

Kapitola 3

Periodický klíč

Slabost Caesarovy šifry spočívá v tom, že existuje pouze 25 možných dešifrování a kryptoanalytik je tak může všechny vyzkoušet. Situaci mu lze zřejmě znesnadnit tak, že zvýšíme počet možností, které musí vzít v úvahu, aby šifru vyřešil. Toho můžeme dosáhnout například tím, že jednotlivá písmena abecedy neposunujeme o stále stejný počet míst, ale tento počet měníme v závislosti na poloze písmene v otevřeném textu. Pochopitelně musí být předem dáno nějaké pravidlo, podle kterého posunutí volíme, aby oprávněný příjemce šifrové zprávy byl schopen text dešifrovat. Na tomto pravidle se musí odesílatel a adresát zprávy dohodnout předem. Jednoduché pravidlo je například používat periodicky několik různých posunutí. Ukážeme si to na jednoduchém příkladě. Při šifrování textu

VESELE VANOCE

budeme místo stále stejného posunutí používat opakovaně posunutí o 14, 18, 4 a 11 míst. Dostaneme tak

VESELE VANOCE

JWPPZW ZLBGGP

Klíčem pro *šifrování* zprávy tak byla posloupnost 14,18,4,11. Pro dešifrování musíme příslušná čísla od šifrové zprávy zase odečíst, a protože počítáme modulo 26, je to totéž, jako periodicky přičítat čísla 12,6,22,15. Tato posloupnost je tak klíč pro *dešifrování* zprávy.

I kdyby měl kryptoanalytik podezření, že byla použita Caesarova šifra se čtyřmi různými posunutími, musel by vyzkoušet nejméně 25^4 možností, ve skutečnosti ještě více, protože aspoň jedno z posunutí může být o 0 míst. V případě krátkých zpráv, jako je tato, použití *hrubé síly*, tj. vyzkoušení všech možností, může být jedinou možností. Je-li ale zpráva dostatečně dlouhá vzhledem k počtu použitých posunutí, tak je poměrně snadné určit jednot-

livá použitá posunutí, jak si brzy ukážeme. Naopak v extrémním případě, kdy se délka zprávy rovná počtu použitých posunutí a velikost těchto posunutí je generována nějakým náhodným procesem, tak je šifra nerozluštitelná.

Zesílení Caesarovy šifry využitím několika různých posunutí je používáno už několik století. Takové šifrovací systémy jsou známé jako *Vigenèrovy šifry*. V češtině je také používán název *šifry s periodickým klíčem*. Pro většinu lidí je snazší pamatovat si slova než posloupnosti čísel. Tomuto slovu pak říkáme *klíčové slovo* nebo také *heslo*. To sice poněkud redukuje množství možných klíčů, je to ale cena, kterou je nutné zaplatit za snazší zapamatování si klíče. Jednotlivá písmena klíčového slova pak nahradíme čísly podle převodové tabulky ze závěru první kapitoly. Tak například posloupnost 14,18,4,11 dostaneme z klíčového slova OSEL. Šifrování pak probíhá tak, že nad otevřený text napíšeme opakovaně heslo a obě posloupnosti sečteme tak, jako jsme to dělali v Příkladu 1.2.

OSELOS	ELOSEL
VESELE	VANOCE
JWPPZW	ZLBGGP

Vigenèrovy šifry jsou dosti speciálním případem *polyalfabetických systémů*. Tyto systémy používají mnoho různých abeced k náhradě jednotlivých písmen otevřeného textu. Počet těchto abeced může být od dvou do mnoha tisíc. Například německý šifrovací systém Enigma používal 16 900 různých abeced a navíc šlo o jednoduché záměny, nikoliv o pouhá posunutí abecedy jako v případě Caesarovy nebo Vigenèrovy šifry.

Jak vyřešit Vigenèrovu šifru, tj. šifru s periodickým heslem

První krok při řešení Vigenèrovy šifry spočívá v určení délky klíče. Je-li k dispozici dostatečně dlouhý šifrový text, pak postupujeme tak, že v něm hledáme opakující se *polygramy*, stejné kombinace aspoň dvou po sobě jdoucích písmen. Pokud tyto kombinace opravdu odpovídají stejným slovům nebo částem slov v otevřeném textu, pak jejich vzdálenost v šifrovém textu musí být násobkem délky klíče. Tímto způsobem najdeme délku klíče, případně ji redukuje na několik málo možností. Čím delší je opakovaný polygram, který v šifrovém textu najdeme, tím lépe. Ale i opakované bigramy mohou být k užitku.

Příklad 3.1 *Následující šifrový text vznikl z anglického textu, ve kterém byly mezery nahrazeny písmenem Z, pomocí Vigenèrovy šifry, tj. šifry s periodickým klíčem. Zjistěte délku klíče, klíč a původní otevřený text.*

HQEOT FNMKP ELTEL UEZSI KTFYG STNME GNDGL PUJCH QWFEX FEEPR
 PGKZY EHHQV PSRGN YGYSL EDBRX LWKPE ZMYPV EWLFG LESVR PGJLY
 QJGNY GYSLE XWYYP SRGFY KECVF XGFMV ZEGKT LQOZE LUIKS FYLXK
 HQWGI LF

Řešení. Prohlédneme šifrový text a zjistíme, že šest bigramů se vyskytuje aspoň třikrát:

- EL na místech 11, 14 a 140,
- FY na místech 23, 119 a 146,
- GN na místech 31, 64 a 103,
- HQ na místech 1, 40, 58 a 151,
- LE na místech 70, 91 a 109,
- YG na místech 24, 66 a 105.

Další zkoumání ukazuje, že bigram GN na místech 64 a 103 je v těchto dvou případech začátkem opakovaného oktogramu

GNYGYSLE

Je velmi nepravděpodobné, že by opakované oktogramy vznikly náhodně (Někdy ale k těmto velmi nepravděpodobným událostem dochází. Válečný kryptoanalytik Jack Good z Bletchley Park, kterého jsme zmínili v úvodní části přednášky, objevil opakovaný oktogram ve dvou různých zprávách. Pravděpodobnost, že by k něčemu takovému došlo náhodně, byla menší než $1 : 2^{10}$. Přesto se to stalo!), a tak budeme předpokládat, že toto opakování je významné. Vzdálenost mezi těmito oktogramy je $103 - 64 = 39 = 3 \cdot 13$, proto můžeme celkem spolehlivě předpokládat, že délka klíče je buď 3 nebo 13 nebo 39. Podíváme se proto ještě na vzdálenosti dalších opakujících se bigramů:

- výskyty EL mají vzdálenosti 3 a $126 = 3 \cdot 42$,
- výskyty HQ mají vzdálenosti 39, 18 a 93, všechny jsou násobkem 3.

To naznačuje, že 3 je zdaleka nejpravděpodobnější délkou klíče.

Další krok spočívá v nalezení samotného klíče. Pracujeme s hypotézou, že při šifrování byla použita tři různá posunutí. První posunutí bylo použito na písmena na místech 1, 4, 7, 10, atd., druhé posunutí na místa 2, 5, 8, 11,

atd., a třetí posunutí na místa 3, 6, 9, 12, atd. Šifrový text si proto napíšeme do třech sloupců podle následujícího vzoru.

1	2	3
H	Q	E
O	T	F
N	M	K
P	E	L
T	E	L
U	E	Z
S	I	K
T	F	Y
G	S	T
N	M	E
⋮	⋮	⋮

První sloupec má 53 písmen, druhý a třetí mají po 52 písmenech. Nyní spočítáme počty výskytů jednotlivých písmen v každém z těchto tří sloupců a zapíšeme je do tabulky:

sloupec	A	B	C	D	E	F	G	H	I	J	K	L	M	⋯	Z
1	0	1	0	0	0	3	13	4	0	0	1	7	1	⋯	1
2	0	0	0	0	13	6	0	0	3	2	2	1	2	⋯	1
3	0	0	2	2	4	1	1	1	0	1	5	5	1	⋯	3

Pokud by frekvence jednotlivých písmen byly náhodné, mohli bychom očekávat, že počty výskytů jednotlivých písmen v jednotlivých řádcích jsou přibližně 2. Protože ale čísla v jednotlivých řádcích odpovídají frekvencím písmen v přirozeném jazyce, můžeme očekávat počty od 0 do 10, přičemž nejvyšší výskyt s největší pravděpodobností odpovídá tomu písmenu v šifrovém textu, který šifruje písmeno Z nahrazující mezeru v otevřeném textu. To je proto, že každý sloupec je tvořený písmeny otevřeného textu posunutými o stejný počet písmen.

V prvním sloupci je nejčastější písmeno G, o kterém tedy předpokládáme, že je šifrovým ekvivalentem Z. První posunutí je tak pravděpodobně o 7 písmen. V tom případě by nejfrekventovanějšímu písmenu E v otevřeném textu mělo odpovídat písmeno L v šifrovém textu, které je skutečně druhým nejčastějším písmenem v prvním sloupci. To dále podporuje naši hypotézu, že první posunutí je o 7 písmen.

Ve druhém sloupci je nejčastějším písmenem E, které je tak vhodným kandidátem při šifrovou období písmena Z, tj. mezery, v otevřeném textu. To znamená, že druhé posunutí je pravděpodobně o 5 míst. Druhá dvě nejčastější šifrová písmena v druhém sloupci jsou F a Q, ze kterých se při posunutí o 5 míst zpět stanou písmena A a L, která jsou skutečně hojně frekventovaná v otevřených textech. Naproti tomu šifrová verze jiného hojně frekventovaného písmene v otevřených textech, písmene E, by po posunutí o pět míst vpřed byla J, které se ve druhém sloupci vyskytuje pouze dvakrát. Naše hypotéza, že druhé posunutí je o 5 míst dopředu, tak není příliš přesvědčivě podepřená.

Ve třetím sloupci není žádné šifrové písmeno příliš časté a tak nemáme žádný rozumný odhad, jaké je třetí posunutí. Nejčastějšími písmeny ve třetím sloupci jsou Y, K a L, jedno z nich bude patrně nahrazovat otevřenou mezeru Z, nevíme ale které. Zkusíme si proto napsat počátek šifrového textu a k němu odpovídající otevřená písmena, která dostaneme našimi odhady velikosti prvního a druhého posunutí.

HQEOTFNMKPELTELUEZSIKTFYGSTNMEGNDGLPUJCHQWFEXFEEPRPGKZY
AL.HO.GH.I .M .N .LD.MA. N.GH. I. G.NE.AL.Y .Y .IM. F.R

První slovo vypadá jako ALTHOUGH, a pokud tomu tak je, potom otevřené T je při třetím posunutí nahrazeno šifrovým E, což znamená posunutí o 11 míst dopředu. V tom případě by otevřené mezeře, tj. písmenu Z, odpovídalo šifrové písmeno K, které je skutečně jedním z nejpravděpodobnějších kandidátů pro šifrovou období mezery při třetím posunutí. Můžeme tedy shrnout, že *šifrovací* klíč je 7,5,11 a *dešifrovací* klíč pak 19,21,15. Pokud tento dešifrovací klíč použijeme, dostaneme otevřený text

ALTHOUGH I AM AN OLD MAN NIGHT IS GENERALLY MY TIME FOR
WALKING IN THE SUMMER I OFTEN LEAVE HOME EARLY IN THE
MORNING AND ROAM ABOUT FIELDS AND LANES ALL DAYS

což je počátek jednoho z románů Charlese Dickense. □

Stejně jako v případě jednoduché záměny jsme postupovali tak, že jsme zvolili několik nejpravděpodobnějších možností a potom jsme otevřený text doplnili na základě znalosti přirozeného jazyka. Také tentokrát nám při tom pomohla skutečnost, že jedno z písmen v otevřeném textu nahrazovalo mezeru.

Následující příklad, který opět připravil Mgr. Pavel Vondruška, ale ukazuje, že řešení šifry s periodickým klíčem je dosti jednoduché a vhodně napsaný program je zvládne během vteřiny. A to nemusíme ani umět jazyk, ve kterém je otevřená zpráva napsána. Stačí pouze vědět, o jaký jazyk jde.

Index coincidence

Důležitým pojmem je *index coincidence*, zkráceně IC. To je pravděpodobnost, že se u dvou textů v daném jazyce vyskytnou na stejném místě stejná písmena. Tato empiricky zjištěná (a v případě náhodného textu teoreticky vypočítaná) pravděpodobnost je v klasické kryptoanalýze hojně využívána. Následující tabulka ukazuje hodnoty IC pro některé vybrané jazyky při použití mezinárodní abecedy a bez mezer, pouze v případě ruštiny jde o abecedu s 32 znaky.

Jazyk	IC
Angličtina	0,0676
Francouzština	0,0801
Němčina	0,0824
Italština	0,0754
Španělština	0,0769
Ruština	0,0470
Čeština	0,0577
Slovenština	0,0581
Náhodný text	0,0385 = 1:26

Příklad 3.2 *Otevřený text je v češtině a je napsán bez mezer. Byla použita Vigenèrova šifra. Najděte otevřený text.*

ZFFLN QATOO AVFTS GQKZN MUXXB FJVVZ FBEPO FQKTN ADTCB OFLEB
 UQKZQ ATNKP HBGTV EQXNI GOTXB FFFLU UDDPP XZFAJ MEXGF PFODB
 WQKPT LLHFN MOBKE MTXGF ELNEF OOHDU UTHFU QABNJ BPSOF VJLEB
 HBCTP BSTGE AWRXJ YBMPN MUBVZ

Řešení. Celý text má 175 znaků. Nejdříve najdeme všechny opakované trigramy.

Trigram	poprvé	podruhé	rozdíl
FFL	2	77	75
LEB	48	148	100
NMU	20	170	150
QAT	6	55	49
QKZ	17	52	35
XBF	24	74	50
XGF	93	118	25

Vidíme, že s výjimkou trigramu QAT jsou všechny ostatní opakované tri-gramy ve vzdálenostech, které jsou násobkem 5. Nejpravděpodobnější délkou klíče je proto 5. Šifrový text si napíšeme do pěti sloupců.

1	2	3	4	5
Z	F	F	L	N
Q	A	T	O	O
A	V	F	T	S
G	Q	K	Z	N
M	U	X	X	B
F	J	V	V	Z
F	B	E	P	O
F	Q	K	T	N
A	D	T	C	B
O	F	L	E	B
U	Q	K	Z	N
A	T	A	K	P
H	B	G	T	V
E	Q	X	A	I
G	O	T	X	B
F	F	F	L	U
U	D	D	P	P
X	Z	F	A	J
M	E	X	G	F
P	F	O	D	B
W	Q	K	P	T
L	L	H	F	N
M	O	B	K	E
M	T	X	G	F
E	L	A	E	F
O	O	H	D	U
U	T	H	F	U
Q	A	B	A	J
B	P	S	O	F
V	J	L	E	B
H	B	C	T	P
B	S	T	G	E
A	W	R	X	J
Y	B	M	P	N
M	U	B	V	Z

Četnosti písmen ve sloupcích

	1	2	3	4	5
A	4	2	0	1	0
B	2	4	3	0	6
C	0	0	1	1	0
D	0	2	1	2	0
E	2	1	1	3	2
F	4	4	4	2	4
G	2	0	1	3	0
H	2	0	3	0	0
I	0	0	0	0	1
J	0	2	0	0	3
K	0	0	4	2	0
L	1	2	2	2	0
M	5	0	1	0	0
N	0	0	2	2	5
O	2	3	1	2	2
P	1	1	0	4	3
Q	2	5	0	0	1
R	0	0	1	0	0
S	0	1	1	0	1
T	0	3	4	4	1
U	3	2	0	0	3
V	1	1	1	2	1
W	1	1	0	0	0
X	1	0	4	3	0
Y	1	0	0	0	0
Z	1	1	0	2	2

Jako poslední krok řešení musíme rozhodnout, jaká posunutí abecedy byla použita v jednotlivých sloupcích. Tentokrát použijeme rovněž následující tabulku frekvencí pěti nejméně používaných písmen v jednotlivých jazycích. Tabulka nejčastěji používaných písmen ve stejných jazycích je v druhé kapitole.

Angl.	Franc.	Něm.	Čeština	Slov.
K: 0,41	Y: 0,21	P: 0,54	G: 0,48	G: 0,40
Q: 0,17	J: 0,19	J: 0,16	F: 0,33	F: 0,31
X: 0,17	Z: 0,07	Q: 0,01	W: 0,06	W: 0,06
J: 0,16	K: 0,00	X: 0,00	X: 0,04	X: 0,03
Z: 0,05	W: 0,00	Y: 0,00	Q: 0,00	Q: 0,00
Σ : 0,96	Σ : 0,47	Σ : 0,71	Σ : 0,91	Σ : 0,80

Pro ruční řešení si připravíme pomůcku v podobě proužku papíru, na který si do sloupce napíšeme jednotlivá písmena abecedy a vyznačíme si červeně pět nejčastěji se vyskytujících písmen v českých textech a modře pět nejméně často se vyskytujících písmen. Tento proužek papíru budeme postupně posunovat podél tabulky výskytu písmen v jednotlivých sloupcích a spočteme četnosti jednotlivých písmen vedle červených míst a od nich odečteme četnosti písmen vedle modrých míst. Dostaneme tak následující tabulku. V prvním sloupci je to písmeno otevřeného textu, které je vedle písmena A v šifrovém textu. Naším cílem je najít takové posunutí, při kterém je tento rozdíl maximální. Můžeme dokonce odhadnout, kolik asi by mělo toto maximum být. Protože součet frekvencí nejčastějších pěti písmen v českých otevřených textech, tj. písmen E, A, O, I a N je 41,07%, tj. 0,4107 a součet frekvencí pěti nejméně častých písmen je 0,91%, tj. 0,0091, je tento rozdíl pro jedno písmeno v průměru $0,4107 - 0,0091 = 0,4016$. Protože naše sloupce mají vždy 35 písmen, mělo by se maximum rozdílů v jednotlivých sloupcích rovnat přibližně $35 \cdot 0,4016 = 14,056$.

	1	2	3	4	5
A	-2	-4	-5	0	5
B	-1	-5	1	-5	-3
C	1	-7	1	-3	-8
D	5	3	9	0	1
E	0	-2	-5	1	-3
F	-5	-3	-5	3	0
G	-1	-6	-1	-4	1
H	2	-2	15	4	-2
I	6	4	1	-3	-1
J	2	3	-3	2	2
K	-1	8	-3	-2	0
L	-12	-3	2	8	-4
M	-1	1	-5	-8	3
N	7	-1	-3	-3	7
O	15	7	-2	2	1
P	-2	-2	-1	12	-2
Q	-7	-4	6	-1	-1
R	-8	-6	1	-8	-3
S	0	-1	-6	-7	-10
T	0	3	-3	-2	0
U	0	3	5	6	6
V	0	3	3	1	7
W	-2	-8	-3	0	-8
X	0	1	3	-3	-6
Y	1	4	2	8	0
Z	3	14	-1	2	18

Tučně vyznačená maxima se skutečně kolem odhadnuté velikosti průměrného maximálního rozdílu 14 pohybují. Šifrovému A tak v prvním sloupci nejspíše odpovídá otevřené O, ve druhém sloupci otevřené Z, ve třetím otevřené H, ve čtvrtém otevřené P a v pátém otevřené Z. V prvním případě tak otevřený text dostaneme z šifrového posunutím o 14, v druhém sloupci posunutím o 25, ve třetím posunutím o 7, ve čtvrtém posunutím o 15 a v pátém posunutím o 25. Klíč pro dešifrování je proto 14,25,7,15,25. Pokud tedy každou pěticí šifrového textu posuneme podle tohoto klíče, dostaneme otevřený text

NEMAME ZADNOU MIRU PRO MATEMATICKY TALENT PRIMOCARA CESTA I
PRO POSUZOVANI USPECHU NA MATEMATICKE OLYMPIADE VEDE VSAK
PRES ZKOUMANI ZDA SE VE SKUTECNOSTI SOUTEZICI POZDEJI

STAVAJI OPRAVDVYMI MATEMATIKY

Klíč pro šifrování tedy byl 12,1,19,11,1. Pokud jej převedeme na písmena, pak je to MBTLB. \square

Všimněte si rovněž, že od okamžiku nalezení délky klíče jsme už postupovali zcela mechanicky pouze s využitím tabulky četností písmen v příslušném jazyce. Ve skutečnosti nalezení délky klíče lze také udělat pomocí algoritmu založeného na výpočtu jistých statistických vlastností šifrovaného textu.

Srovnání Vigenèrovy a Caesarovy šifry

Přestože je řešení Vigenèrovy šifry o něco obtížnější než řešení Caesarovy šifry, na kterou stačí pouze hrubá síla, má tato šifra neodstranitelnou slabinu spočívající v tom, že jakmile se nám podaří odhalit šifrovou verzi jednoho často používaného písmene (případně mezery), dostaneme okamžitě všech zbývajících 25 písmen posunuté abecedy. Tato slabina zmizí, pokud místo posunuté abecedy o počet míst daných klíčem připustíme použití libovolné permutace abecedy, v každém sloupci jiné a nezávislé na permutacích v ostatních sloupcích. Tato úprava ale vede ke dvěma jiným problémům.

1. Jak takové nezávisle permutované abecedy získáme?
2. Jak je sdělíme zamýšlenému adresátovi zprávy tak, aby se je nikdo neoprávněný nedozvěděl?

Tyto otázky mají zásadní důležitost, protože jsou-li permutace abecedy získané nějakou jednoduchou metodou, například posunutím jako u Caesarových šifer, pak kryptoanalytik brzy tuto metodu odhalí a tím si luštění textu výrazně usnadní. Na druhou stranu ale zamýšlený adresát musí vědět, jaké permutace abecedy používáme, případně jak je může sám vytvořit. Existuje řada způsobů, jak oba problémy vyřešit a o některých si řekneme později. Na tomto místě si uvedeme dvě možná řešení problému (2). Před tím si ale uvedeme dva nové pojmy, které mají v kryptologii obecný význam.

Indikátory

Jestliže si odesílatel zprávy může v průběhu šifrování zvolit nějaké parametry šifry, které potřebuje sdělit příjemci, sdělí mu je pravděpodobně také v zašifrované podobě, jako *indikátor*, který samotné zprávě předchází nebo ji následuje, případně je skrytý uvnitř zprávy. Tak například klíč pro zašifrování 7,5,11, který jsme použili v prvním příkladu této kapitoly, musí být adresátovi nějak sdělen. Ať už jako nějaká dříve poslaná informace, nebo zašifrovaný nějakou jinou předem domluvenou šifrou, nebo skrytý uvnitř

zprávy na nějakém předem domluveném místě. V tomto případě samotný klíč 7,5,11 je indikátorem. Je ale velmi nepravděpodobné, že by byl adresátovi sdělován přímo v této podobě.

Hloubka

Pokud byly dvě nebo více šifrovaných zpráv zašifrovány stejným systémem se stejnými parametry (částmi, klíči, nastavením), pak říkáme, že *souvisí do hloubky* (aglicky je to *in depth*, na českém termínu se teprve budeme muset domluvit). Tak například, jsou-li dvě zprávy zašifrovány Vigenèrovou šifrou se stejným klíčem, pak souvisí do hloubky. Pokud jsou ale klíče různé, tak i kdyby měly oba klíče stejnou délku, tak zprávy do hloubky nesouvisí. Jsou-li oba klíče stejné délky a na některých místech jsou stejná písmena, pak říkáme, že šifrované zprávy souvisí částečně. To nemusí být pravda u jiných šifrovacích systémů, kde i ta nejmenší změna indikátorů způsobí, že šifrované zprávy spolu nesouvisí ani částečně. To, jestli kryptoanalytik dokáže využít skutečnosti, že přišel na to, že dvě nebo více šifrovaných zpráv souvisí do hloubky, záleží na systému, který byl k jejich šifrování použit. V případě Vigenèrovy šifry by toho nepochybně uměl využít. V případech jiných šifer takové zjištění příliš užitečné být nemusí. Velmi zhruba řečeno, pokud šifrovací systém je založený na záměně písmeno-za-písmeno, pak je možné odhalit, že dvě šifrované zprávy souvisí do hloubky a kryptoanalytikovi bude toto zjištění k užítku. Pokud ale šifrování používá dvě nebo více písmen otevřeného textu v jednom kroku, pak souvislost do hloubky, pokud ji lze vůbec rozpoznat, nemusí být příliš užitečná.

Jak poznat souvislost do hloubky

Jak kryptoanalytik pozná, že dvě zprávy souvisí do hloubky? Pokud byly vytvořeny stejným systémem a mají stejné indikátory, pak pravděpodobně souvisí do hloubky. Musím říkat *pravděpodobně*, nikoliv *jistě*, protože v intervalu mezi odesláním obou zpráv mohlo dojít k nějaké periodické změně šifrovacího systému. To se dělo například v německém systému Enigma, pokud jedna zpráva byla odeslána těsně před půlnocí a druhá krátce po půlnoci.

Jsou-li indikátory skryté, nemusí existovat žádný vnější důkaz, že zprávy souvisí do hloubky. Jak by na to kryptoanalytik vůbec mohl přijít? Pokud předpokládáme, že při šifrování je použita záměna písmeno-za-písmeno, napíšeme si obě zprávy pod sebe tak, aby začínaly současně, a spočítáme, na kolika místech mají stejná písmena. To znamená, že spočítáme IC, index koincidence. Pokud dva šifrované texty nesouvisí do hloubky, pak pravděpodobnost, že šifrované písmeno jedné zprávy je stejné jako šifrované písmeno

druhé zprávy na témže místě, je 1 : 26. Pokud obě šifrové zprávy souvisí do hloubky, pak se tato pravděpodobnost rovná indexu koincidence příslušného jazyka otevřených zpráv. Tento index jsme si již uvedli, v případě angličtiny je to zhruba 1 : 15, v případě češtiny nebo slovenštiny 1 : 17, v případě francouzštiny nebo němčiny 1 : 13. Pokud tedy šifrové zprávy nesouvisí do hloubky, můžeme očekávat, že se písmena na stejných místech objeví zhruba čtyřikrát za každých 100 písmen šifrových textů. Pokud souvisí do hloubky, pak se tato shoda bude objevovat přibližně šestkrát až sedmkrát v každých 100 písmenech, jsou-li oba otevřené texty v angličtině. Shody v bigramech nebo ještě delších polygramech jsou ještě významnější náповědou, že oba texty souvisí do hloubky. Jak ale ukazuje výše zmíněný příklad Jacka Gouda z doby druhé světové války, ani nalezení shodného oktogramu nemusí znamenat, že oba šifrové texty skutečně souvisí do hloubky.

Čím delší je kratší z obou zpráv, tím spolehlivější je předpověď založená na indexu koincidence.

Kolik textu potřebujeme pro rozluštění Vigenèrovy šifry

V prvním příkladu jsme měli text délky 157, ve kterém jsme našli řadu opakovaných bigramů, několik trigramů a dokonce jeden opakovaný oktogram (což je velmi nepravděpodobné). Klíč měl délku 3 a z pouhé frekvence mezery jsme ze sloupců délky zhruba 50 správně poznali ve dvou případech velikost posunutí, což v případě klíče délky 3 stačilo ke snadnému rozluštění celého textu.

Ve druhém příkladu jsme měli text délky 175 a klíč délky 5. Také v tomto případě počet písmen stačil k nalezení délky klíče a 35 písmen v každém sloupci stačilo k pohodlnému rozluštění celého šifrového textu. Čili lze říct, že text délky 50 krát délka klíče je možné snadno rozluštit, pokud známe pouze šifrový text. Pokud známe otevřený text a současně jeho šifrovou podobu, pak stačí dvakrát délka klíče.

Vigenèrova šifra je tak bezpečná pouze v případě velmi krátkých textů nebo velmi dlouhých klíčů, jejichž délka se blíží délce celého textu.

Cvičení 3.1 *Rozluštěte následující šifrový text, který byl vytvořen z anglického textu bez mezer Vigenèrovou šifrou. Napřed najděte délku klíče a potom také samotný klíč.*

```

ZZLES FMDCU LQMEW SGWLM XHZUY ZRJKU SGKBM GNBEU VPJCT VNGVW
HPLOY VLBAM RIIZN UJVKH XADVV GBQWX OOTKM RSEMV THMEU SNZMS
FHPPB KTQKK IZVPU ABGVT CWXKE FZNLY YVINE UTOGP MGCPM ESYBZ
OAVHG QDYOD ITKBC SUGPH VDGVP QDVLB NPFCP NYZQX QULBK GMIXI

```

BVCHR FYYWD OPEGL EGVCA QWMUE XBWXG KIIGH RTJIU WYYJB BSPPS
 VLTDO PLJNL DYODI TKBCS UGPHV DGVPQ DVJYN PFVPN YZQXW EMGYO
 GEFCH CMOEI VLGQE TWBWX GFANB RWECC KWLOK LRYGZ RHSKV EAVAB
 SVKLC XYWBA JPARK ZRGEW MBRZE RAWJR AGTZZ SENRP

Upozornění. Při přepisu šifrovaného textu bylo několik málo písmen chybně opsáno.

Statistické testy

Index koincidence IC je příkladem statistického testu, které jsou používány při řešení problému identifikace, při zjišťování, jaký šifrovací systém byl při šifrování zprávy použit. Ukážeme si nyní několik z těchto statistických testů, které vycházejí z indexu koincidence. Poté se podíváme, jak tyto testy použít při zjišťování délky klíče v případě Vigenèrovy šifry.

Začneme definicí Kroneckerova symbolu δ :

$$\delta(x, y) = \begin{cases} 1 & \text{pokud } x = y, \\ 0 & \text{pokud } x \neq y. \end{cases}$$

Jsou-li nyní $T = t_0 t_1 \cdots t_{M-1}$ a $T' = t'_0 t'_1 \cdots t'_{M-1}$ dva texty stejné délky M , tak můžeme definovat *index koincidence* těchto dvou textů jako číslo

$$Kappa(T, T') = \frac{1}{M} \sum_{v=0}^{M-1} \delta(t_v, t'_v).$$

Můžeme si spočítat index koincidence následujících dvou anglických textů délky 180.

T : THEPRECEDINGCHAPTERHASINDICATEDHOWAMONOALPHABETIC
 T' : WOULDSEEMTHATONEWAYTOOBTAINGREATERSECURITYWOULDBET
 * * *
 IPHERCANBESOLVEEVENIFTHEORIGINALWORDLENGTHSARECON
 OUSEMORETHANONEALPHABETINENCIPHERINGAMESSAGETHEGEN
 * * * * *
 CEALDANDTHESUBSTITUTIONALPHABETISRANDOMITISPOSSIB
 ERALSYSTEMCOULDBEONETHATUSESANUMBEROFDIFFERENTIALPH
 * * *
 LETOFINDASOLUTIONBYUSINGFREQUE
 ABETSFORENCIPHERMENTWITHANUNDE
 * *

Index koincidence těchto dvou textů se proto rovná

$$Kappa(T, T') = 14/180 = 0,0778,$$

což je zhruba o jedno procento více než je index koincidence anglického jazyka.

Očekávanou hodnotu indexu koincidence nějakého jazyka můžeme spočítat z frekvenční tabulky výskytu jednotlivých písmen v tomto jazyce. Označíme-li p_i pravděpodobnost výskytu i -tého písmene a_i v textech napsaných v daném jazyce, pak pravděpodobnost, že se i -té písmeno a_i vyskytne na témže místě ve dvou textech T a T' stejné délky, se rovná p_i^2 . Má-li abeceda příslušného jazyka N písmen (pro mezinárodní abecedu platí $N = 26$), potom očekávaná hodnota indexu koincidence těchto dvou textů se rovná

$$Kappa(T, T') = \sum_{i=0}^{N-1} p_i^2.$$

Takto jsme získali hodnoty indexu koincidence některých jazyků uvedené v této kapitole z frekvenčních tabulek výskytu jednotlivých písmen v odpovídajících jazycích, které jsou v předchozí kapitole o jednoduché záměně.

Testy *Chi*, *Psi* a *Phi*

Je-li $T = t_0 t_1 \cdots t_{M-1}$ text v nějakém jazyce, můžeme definovat jeho posunutí o r písmen jako text

$$T^{(r)} = t_{M-r} t_{M-r+1} \cdots t_{M-1} t_M t_{M+1} \cdots t_{M-r-1},$$

kde každý index $M - r + u$ interpretujeme modulo M . Posunutí o r písmen dostaneme tedy tak, že posledních r písmen textu T přesuneme na počátek (a zachováme jejich pořadí).

Jsou-li nyní $T = t_0 t_1 \cdots t_{M-1}$ a $T' = t'_0 t'_1 \cdots t'_{M-1}$ dva texty stejné délky M , spočítáme *průměrný index koincidence* $Chi(T, T')$ těchto dvou textů jako

$$Chi(T, T') = \frac{1}{M} \sum_{r=0}^{M-1} Kappa(T^{(r)}, T').$$

Odvodíme si nyní jednoduchou formulku, která nám umožní spočítat průměrný index koincidence dvou textů T a T' stejné délky, známe-li počty výskytů m_i a m'_i pro i -té písmeno a_i v obou textech T a T' pro každé $i = 0, 1, \dots, N - 1$. K tomu budeme potřebovat následující pomocnou funkci

$$g_{i,u} = 1, \text{ pokud } t_u = a_i, \\ 0, \text{ pokud } t_u \neq a_i.$$

Analogicky definujeme pro text T'

$$g'_{i,u} = 1, \text{ pokud } t'_u = a_i, \\ 0, \text{ pokud } t'_u \neq a_i.$$

Snadno si ověříme, že potom platí

$$\delta(t_u, t'_v) = \sum_{i=0}^{N-1} g_{i,u} \cdot g'_{i,v},$$

a také

$$m_i = \sum_{u=0}^{M-1} g_{i,u} \quad \text{a} \quad m'_i = \sum_{v=0}^{M-1} g'_{i,v}.$$

Nyní budeme počítat:

$$\begin{aligned} Chi(T, T') &= \frac{1}{M} \sum_{r=0}^{M-1} Kappa(T^{(r)}, T') = \\ &= \frac{1}{M} \cdot \frac{1}{M} \sum_{r=0}^{M-1} \sum_{v=0}^{M-1} \delta(t_{v-r}, t'_v) = \\ &= \frac{1}{M} \cdot \frac{1}{M} \sum_{u=0}^{M-1} \sum_{v=0}^{M-1} \delta(t_u, t'_v) = \\ &= \frac{1}{M} \cdot \frac{1}{M} \sum_{u=0}^{M-1} \sum_{v=0}^{M-1} \sum_{i=0}^{N-1} g_{i,u} \cdot g'_{i,v} = \\ &= \frac{1}{M} \cdot \frac{1}{M} \sum_{i=0}^{N-1} \sum_{u=0}^{M-1} \sum_{v=0}^{M-1} g_{i,u} \cdot g'_{i,v} = \\ &= \frac{1}{M} \cdot \frac{1}{M} \sum_{i=0}^{N-1} \left(\sum_{u=0}^{M-1} g_{i,u} \right) \cdot \left(\sum_{v=0}^{M-1} g'_{i,v} \right) = \\ &= \frac{1}{M} \cdot \frac{1}{M} \sum_{i=0}^{N-1} m_i \cdot m'_i. \end{aligned}$$

Platí proto

$$Chi(T, T') = \frac{1}{M^2} \sum_{i=0}^{N-1} m_i \cdot m'_i,$$

pro výpočet průměrného indexu koincidence dvou textů tak stačí znát počty výskytů jednotlivých písmen v obou textech. Pro dva anglické texty na předchozí stránce dostaneme následující počty výskytů písmen

	A	B	C	D	E	F	G	H	I	J	K	L	M
T :	15	6	8	9	21	3	4	10	17	0	0	8	2
T' :	15	6	4	5	30	4	4	9	8	0	0	6	6
	N	O	P	Q	R	S	T	U	V	W	X	Y	Z
	14	12	6	1	8	11	14	5	2	2	0	1	0
	15	12	4	0	10	10	17	8	0	4	0	3	0

Podle právě odvozeného vzorce pak dostáváme, že

$$Chi(T, T') = 0,066,$$

což je hodnota mnohem bližší hodnotě indexu koincidence pro anglický jazyk, která se rovná 0,0676, než dříve spočítaná hodnota indexu koincidence $Kappa(T, T') = 0,0778$.

Očekávanou hodnotu průměrného indexu koincidence $Chi(T, T')$ dvou textů T a T' v nějakém jazyce můžeme rovněž spočítat z frekvenční tabulky výskytu jednotlivých písmen v textech v tomto jazyce. Jestliže se pravděpodobnost výskytu písmene a_i rovná p_i , pak v textu délky M se očekávaný počet výskytů písmene a_i rovná $m_i = p_i M = m'_i$. Potom podle právě odvozeného vzorce dostáváme, že očekávaná hodnota průměrného indexu koincidence dvou textů T a T' se rovná

$$Chi(T, T') = \sum_{i=0}^{N-1} \frac{m_i^2}{N^2} = \sum_{i=0}^{N-1} p_i^2.$$

Očekávaná hodnota průměrného indexu koincidence se proto rovná očekávané hodnotě indexu koincidence stejného jazyka.

Pokud je $T = T'$, tj. $m_i = m'_i$ pro každé písmeno a_i , označujeme hodnotu $Chi(T, T) = Psi(T)$. Hodnota

$$Psi(T) = \sum_{i=0}^{N-1} \frac{m_i^2}{M^2}$$

se tak rovná průměrné hodnotě indexu koincidence mezi textem T a všemi jeho cyklickými posunutími $T^{(r)}$. Pro naše dva konkrétní texty potom platí $Psi(T) = 0,0655$ a $Psi(T') = 0,0745$. Očekávaná hodnota $Psi(T)$ se po dosažení $m_i = p_i \cdot M$ proto rovná

$$Psi(T) = \sum_{i=0}^{N-1} \frac{m_i^2}{M^2} = \sum_{i=0}^{N-1} p_i^2$$

a je to tedy opět hodnota indexu koincidence příslušného jazyka.

Nevýhodou výpočtu hodnoty $Psi(T)$ podle předchozího vzorce je skutečnost, že hodnota $Kappa(T^{(0)}, T) = Kappa(T, T) = 1$, a je tak podstatně větší než hodnota $Kappa(T^{(r)}, T)$ pro $r > 0$. Při výpočtu průměrné hodnoty indexu koincidence mezi textem T a všemi jeho posunutími $T^{(r)}$ je proto přirozenější vynechat hodnotu $r = 0$ a uvažovat pouze $M - 1$ případů pro $r = 1, 2, \dots, M - 1$. Tuto hodnotu označujeme $Phi(T)$. Výpočtem pro hodnotu $Phi(T)$ dostaneme vzorec

$$\begin{aligned}
 Phi(T) &= \frac{1}{M-1} \cdot \sum_{r=1}^{M-1} Kappa(T^{(r)}, T) = \\
 &= \frac{1}{M-1} \left(-1 + \sum_{r=0}^{M-1} Kappa(T^{(r)}, T) \right) = \\
 &= \frac{1}{M-1} (-1 + M \cdot Psi(T)) = \\
 &= \frac{1}{M-1} \left(-1 + \sum_{i=0}^{N-1} \frac{m_i^2}{M} \right) = \\
 &= \frac{1}{M-1} \cdot \frac{1}{M} \left(-M + \sum_{i=0}^{N-1} m_i^2 \right) = \\
 &= \frac{1}{M-1} \cdot \frac{1}{M} \left(-\sum_{i=0}^{N-1} m_i + \sum_{i=0}^{N-1} m_i^2 \right) = \\
 &= \frac{1}{M-1} \cdot \frac{1}{M} \left(\sum_{i=0}^{N-1} m_i(m_i - 1) \right).
 \end{aligned}$$

Pro naše dva konkrétní anglické texty T a T' potom dostáváme hodnoty $Phi(T) = 0,0603$ a $Phi(T') = 0,0693$.

Pro výpočet indexu koincidence jednoho textu T tak máme dvě možnosti. Buď spočítáme $Psi(T)$ nebo $Phi(T)$. Bez důkazu (který si snadno můžete doplnit sami) si uvedeme několik vztahů mezi těmito dvěma testy:

$$Psi(T) = \frac{M-1}{M} Phi(T) + \frac{1}{M}, \quad Phi(T) = \frac{M}{M-1} Psi(T) - \frac{1}{M-1},$$

a dále

$$Psi(T) - Phi(T) = \frac{1 - Phi(T)}{M} = \frac{1 - Psi(T)}{M-1}, \quad \text{a proto } Phi(T) \leq Psi(T).$$

Index koincidence jednoho textu T budeme definovat jako hodnotu $Phi(T)$. Stejně jako v případě ostatních statistických testů uvedených v této kapitole

můžeme spočítat očekávanou hodnotu indexu koincidence $Phi(T)$ textu T délky M z hodnot pravděpodobnosti p_i výskytu písmene a_i v příslušném jazyce. Dostaneme tak očekávanou hodnotu

$$Phi(T) = \frac{M}{M-1} \cdot \left(\sum_{i=0}^{N-1} p_i \left(p_i - \frac{1}{M} \right) \right).$$

Čím je text T delší, tj. číslo M větší, tím více se očekávaná hodnota indexu koincidence $Phi(T)$ blíží indexu koincidence příslušného jazyka.

Invariantnost indexu koincidence

Některé šifrovací systémy nemění index koincidence. Zejména to platí pro ty jednodušší systémy, kterými jsme se dosud zabývali. Platí totiž následující věta.

Věta 3.3 *Použijeme-li jednoduchou záměnu se stejnou permutací abecedy k zašifrování dvou textů T a T' stejné délky, dostaneme šifrované texty S a S' . Pro indexy koincidence těchto textů platí*

$$Kappa(T, T') = Kappa(S, S'), \quad Chi(T, T') = Chi(S, S').$$

Pro indexy koincidence jednoho textu T a jeho šifrované verze S platí

$$Psi(T) = Psi(S), \quad Phi(T) = Phi(S).$$

Formální důkaz nebudeme provádět, stačí si uvědomit, že pokud dva texty T a T' mají na stejném místě stejné písmeno, pak při použití stejné permutace abecedy budou mít oba šifrované texty S a S' na téže místě také stejné písmeno (obecně to bude jiné písmeno, než bylo na tomto místě v otevřených textech T a T').

Všimněte si, že z druhé části této věty vyplývá, že index koincidence šifrovaného textu zašifrovaného jednoduchou záměnou můžeme využít k odhadu, v jakém jazyce byl napsán otevřený text. Tak například šifrovaný text z Příkladu 2.2 má index koincidence Phi rovný 0,056, který se příliš neliší od indexu koincidence českých textů 0,0577. Pokud v šifrovaném textu z Příkladu 2.1 vynecháme písmeno Y, které označuje mezeru mezi slovy, má zbývající šifrovaný text index koincidence Phi rovný 0,0721, který se rovněž blíží indexu koincidence anglického jazyka 0,676. Větší rozdíl lze přičíst kratší délce textu. V každém případě je mnohem pravděpodobnější, že text z Příkladu 2.2 je šifrovou verzí českého textu a text z Příkladu 2.1 je šifrovou verzí anglického textu, než naopak.

Podobnou vlastnost invariance má i Vigenèrova šifra.

Věta 3.4 *Použijeme-li Vigenèrovu šifru se stejným klíčem k zašifrování dvou textů T a T' stejné délky, dostaneme šifrované texty S a S' . Pro indexy koincidence těchto textů platí*

$$Kappa(T, T') = Kappa(S, S'), \quad Chi(T, T') = Chi(S, S').$$

Použití indexu koincidence k odhadu délky klíče

V případě Vigenèrovy šifry můžeme pomocí indexu koincidence odhadnout, jakou délku měl použitý klíč. Pokud je délka klíče k , rozpadá se šifrovaný text do k skupin, z nichž v každé z těchto skupin jsou všechna písmena během šifrování posunuta o stejný počet písmen v abecedě. Každá skupina zvlášť tak má index koincidence stejný jako otevřený text v příslušném jazyce. Aritmetický průměr těchto r indexů koincidence se proto také přibližně rovná indexu koincidence příslušného jazyka. Hledáme-li délku klíče použitého při šifrování textu Vigenèrovou šifrou, vyzkoušíme různé délky r možného klíče, spočítáme indexy koincidence jednotlivých skupin a jejich aritmetický průměr. Jako nejpravděpodobnější délku klíče pak zvolíme to r , pro které se průměrný index koincidence nejvíce blíží indexu koincidence příslušného jazyka. Ukážeme si to na šifrovaném textu z Příkladu 3.2. Ten má 175 znaků a jejich výskyty najdeme v následující tabulce.

A	B	C	D	E	F	G	H	I	J	K	L	M
7	15	2	5	9	18	6	5	1	5	6	7	6
N	O	P	Q	R	S	T	U	V	W	X	Y	Z
9	10	9	8	1	3	12	8	6	2	8	1	6

Z této tabulky vypočteme hodnotu indexu koincidence šifrovaného textu Φ jako 0,04683. Tato hodnota odpovídá předpokládané délce klíče 1.

Nyní budeme počítat průměrnou hodnotu indexu koincidence pro předpokládanou délku klíče 2. V tomto případě si rozdělíme šifrovaný text na dvě části. V první budou všechna písmena na lichých místech šifrovaného textu a ve druhé budou písmena na sudých místech šifrovaného textu. První skupina se tak rovná

ZFNAO AFSQZ MXBJV FEOQT ATBFE UKQTK HGVQN GTBFL UDPZA MXFFD
WKTFL MBETG ENFOD UHUAN BSFJE HCPSG ARJBP MBZ

Tato část šifrovaného textu má 88 písmen a jejich výskyt je v následující tabulce

A	B	C	D	E	F	G	H	I	J	K	L	M
6	7	1	3	5	10	4	3	0	3	3	2	4
N	O	P	Q	R	S	T	U	V	W	X	Y	Z
4	3	3	4	1	3	6	4	2	1	2	0	4

Index koincidence Phi této části je proto 0,043365. Druhá část textu (na sudých místech) je

FLQTO VTGKN UXFVZ BPFKN DCOLB QZANP BTEXI OXFFU DPXFJ EGPOB
 QPLHN OKMXF LEOHU TFQBJ POVLB BTBTE WXYMN UV

Tabulka výskytu jednotlivých písmen v druhé části textu je

A	B	C	D	E	F	G	H	I	J	K	L	M
1	8	1	2	4	8	2	2	1	2	3	5	2
N	O	P	Q	R	S	T	U	V	W	X	Y	Z
5	7	6	4	0	0	6	4	4	1	6	1	2

a index koincidence Phi této části je 0,046779. Aritmetický průměr indexů koincidence obou částí je pak 0,4507.

Podobně spočítáme průměrné hodnoty indexu koincidence jednotlivých částí textu pro další předpokládané délky klíče. Celou proceduru lze samozřejmě algoritmizovat a naprogramovat. Dostaneme tak následující hodnoty průměrného indexu koincidence.

Délka klíče	Průměrné Phi
1	0,04683
2	0,04507
3	0,04685
4	0,04172
5	0,05647
6	0,04860
7	0,04476
8	0,03772
9	0,04763
10	0,05253
11	0,04827
12	0,04560
13	0,04832
14	0,04462
15	0,06303

Průměrná hodnota indexu koincidence Φ se nejvíce blíží hodnotě koincidence češtiny, která je 0,0577, pro délku klíče rovnou 5. To je rovněž nejpravděpodobnější délka klíče, kterou jsme dostali zkoumáním vzdáleností opakovaných trigramů, a která je skutečnou délkou klíče použitou při šifrování textu z Příkladu 3.2. Všimněte si rovněž, že další dvě nejbližší hodnoty průměrného indexu koincidence jsou pro délky klíče 10 a 15, což jsou násobky délky skutečně použitého klíče.

Cvičení 3.2 *Následující šifrový text byl vytvořen z českého textu bez mezer Vigenèrovou šifrou. Určete délku klíče, klíč a otevřený text.*

EFCWV SFXWQ JBTWB SDINZ AVDWW APWAA CZMKF SEPWX SRLGC JTWDC
 PKEEO FSPSB ZVSWZ CSCTM KRLGM CDMEY AVHGK QPVFC WFSJY NEERC
 CPMAD SYOQK XLDQW CKRSX TDECU SAVAJ CRVWB INISN GJTQX OYSZM
 JEAUR OKIDW HBUDS BJENW HKUHU MUECC QPZSN BJLGZ XKNAU RJRWQ
 ERTJS ZLRWN DLZAV OWINW ERDSV CESBY ZFKVI NDIHP DJTWZ ZTDAB
 DCUHY DCMKY KEYUR DZPKR PICAM VUECM CVJSU MGIDG ZPZSV IORQF
 BPZSZ CDXMN DMADY HZLDY KFUVY ZFNSI DBDQJ GTXGF AVDSM SDXMA
 XTHWW HFELY BKOZV OGYHM BRLMJ JPHDM IRKRO XTDNP PKIDY NNIDY
 SFZSN IYEVJ DLHWW VFFWL TDKJU ILLGZ XKINS RPPAX PSLWC YOZGS
 QCEVI QSWNC IPLWU CNMCR TIEGU OXDAI BYOMB OWCFY BVSAM OAECQ
 TIYUR ZPTJG ZIYDI JTGCW