regions) and associated with active gene transcription.

H3K27me3 is typically a repressive histone modification of nearby genes.

Histone marks can be detected by the ChIP-seq method that will be explained on the next slide.



Experimentally, histone marks are nowadays ususally detected by the ChIP-seq method (Chromatin Immuno Precipitation followed by sequencing) that is illustrated on the left.

First, DNA is **crosslinked** to bound proteins e.g. by applying **formaldehyde**, see right figure.

Formaldehyde crosslinking is routinely employed for detection and quantification of protein-DNA interactions, interactions between chromatin proteins, and interactions between distal segments of the chromatin fiber.

The DNA-protein mixture is then sheared into ~500 bp DNA fragments by **sonication** (application of ultrasound, induces DNA vibrations) or by digesting the free DNA ends with the enzyme **DNA nuclease**.

Then, an antibody is added to the mixture that is attached to a bead that can later be used to "fish" the antibodies from the sample. One selects for this a particular antibody that binds selectively e.g. to a histone protein carrying a particular histone mark. The antibodies are then "fished" from the solution. Subsequently, the protein-DNA crosslinks are broken up and the DNA is sequenced.

One assumes that all DNA prepared in this way was bound e.g. to the histone protein carrying the particular histone mark.

This experimental strategy is quite labor intensive and costly.

Every histone mark needs to be detected in a separate experiment using a different special antibody.



Now we discuss the output of the final sequencing step of a ChIP-seq experiment.

One obtains sequencing reads that belong to the DNA sequences that were "protected" by the protein of interest (e.g. a histone protein) against digestion by DNA nuclease or against DNA breakage during sonication.

Thus, one can assume that these DNA sequences bind specifically to the protein of interest. Of course, these regions will not only consist of the DNA stretch that makes physical contacts with the protein. The regions will extend a bit further. The sequencing reads may also contain further regions that are included by accident (experimental noise or unspecific binding events).

Some of this noise can be suppressed by performing several replicate experiments.

One checks which regions show a higher coverage (enrichment) over the background of the full genome.



MACS is a very popular tool to detect peaks in ChIP data.

It considers the average read coverage in a window relative to the background.

The Poisson distribution (compare V4) is a statistical distribution that is often used to model stochastic processes.

Here, one assumes that obtaining NGS reads from a genomic sample is such a stochastic process.

Regions in the upper tail of the distribution (default 10⁻⁵) are reported as peaks. Needed for this is an estimate of the lambda parameter.

MACS does not use a uniform lambda for the full sample, but a local lambda for the local segment.

	GEM	BCP (TF)	BCP (Histone)	MUSIC	MACS2	ZINBA	TM
ocating the potential peaks							
High resolution	Yes	Yes	No	Yes	Yes	No	Yes
ChIP and input sample signals combined	No	No	No	No	No	Yes	Yes
Aultiple alternate window sizes	Yes	Yes	Yes	Yes	No	No	No
Ise of variability of local signal	Yes	Yes	Yes	No	Yes	Yes	No
anking of peaks							
inomial test	Yes	No	No	Yes	No	No	No
oisson test	No	Yes	No	No	Yes	No	No
Iormalized difference score	No	No	No	No	No	No	Yes
Ise of underlying genome sequence	Vee	No	No			NI-	N.L.
of or an activity for the body of the	I Co	140	INO	NO	No	NO	NO
Representative selection	from ov	ver 30 exi	sting tools.	No No	No No	No Yes	No
Representative selection	from ov	ver 30 exi	Yes	No No	No No	NO Yes	No
Representative selection	from ov	ver 30 exi	Yes Sting tools. Park J, Natur Thomas et al	re Reviews . Brief Bioir	No No Genetics, 1	^{NO} Yes 10, 669 (20 41–450 (2	009) 017).

Thomas et al.: https://www.ncbi.nlm.nih.gov/pubmed/27169896

This is a comparison of several tools that are used to identify ChIP-seq peaks.

GEM is a 2-step method. In the second step, GEM also considers the motif content of the analyzed sequences (red circle).



Presented here is a protocol to generate synthetic ChIP data.

Link to Zhang et al. paper:

https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.100015 8

Gap regions in the UCSC genome assembly are excluded. Also, repetitive regions are excluded (row 2).

In row 3 row, we place synthetic transcription factor binding sites that should be detected by the ChIP-seq protocol.

In row 4, we select a suitable (blue colored) probability distribution for the expected read coverage (looks like a Poisson distribution) of the background and assign a coverage to each sequence region. Based on this distribution, many regions will get an average (low) coverage. Few regions will get a high coverage (darker blue).

For the binding sites, we use a different (green) probability distribution for their coverage (row 5).

In row 6, the coverage of each binding site is adjusted to follow somehow a Gaussian profile.

Finally, in the bottom row, we generate synthetic sequence reads. Their coverage matches the previously assigned coverage values.



It makes quite a difference whether one assumes a uniform background or varying backgrounds. For a uniform background, every nucleotide position in the background is given one as its sampling weight. For a varying background, every adjacent 1-kb block in the background is given a random weight drawn from a pre-specified underlying distribution and all nucleotide positions in a block are assigned the same weight.

The authors distinguished 4 regions of varying tag counts: low / medium / high / ultra-high. Tag clusters with low and high (including ultrahigh) tag counts are almost certain to be background and binding sites, respectively. Because there is a mixture of signals, the true identities of the clusters with medium tag counts are much less certain, and thus some form of thresholding is necessary.



(left) The **sensitivity** is also called **true positive rate** (TRP) or **recall**. TRP = TP / P = TP / (TP + FN).

The more peaks exist (x-axis from the left to right: 10^1 , 10^2 , 10^3 , 10^4), the better all methods perform in terms of sensitivity.

(Middle) Precision PPV is also called positive predictive value. PPV = TP / (TP + FP)

Precision measures how many identified peaks are correct. Here, the performance decreases steadily from left to right. The more peaks exist, the more difficult it is to detect the correct ones.

(Right) The F1-score is a measure that combines sensitivity and precision. It is the harmonic mean of precision and sensitivity F1-score = $2 \cdot PPV \cdot TPR$ /(PPV + TPR)

Consequently, it shows an optimal performace near 10^2 reads.

This comparison was done based on a simulated data set for which the correct answer is known.



Here, ChIP-seq was used to identify binding positions of the transcription factor Tbx5 on the genomic DNA. This example shows real data.

The precise binding motif where a transcription factor binds to DNA is known for many transcription factors including Tbx5.

One can identify such motifs e.g. with the MEME tool by checking for often occurring DNA strings in the ChIP-data for this transcription factors.

Here, several methods can identify the precise location of about half the Tbx5 binding positions to about 10 bp and even more to about 100 bp.



This is the cumulative distribution of the plot on the previous slide. About 90% of the regions are detected within 1000 bp.



H3K36me3 is a mark that is characteristic for actively transcribed genes.



Summary by Thomas et al.



In the second example of this lecture, we will discuss the task of identifying peaks in mass spectroscopy data.

We have already introduced the basic principle of MS in lecture V2.

This is a quick reminder of the main principles.



Shown here is the MS spectrum of the simple alkane molecule pentane shown at the bottom.

A carbon atom has mass 12 Da, a hydrogen has mass 1 Da.

Hence, the mass of an intact pentane molecule with 5 carbon atoms and 12 hydrogens is $5 \times 12 + 12 \times 1 = 72$ Daltons.

This is the right-most peak in the upper spectrum. Apparently, this molecule was detected with charge z = 1, giving a m/z ratio of 72.

Also detected are peaks at 57 Da (4 carbons with 9 hydrogens – meaning that one of the terminal carbon atoms has 3 hydrogens attached to it, the other one has 2 hydrogens),

43 Da (3 carbons with 7 hydrogens), and at 29 Da (2 carbons with 5 hydrogens).

The peak at 43 Da is highest showing that ionization of pentane mostly produces fragments with 3 carbon atoms.



This is the main protocol for processing of raw MS m/z data and identification of peaks.

First, the raw data is smoothened (a -> b). This suppresses many small intensity peaks.

Then, (b -> c) a baseline signal is removed (this is high (4000 to 6000 intensities) at small m/z values, and converges to an intensity of around 1000 for large m/z.

This step makes sure that one can identify peaks against a uniform background intensity of 0.

S: B: P: Table I: Open sou	smoothing baseline o peak findi irce softwa	g strategy correction ng strate re packag	/ i strategy gy ;es for MS da	Peak detecti	
Program	s	В	Р		P1: SNR
Cromwell [12]	S7	BI	P1, P4	• Smoothing	P2: Detection/Intensity threshold
LCMS-2D [20]		B5	P1, P2	S1: Moving average filter	P3: Slopes of peaks
LIMPIC [21]	S4	B2	PI, P3	S2: Savitzky-Colay filter	P4: Local maximum
LMS [22]	S3	B2	P1, P4	S3: Gaussian filter	P5: Shape ratio
MapQuant [16]	\$1,\$2,\$3		P7	S4: Kaiser window	P6: Ridge lines P7: Model-based criterion
CWT [10]	S5	B4	P1, P6	S6: Discrete Wavelet Transform	P8: Peak width
msInspect [23]	S6	B2	P5	S7: Undecimated Discrete Wavelet	Fransform
mzMine [24]	SI, S2	-	PI, P2, P8	Baseline Correction	
OpenMS [15]	S5	B4	P7	B1: Monotone minimum	
PROcess [13]	SI	B2, B3	P1, P2, P5	B2: Linear interpolation	
PreMS [25]	S7	BI	P1, P4	B3: Loess B4: Continuous Wavelet Transform	Yang et al. BMC Bioinformatics (2009) 10 :4
XCMS [8]	S3		P1, P4	B5: Moving average of minima	18

In this benchmark, the authors compared 12 tools that use various strategies for smoothing (S), baseline correction (B) and for peak finding (P).



A typical approach for smoothing of the raw data is to replace actual values y(n) or y(t) by averages taken over a local region.

The simplest approach is a **"moving average filter**". Here, one simply adds the values of the *k* values to the left and the *k* values to the right to the central value and divides the sum by 2k + 1.

This average is then assigned as smoothened value to the central data point.

An alternative is applying a **Gaussian filter** that takes into account essentially all data points from -infinity to +infinity, but weights the contribution of each point by the negative exponential of its quadratic distance *t* to the central point (as in the Gaussian distribution). Again, this weighted average is assigned as smoothened value to the central data point.



Another smoothing method is to weight neighboring data points by a so-called Mexican-hat wavelet, see figure.

This belongs to the so-called continuous wavelet transforms (CWT).



Now we introduce different methods for identifying peaks in the smoothened data.

The SNR method tries to identify peaks as "signals" relative to the normal fluctuation ("noise") of the data.

The noise is identified e.g. as the area including most (95%) of the data points or as MAD (see lecture 4, slide 22).

The "Slopes of peaks" method inspects the shape of any peak.

Left slope and right slope need to be steeper (i.e. the first derivative of the signal) than a certain threshold.

This criterion was likely developed to prevent detection of very broad and slowly rising mountains.



A local maximum is simply the largest data point among all its neighbors.

The shape ratio requires that the peak area should exceed a certain threshold. This excludes peaks that appear like sharp needles.



The left example tests how well different peak detection methods can identify peaks in synthetically generated data.

The right example is an experimental benchmark data set of 246 given proteins that have been digested by trypsin.

On both examples, CWT (detecting a Mexican hat profile) worked best.



Now we will discuss a related example, detected peaks in 2D data from MS.

Precisely, the field of breathomics attempts to identify organic compounds in exhaled breath.

The aim is - as can be expected - to use this method as early detection for diseases of the individual.

Shown here is how the exhaled breath is analyzed by a MS instrument and then processed in several steps of data analysis.



If the sample contains many different species, their MS signals could largely overlap if we try to analyze them only in a 1D m/z spectrum.

Therefore, breathomics separates the data in two dimensions.

Along the y-axis, we plot the retention time how fast a substance passes a capillary column. One uses a 17 cm long, 3mm diameter column that contains about 1000 thin capillaries. This architecture largely increases the surface of the capillary walls. The walls are coated with a thin "stationary phase", often a silica polymer.

Along the x-axis, we plot a kinetic property measured by the mass spectrometer.



The reduced mobility K of an ion drifting through a buffer gas is related to the square root of the charge over mass ratio, see eq. (1).

Instead of the mass of the ion, one considers the "reduced mass" that is combined from the ion mass and the mass of the gas molecules in the buffer gas inside the mass spectrometer.

The details of converting K into the inversed reduced ion mobility are not relevant for us here.



This figure shows the raw data of an IM spectrum-chromogram from which we want to identify the peaks of individual organic molecules.

Remember, plotted on the y-axis is the retention time through the MCC capillary column in seconds. Compounds that pass quickly, will show up at the bottom (short retention times).

Plotted on the x-axis are signals with different reduced inverse mobilities. The MS measurements are carried out sequentially for different retention time points.

This spectrum is provided to us as an r x t matrix.

The brightest peak of the spectrum (colored in yellow) is a peak at x = 0.5 that is present at all retention times.

This RIP peak belongs to the ions of the carrier gas in the MS spectrometer and is not relevant for us.

The other yellow and red peaks shown right of the RIP are only present at one retention time.



These are different steps of breathomics analysis.

In step 1, the RIP peak is removed from the spectrum.

In step 2, the signal is denoised (smoothed) and the baseline is subtracted.

In step 3, the peaks of interest are identified, here marked by boxes.



This is a flowchart presented in the PhD thesis of Dr. Ann-Christin Hauschild who worked on this topic in the group of Dr. Jan Baumbach.

Jan Baumbach was previously a young group leader at CBI and is now a full professor at TU Munich.



Humans are best able to identify the most interesting peaks in such a complicated spectrum.



Dr. Hauschild compared different algorithms and their ability to precisely identify peaks.

A simple "local maximum search" identifies central points as peaks with higher intensity than that of all 8 neighboring points.

Even very tiny differences would then be reported as local maximum.

Therefore in a second step, "significant" maxima are identified as those points that are higher at least by a given minimal intensity threshold than their neighbors.

Merged	peak cluster localizatio	n (MPCL)	
The MPCL consists of tw	wo phases: (1) clustering and (2)	merging.	
(1) each data point in the either peak or non-peal	e chromatogram is assigned to o k .	ne of 2 classes,	
For this, one uses a clus distance metric of the in	stering method that is based e.g. tensity values.	on the Euclidean	
(2) neighboring data po belong to the same peal	ints that are both labeled as pea k and are merged together .	k can be assumed to	
(3) each peak of the ana point , i.e. the data point in this peak region.	alyzed measurement is character t, which has the smallest mean d	ized by its centroid istance to all other points	i
		PhD thesis Ann-Christin Hausc Saarland University (2016) PhD dissertation Sabine Bader Dortmund University (2008)	hild,
V5	Processing of Biological Data - WS 2021/22		32

Also clustering can be used to identify peaks.

	Watershed algorithm	
Here, the IMS chromatogr	am is treated like a landscape ir	ncluding hills and valleys.
The algorithm starts with a	a water level above the highest in	ntensity followed by a
continuous lowering of the maxima.	e level while uncovering more and	d more of the local
At each step, the new unc labeled neighbors. Those peak and receive a new la	overed data points are annotated data points that remain unlabeled bel.	d by the label of adjacent d are identified as a new
The highest data point am coordinate.	ong a set of new labeled position	ns denotes the peak
The algorithm stops if all d threshold.	lata points are labeled or the leve	el drops below a given
1/5	Processing of Biological Data WS 2021/22	PhD thesis Ann-Christin Hauschild, Saarland University (2016)
v 3	Frocessing of Biological Data - WS 2021/22	2 33

The watershed algorithm is a widely used algorithm in image processing: https://en.wikipedia.org/wiki/Watershed_(image_processing).

This is an overview of the algorithm when it is applied for peak detection.

Watershed algorithm: impleme	entation
The watershed algorithm can be implemented as a prior points (having 2D coordinates) are sorted by magnitude.	ity queue where all data
(1) The largest data point is extracted and labeled first.	
(2 - n) This is followed by the next largest point in the que	eue and so on.
- Each point drawn out of the queue is compared with its	neighbors in space.
 If the neighbors are of equal or larger value, the extrac same label as its largest neighbor. 	cted point is given the
(comment: if of equal value, neighbor has not necessarily been labeled	.)
 In contrast, if the data point is larger than its neighbors not been labelled sofar), the data point is given a new lab part of another peak. 	(i.e. the neighbors have bel to indicate that it is
(n + 1) This procedure is repeated until the queue is emp	ty.
V5 Processing of Biological Data - WS 2021/22	Latha et al. Journal of Chromatography A, 1218 (2011) 6792– 6798 34

The Watershed algorithm was adapted for 2D chromatographic peak detection by S. Reichenbach, M. Ni, V.V.A. Kottapalli, Chemom. Intell. Lab. Syst. 71 (2004) 107.

Peak model estimation	l				
In the PME method, the expectation maximization (EM) algorithm is used to optimize the parameters of a mixture model from a given set of starting values.					
The algorithm requires a given set of "seed" coordinates modeled.	for each peak to be				
In general, any peak detection method is suitable to prov However, the quality of the results strongly depends on the approach.	ide these initial " seeds' ne chosen seed-ing				
Utilizing the EM algorithm, each peak is described by a n of two shifted Gaussian distributions and an additional pe	nodel function consisting eak volume parameter.	g			
Finally, the set of model functions plus a noise componer MCC/IMS measurement.	nt describe the whole				
V5 Processing of Biological Data - WS 2021/22	PhD thesis Ann-Christin Haus Saarland University (2016)	schild, 35			

The PME method will not be explained in detail here.

LMS : Automated Iod	cal b	reatho	mic	S				
maxima search	Table 6.1: Number merging the peak	er of peaks o lists (postpro	letected ocessing).	by all m	ethods.	Numb	er of peak c	usters after
			Method	# Pea	uks #	Peak C	lusters	
Mi CE : Automated		Manual (Vis	ualNow)	16	61		41	
peak detection via		LMS	8 (PeaX)	14	77		69	
merged peak cluste	r	MPCL (Vis	ualNow)	42	92		88	
leadization support	ad	PM	(IPHEX) (PosX)	50 13	97		420	
localization support	ea _	1 311	(i eax)	10	00		03	
by VisualNow								
	Table 6.2: Overla	ap of the five	e peak d	etection	method	s. The	overlap of th	e peak list
	A (row) and peak	: list B (colu	mn) is d	efined as	the nu	mber of	peaks in V t	hat can be
WST : Automated pe	ak ^{mapped to at leas}	st one peak i	n W. No	ote that t	he resu	lting ma	pping count	able is not
detection via water	symmetric.							
	Metho	d Manual	LMS	MPCL	WST	PME	Software	
shed transformatio	Manu	al 1661	911	1522	1184	791	VisualNow	
implemented in IPHE	x LM	S 868	1477	1096	1074	1128	PeaX	
	MPC	L 2667	2233	4292	2341	2082	VisualNow	
	WS	Γ 1112 Γ 797	1009	1157	5697	912 1959	IPHEX Doo V	
DME : Beek medel	1.1/1	E 131	1080	900	920	1998	reax	
PIVIE : Peak model								
estimation approact	ו							
by the PeaX tool								
by the reaction.					PhD) thesis	Ann-Christin	Hauschild,
VE	Dessertion	of Dielester	Data N	UE 2021/	Saa	riand U	niversity (20	16) p.95
V0	Processing	of Biologica	Data - V	vs 2021/	22			36

Thesis: https://publikationen.sulb.uni-saarland.de/handle/20.500.11880/26718

Table 6.1: The number of identified clusters varies between 41 and 88 except for the Watershed algorithm WST

Table 6.2: The overlap between the peaks identified by different methods is quite reasonable.



Testing of the peak annotation was performed using samples containing known reference molecules.

This is similar to the spike-in protocol presented in lecture #4, slide 38.



Signals #5 and #14 - #17 were not part of the reference analyte mixture, but could be clearly identified as decanal, n-nonan and heptanal.

They are components in many fragrances and could have entered the IMS from the room air.



It would be great if one could use breathomics for detection of complicated diseases.

Obvious candidates that may affect the composition of exhaled breath are lung diseases.



The software was tested on a public MCC/IMS dataset of COPD patients and healthy controls.



This study is described in https://www.mdpi.com/2218-1989/5/2/344/htm.

In the spectra, characteristic peaks of 120 volatile organic compounds were identified that are present in at least three of the patients' measurements.

Then, the 120 metabolites were clustered by hierarchical agglomerative clustering (HAC) and Pearson correlation.

By a suitable clustering threshold, the set of metabolites was split into 40 subsets, one for each cluster of correlating metabolites.

All clusters with less than three compounds were excluded, yielding a total of 14 metabolite sets.

Using this data, COPD could be separated from healthy samples with good success (85-95% success).

A Random Forest classifier achieved the highest accuracy.

Summary						
Peak detection is a free	uent task in diverse areas of biology.					
The challenge is posed by the noisy nature of biological data and the irregular shape of peaks.						
Testing and benchmarking of methods is typically done with synthetic (artificially generated) data.						
Peak detection and jud	ging their significance are equally important tasks.					
V5	Processing of Biological Data - WS 2021/22	42				

Today, we discussed examples ranging from identification of the peaks of certain histone marks over 1D mass spectroscopy to 2D MCC/IMS-based breathomics analysis.

These examples illustrated that one needs to adapt various peak identification methods to the data type and problem being studied.