

Shahmukhi Corpus Report

I. Character Analysis:

The character analysis has been performed on the Shahmukhi corpus at sentence level and below is the list of all characters in the corpus with percentage of their occurrence.

Shahmukhi Char	Unicode Decimal	Percentage	SN
ا	1575	13.21	1
س	1740	10.16	2
و	1608	7.59	3
ن	1606	6.53	4
ع	1746	5.53	5
و	1729	4.85	6
ج	1585	4.75	7
د	1583	4.55	8
ه	1722	4.42	9
ک	1705	4.12	10
ل	1604	3.32	11
ت	1578	3.24	12
ڌ	1726	3.04	13
س	1587	2.97	14
ڻ	1605	2.44	15
ڦ	1576	1.99	16
ڦ	1662	1.98	17
ڦ	1580	1.91	18
ڱ	1711	1.59	19
ڦ	1616	1.53	20
ڦ	1670	1.46	21
ڦ	1615	1.4	22
ڦ	1569	1.1	23
ڦ	1681	0.78	24
ڦ	1657	0.71	25
ڦ	1570	0.7	26
ڦ	1672	0.66	27
ش	1588	0.53	28
ف	1601	0.35	29
ح	1581	0.34	30
ع	1593	0.34	31
ق	1602	0.33	32
ز	1586	0.31	33
ڙ	1582	0.3	34

ص	1589	0.22	35
ـ	1614	0.17	36
ط	1591	0.13	37
غ	1594	0.09	38
ض	1590	0.08	39
ـ	1617	0.08	40
ظ	1592	0.07	41
ـ	1584	0.05	42
ـ	1579	0.05	43
ـ	1648	0.01	44
الله	65010	0.01	45
ـ	1556	0.01	46
ـ	1554	0.01	47
ـ	1611	0	48
ـ	1552	0	49
ـ	1574	0	50
ـ	1654	0	51
ـ	1688	0	52
ـ	1553	0	53
ـ	1555	0	54
ـ	1731	0	55
ـ	1610	0	56
ـ	1550	0	57
ـ	1572	0	58
ـ	1537	0	59
ـ	1620	0	60
ـ	1619	0	61
ـ	1622	0	62
ـ	1613	0	63
ـ	1612	0	64
ـ	1618	0	65
	65279	0	66

II. Unigram or Word Analysis:

The word analysis has been performed on the Shahmukhi corpus at sentence level and below is the list of top 100 words.

Shahmukhi Word	Percentage	SN
ت	2.54	1
ا	2.05	2
د	1.75	3
دی	1.61	4
دا	1.41	5
نون	1.29	6
ک	1.26	7
سی	1.13	8
وچ	0.95	9
نہیں	0.94	10
وی	0.94	11
اوہ	0.81	12
نال	0.79	13
ای	0.72	14
میں	0.7	15
اپہ	0.66	16
توں	0.66	17
نیں	0.65	18
تان	0.58	19
نے	0.57	20
نہ	0.54	21
وچ	0.52	22
اوس	0.47	23
پر	0.47	24
ہے	0.45	25
اک	0.44	26
کر	0.42	27
کوئی	0.42	28
کہ	0.41	29
اپس	0.41	30
ہو	0.4	31

اوہناں	0.38	32
گل	0.34	33
سن	0.34	34
گیا	0.34	35
لے	0.33	36
آپنے	0.24	37
گئے	0.24	38
مینوں	0.24	39
اُتے	0.23	40
پے	0.22	41
اوہدے	0.22	42
جے	0.21	43
میرے	0.2	44
جو	0.2	45
گھر	0.2	46
کیہ	0.2	47
اوہنوں	0.19	48
ہووے	0.19	49
آ	0.19	50
ہویا	0.19	51
ہور	0.18	52
اک	0.18	53
جی	0.18	54
گئے	0.17	55
کچھ	0.17	56
جاندا	0.17	57
جا	0.17	58
پھیر	0.17	59
دیان	0.16	60
آپ	0.16	61
بے	0.16	62
لگ	0.16	63
ہوندا	0.16	64
آپنی	0.16	65
ہر	0.15	66

کم	0.15	67
ہن	0.15	68
ہُن	0.15	69
آکھیا	0.15	70
آن	0.15	71
ہوئی	0.15	72
جوئیں	0.15	73
پیا	0.15	74
ول	0.14	75
اپناں	0.14	76
کرن	0.14	77
کسے	0.14	78
دیاں	0.14	79
جدوں	0.14	80
والے	0.14	81
بولی	0.14	82
ربیا	0.14	83
اسیں	0.14	84
کیتا	0.13	85
اندر	0.13	86
ویلے	0.13	87
دو	0.13	88
لے	0.13	89
بته	0.13	90
ہوئے	0.12	91
مان	0.12	92
کول	0.12	93
اوہنے	0.12	94
بندے	0.12	95
پتہ	0.11	96
اج	0.11	97
سارے	0.11	98
اوہبڈی	0.11	99
میری	0.11	100

III. Bi-Gram Analysis:

The bi-gram analysis has been performed on the Shahmukhi corpus at sentence level and below is the list of top 100 bi-gram words.

Previous word	word	Frequency	SN
	میں	8473	1
	اوہ	7923	2
	اپہ	7083	3
	پر	5917	4
	اوں	4744	5
	تے	4526	6
	ایس	4187	7
اے	تے	4181	8
نہیں	سی	3804	9
کر	کے	3759	10
جاندا	اے	2941	11
	اوہناں	2598	12
	اک	2563	13
سی	تے	2512	14
ہوندا	اے	2446	15
اے	کہ	2416	16
	جے	2302	17
	مینوں	2051	18
اوہناں	نوں	2033	19
اوہناں	دے	1843	20
ہو	گیا	1842	21
	بُن	1838	22
تے	اوہ	1793	23
	پھیر	1754	24
ہو	کے	1654	25
سی	کہ	1587	26
	کوئے	1569	27
	اوہنے	1553	28
	جدوں	1547	29
ای	نہیں	1544	30
	اسیں	1517	31

سکدا	اے	1511	32
بوندی	اے	1487	33
اوہناں	دی	1472	34
	اپہناں	1434	35
اوہس	نوں	1432	36
جاندی	اے	1406	37
	میرے	1393	38
دی	گل	1390	39
نیں	تے	1375	40
وی	نہیں	1355	41
	توں	1340	42
گیا	اے	1329	43
لے	کے	1299	44
دے	نال	1298	45
	اج	1203	46
ہو	گءی	1195	47
	اوہدے	1188	48
	کیہ	1176	49
اوہس	دے	1172	50
	جیوین	1151	51
اے	پر	1147	52
	ہر	1146	53
	نہیں	1142	54
	نہ	1113	55
گیا	سی	1113	56
جا	کے	1097	57
رہیا	سی	1054	58
ہویا	اے	1048	59
اوہناں	دا	1042	60
رہی	سی	1026	61
ویکہ	کے	1025	62
تے	میں	1025	63
تے	اوہس	1005	64
پتھ	نہیں	1002	65
تے	اوہناں	1002	66

	اوہنون	993	67
نوں	وی	982	68
جاندے	نیں	972	69
اے	پءےی	958	70
	میرا	957	71
آ	کے	953	72
تے	پھیر	952	73
اے	جو	951	74
لگ	پیا	944	75
مان	بولی	943	76
	اک	937	77
نال	ای	936	78
کوئےی	نہیں	926	79
اوس	نے	919	80
آپنے	آپ	916	81
سی	پر	911	82
تار	اوہ	900	83
	بے	886	84
گءی	اے	882	85
اوس	دی	880	86
گءی	سی	872	87
	سادے	863	88
سن	تے	857	89
	میری	855	90
ہویا	سی	845	91
کہ	اوہ	842	92
ہو	گءے	841	93
گیا	تے	838	94
اپہ	گل	838	95
ہو	سکدا	837	96
ایس	لءی	828	97
گل	اے	818	98
	ایتهے	814	99
کردا	اے	802	100

IV. Tri-Gram Analysis:

The tri-gram analysis has been performed on the Shahmukhi corpus at sentence level and below is the list of top 100 tri-gram words.

Previous to Previous Word	Previous Word	word	Frequency	SN
		میں	8473	1
		اوہ	7923	2
		ابہـ	7083	3
		پر	5917	4
		اوس	4744	5
		تے	4526	6
		ایس	4187	7
		اوہناں	2598	8
		اک	2563	9
		جـ	2302	10
		مینوں	2051	11
		بُن	1838	12
		پھیر	1754	13
		کوئی	1569	14
		اوہنےـ	1553	15
		جدوں	1547	16
		اسیں	1517	17
		ابہناں	1434	18
		میرےـ	1393	19
		توں	1340	20
		اج	1203	21
		اوہدےـ	1188	22
		کیہ	1176	23
		جبوں	1151	24
		ہر	1146	25
		نہیں	1142	26
		نـ	1113	27
		اوہنوں	993	28
		میرا	957	29
		اک	937	30
		بـ	886	31

		سالنے	863	32
		میری	855	33
		ایتھے	814	34
		اپدے	802	35
		کجھ	781	36
		بن	764	37
		کیوں	752	38
		سو	708	39
		اوہدی	705	40
		گل	687	41
		ٹھیں	681	42
		جد	669	43
		بس	669	44
		مان	657	45
		تار	648	46
		آپنے	646	47
		ایسے	644	48
		جو	642	49
		ابو	625	50
		گھر	617	51
		اس	606	52
		جبڑا	597	53
		فیر	576	54
		اویدا	576	55
		ہور	572	56
		کسے	562	57
		تیرے	557	58
		ہان	556	59
		سانوں	549	60
ہو	سکدا	اے	548	61
		پنجابی	544	62
		سِمن	539	63
		اسان	534	64
		دو	529	65
		مُڑ	522	66

		وج	511	67
		إنج	501	68
		سارے	500	69
		لوك	497	70
	اوس	آکھیا	495	71
		اوٹھے	489	72
		سبه	485	73
		آپ	481	74
		على	478	75
		ٿوں	474	76
		ڈاڪٹر	466	77
اپنے		آپ	465	78
		تینوں	462	79
ہو	جاندا	اے	446	80
		جس	435	81
ایس	کر	کے	422	82
		پاکستان	414	83
		دے	413	84
		جس	411	85
		نالے	410	86
		جبھڑے	401	87
		آپنی	401	88
		رب	398	89
		پنجاب	394	90
	ایس	لءی	393	91
		اتے	389	92
		پته	387	93
		کءی	385	94
		کجه	382	95
		نال	375	96
		بندے	375	97
		وج	371	98
	اوہناں	دے	370	99
		شاه	366	100

© Advanced Centre for Technical Development of Punjabi Language,
Literature and Culture, Punjabi University Patiala, Punjab India.

PIN- 147 002

Phone: +91-0175-3046171, 3046172

Fax: +91-0175-3046313

email: gslehal@yahoo.com