

K-NN SINIFLANDIRMA ALGORİTMASININ PARAMETRELERİNİN ZEKİ OPTİMİZASYON TEKNİKLERİYLE İYİLEŞTİRİLMESİ

Erdem Korhan AKÇAY*Eskişehir Teknik Üniversitesi

Özet

Makine öğrenmesi alanında verilerin sınıflandırılması algoritmalarının içinden sıkça kullanılan K-NN algoritmasının daha da güçlendirilmesi ve daha doğru sonuçları elde etmek için parametre optimizasyonun sağlanması önemli bir adımdır. Parametre optimizasyonunun en etkin şekilde kullanılabilmesi için zeki optimizasyon tekniklerinden olan Genetik Algoritması, Parçacık Sürü Algoritmaları ile K-NN algoritmasıyla birlikte çalıştırılmasıyla güçlü sonuçlar elde edilecektir. Genetik algoritma ve parçacık sürü optimizasyonu, K-NN algoritmasının kullanıcı tarafından belirlenen parametrelerini otomatik olarak ayarlayarak, daha iyi sınıflandırma sonuçları elde etmeyi hedefler.

1. Giriş

Sınıflandırma problemleri, makine öğrenimi alanında büyük bir öneme sahiptir. Bu problemlerde, veri noktaları farklı sınıflara atanırken, karar verme sürecinde kullanılan algoritmalara ihtiyaç duyulur. K-NN (k-en yakın komşu) sınıflandırma algoritması, basit ancak etkili bir yöntem olarak bilinir. Temel fikri, bir veri noktasını sınıflandırmak için çevresindeki k en yakın komşusuna bakmaktır. Ancak, K-NN algoritması birkaç parametreye sahiptir ve bu parametrelerin optimal değerlerini belirlemek her zaman kolay değildir.

Bu bağlamda, zeki optimizasyon teknikleri, K-NN sınıflandırma algoritmasının parametrelerinin iyileştirilmesinde önemli bir rol oynamaktadır. Zeki optimizasyon teknikleri olarak, genetik algoritmalar, parçacık sürü optimizasyonu, bulanık mantık gibi doğal veya yapay zeka temelli yöntemleri içerir. Bu teknikler, karmaşık ve çok boyutlu parametre arama uzaylarında optimal çözümleri bulmak için kullanılır.

K-NN sınıflandırma algoritmasının temel parametreleri arasında komşuluk sayısı (k), veri noktalarının benzerlik ölçüsü ve

ağırlıklandırma faktörü yer alır. Komşuluk sayısı, sınıflandırma doğruluğunu etkileyen kritik bir faktördür. Veri noktaları arasındaki benzerlik ölçüsü, K-NN'nin sınıflandırma kararlarını doğrudan etkiler. Ağırlıklandırma faktörü ise her bir komşunun sınıflandırmaya olan katkısını belirler.

Bu makalenin amacı, K-NN sınıflandırma algoritmasının parametrelerini iyileştirmek için zeki optimizasyon tekniklerinin nasıl kullanılabileceğini incelemektir. Zeki optimizasyon teknikleri, parametre arama uzayında etkin bir şekilde gezinerek en iyi parametre değerlerini bulmayı hedefler. Bu sayede, K-NN algoritmasının sınıflandırma performansı artırılabilir ve daha doğru sonuçlar elde edilebilir.

Bu çalışmanın bulguları, literatürdeki bazı önemli çalışmaların sonuçlarına da atıfta bulunarak sunulmuştur. Örneğin, "Classification of Heart Disease Using K-Nearest Neighbor and Genetic Algorithm" adlı çalışmada, jenerasyon sayısı 20 olarak alındığında, göğüs kanseri veri setinde %94 doğruluk oranı elde edilmiştir. Ayrıca, "Particle Swarm Optimization Feature Selection for Breast Cancer Prediction" adlı başka bir çalışmada ise göğüs kanseri veri setinde %97.54 ve %96.66 doğruluk oranlarına ulaşılmıştır.

detaylı analizler ve sonuçların daha geniş bir değerlendirmesi sunulacaktır.

2. Veri Seti

Uygulamanın yapıldığı veri seti Python programlama dilinin makine öğrenmesi uygulamaları amacıyla oluşturulmuş paketi “**sklearn**”in içerisinde barındıran “**Göğüs kanseri**” veri seti kullanılmıştır. Veri setini incelediğimizde **Şekil 1**’de genel bir bakışını görmekteyiz. İncelediğimizde 31 değişkenimiz ve her bir değişkenin 568 değeri bulunmaktadır.

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100	101	102	103	104	105	106	107	108	109	110	111	112	113	114	115	116	117	118	119	120	121	122	123	124	125	126	127	128	129	130	131	132	133	134	135	136	137	138	139	140	141	142	143	144	145	146	147	148	149	150	151	152	153	154	155	156	157	158	159	160	161	162	163	164	165	166	167	168	169	170	171	172	173	174	175	176	177	178	179	180	181	182	183	184	185	186	187	188	189	190	191	192	193	194	195	196	197	198	199	200	201	202	203	204	205	206	207	208	209	210	211	212	213	214	215	216	217	218	219	220	221	222	223	224	225	226	227	228	229	230	231	232	233	234	235	236	237	238	239	240	241	242	243	244	245	246	247	248	249	250	251	252	253	254	255	256	257	258	259	260	261	262	263	264	265	266	267	268	269	270	271	272	273	274	275	276	277	278	279	280	281	282	283	284	285	286	287	288	289	290	291	292	293	294	295	296	297	298	299	300	301	302	303	304	305	306	307	308	309	310	311	312	313	314	315	316	317	318	319	320	321	322	323	324	325	326	327	328	329	330	331	332	333	334	335	336	337	338	339	340	341	342	343	344	345	346	347	348	349	350	351	352	353	354	355	356	357	358	359	360	361	362	363	364	365	366	367	368	369	370	371	372	373	374	375	376	377	378	379	380	381	382	383	384	385	386	387	388	389	390	391	392	393	394	395	396	397	398	399	400	401	402	403	404	405	406	407	408	409	410	411	412	413	414	415	416	417	418	419	420	421	422	423	424	425	426	427	428	429	430	431	432	433	434	435	436	437	438	439	440	441	442	443	444	445	446	447	448	449	450	451	452	453	454	455	456	457	458	459	460	461	462	463	464	465	466	467	468	469	470	471	472	473	474	475	476	477	478	479	480	481	482	483	484	485	486	487	488	489	490	491	492	493	494	495	496	497	498	499	500	501	502	503	504	505	506	507	508	509	510	511	512	513	514	515	516	517	518	519	520	521	522	523	524	525	526	527	528	529	530	531	532	533	534	535	536	537	538	539	540	541	542	543	544	545	546	547	548	549	550	551	552	553	554	555	556	557	558	559	560	561	562	563	564	565	566	567	568	569	570	571	572	573	574	575	576	577	578	579	580	581	582	583	584	585	586	587	588	589	590	591	592	593	594	595	596	597	598	599	600	601	602	603	604	605	606	607	608	609	610	611	612	613	614	615	616	617	618	619	620	621	622	623	624	625	626	627	628	629	630	631	632	633	634	635	636	637	638	639	640	641	642	643	644	645	646	647	648	649	650	651	652	653	654	655	656	657	658	659	660	661	662	663	664	665	666	667	668	669	670	671	672	673	674	675	676	677	678	679	680	681	682	683	684	685	686	687	688	689	690	691	692	693	694	695	696	697	698	699	700	701	702	703	704	705	706	707	708	709	710	711	712	713	714	715	716	717	718	719	720	721	722	723	724	725	726	727	728	729	730	731	732	733	734	735	736	737	738	739	740	741	742	743	744	745	746	747	748	749	750	751	752	753	754	755	756	757	758	759	760	761	762	763	764	765	766	767	768	769	770	771	772	773	774	775	776	777	778	779	780	781	782	783	784	785	786	787	788	789	790	791	792	793	794	795	796	797	798	799	800	801	802	803	804	805	806	807	808	809	810	811	812	813	814	815	816	817	818	819	820	821	822	823	824	825	826	827	828	829	830	831	832	833	834	835	836	837	838	839	840	841	842	843	844	845	846	847	848	849	850	851	852	853	854	855	856	857	858	859	860	861	862	863	864	865	866	867	868	869	870	871	872	873	874	875	876	877	878	879	880	881	882	883	884	885	886	887	888	889	890	891	892	893	894	895	896	897	898	899	900	901	902	903	904	905	906	907	908	909	910	911	912	913	914	915	916	917	918	919	920	921	922	923	924	925	926	927	928	929	930	931	932	933	934	935	936	937	938	939	940	941	942	943	944	945	946	947	948	949	950	951	952	953	954	955	956	957	958	959	960	961	962	963	964	965	966	967	968	969	970	971	972	973	974	975	976	977	978	979	980	981	982	983	984	985	986	987	988	989	990	991	992	993	994	995	996	997	998	999	1000
1	1	10	10	10	10	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																			

Şekil 1

Değişkenleri incelediğinde:

radius: Tümörün ortalama yarıçapı,

texture: Tümörün ortalama dokusu.

perimeter: Tümörün ortalama çevresi,

area: Tümörün ortalama alanı,

smoothness: Tümörün ortalama düzgünlüğü

compactness: Tümörün ortalama sıklığı.

concavity: Tümörün ortalama çöküklüğü.

concave points: Tümördeki ortalama çukur nokta sayısı.

symmetry: Tümörün ortalama simetrisi.

fractal dimension: Tümörün ortalama fraktal boyutu.

Göğüs kanseri veri setinde hedef değer, iyi huylu veya kötü huylu olup olmadığını belirten 0 veya 1'dir

3. Yöntem

3.1. K – En Yakın Komşu Algoritması

K-NN algoritması, en temel örnek tabanlı öğrenme algoritmaları arasındadır. Örnek tabanlı öğrenme algoritmalarında, öğrenme işlemi eğitim setinde tutulan verilere dayalı olarak gerçekleştirilmektedir. Yeni karşılaşılan bir örnek, eğitim setinde yer alan örnekler ile arasındaki benzerliğe göre sınıflandırılmaktadır. K-NN

sınıflandırılmaktadır. K-NN algoritmasında, eğitim setinde yer alan örnekler n boyutlu sayısal nitelikler ile belirtilir. Her örnek n boyutlu uzayda bir noktayı temsil edecek biçimde tüm eğitim örnekleri n boyutlu bir örnek uzayında tutulur. Bilinmeyen bir örnek ile karşılaşıldığında, eğitim setinden ilgili örneğe en yakın k tane örnek belirlenerek yeni örneğin sınıf etiketi, k en yakın komşusunun sınıf etiketlerinin çoğunluk oylamasına göre atanır.

3.1.2. K-NN Algoritmasının Genel İşleyişi

K-En Yakın Komşular algoritmasının performansı için kritik öneme sahip noktalardan birisi örnekler arası yakınlığın nasıl ölçümleneceğidir. Yakınlık, Manhattan, Minkowski, Öklid uzaklığı gibi uzaklık ölçütleri kullanılarak hesaplanır. K-NN algoritması az sayıda parametre gerektirmesi ve basit bir yapıya sahip olması önemli noktalarındır. K-NN algoritmasında, ortaklık yapısının çoğunluk oylamasına dayalı olarak belirlenmesi, simetrik olmayan dağılıma sahip veri setlerinde sıklıkla görülen sınıfların, yeni örneklerin sınıf etiketlerinin belirlenmesinde daha baskın bir role sahip olmalarına neden olmaktadır. Bu nedenle, temel K-NN algoritmasının uzaklık

ölçütünün etki değerine farklı şekillerde ağırlık değeri atayan yöntemler bulunmaktadır.

K-NN algoritması, büyük eğitim setlerinin varlığında son derece etkili sonuçlar üretebilir. K-NN algoritması, ilgisiz değerlerin varlığında dahil sınıflandırma modeli oluşturabilir. Bu durumlarda eğitim süresi önemli ölçüde artmaktadır. K-NN algoritmasının basit yapısına rağmen, hesaplaması yüksek maliyetlere sahip olmaktadır. Özellikle büyük eğitim veri setleri için, sınıf etiketini belirlemek istenen örnek ile veri setindeki örnekler arasındaki farklılıkları belirlemek oldukça maliyetli olabilmektedir. K-NN algoritması, temel bileşenler analizi gibi boyut azaltma teknikleri veya arama ağaçları gibi daha güçlü veri yapıları kullanılarak bu maliyeti azaltabilir. Bunun yanı sıra, K-NN algoritması, komşu sayısı ve uzaklık ölçütü gibi parametrelere karşı hassas bir yapıya sahiptir ve çok boyutlu veri setlerinde çalışması etkili olmaz.

3.1.3. K-NN Parametreleri

Uzaklık ölçütü, komşu sayısı (k) ve ağırlıklandırma yöntemi, K-NN algoritmasının performansında önemli parametrelerdir. Bu parametreler alt bölümlerde açıklanmaktadır.

3.1.3.1 Uzaklık Ölçütleri

K-En Yakın Komşular algoritmasının kullandığı uzaklık ölçütleri olarak başlıca, Manhattan, Minkowski, Öklid uzaklıkları kullanılmaktadır.

3.1.3.1.1 Manhattan Uzaklığı

Manhattan uzaklığı, boyutu n olan iki nokta arasındaki uzaklığın mutlak değerlerince toplamıdır. Diğer bir deyişle, A ve B noktaları arasındaki Manhattan uzaklığı'nı hesaplamak istersek, $A = (x_1, \dots, x_n)$ ve $B = (y_1, \dots, y_n)$ olmak üzere,

$$\left[\sum_{i=1}^n |x_i - y_i| \right]$$

Formülüyle hesaplanır.

3.1.3.1.2 Minkowski Uzaklığı

Öklid uzayında tanımlı bir dizi olan Minkowski uzaklığı. Sınıflandırma, makine öğrenmesi, veri madenciliği alanlarında yoğun bir şekilde kullanılan Öklid uzaklığı, Uzaklık ölçütlerinin genelleştirilmiş halidir. Diğer bir deyişle, A ve B noktaları arasındaki Minkowski uzaklığı'nı hesaplamak istersek, $A = (x_1, x_2, \dots, x_n)$ ve $B = (y_1, y_2, \dots, y_n)$ olmak üzere,

$$\left[\sum_{i=1}^n |x_i - y_i|^p \right]^{1/p}$$

Minkowski uzaklığı, genel bir formül olarak tanımlanmakta olup, p'nin çeşitli değerleri için farklı uzaklık ölçütlerini tanımlar. Minkowski ölçütünde p=2 durumunda Öklid uzaklığı elde edilir, p=1 durumunda Manhattan uzaklığı, aynı zamanda $n \rightarrow \infty$ durumunda Chebyshev uzaklığı elde edilir.

3.1.3.1.3 Öklid Uzaklığı

Öklid uzaklığı, sınıflandırma algoritmalarında yaygın olarak kullanılan uzaklık ölçütüdür. Öklid uzaklığı, iki nokta arasındaki doğrusal uzaklığı belirtir. Diğer bir deyişle A ve B noktaları arasındaki Öklid uzaklığı'nı hesaplamak istersek,

$$A = (x_1, x_2, \dots, x_n) \quad \text{ve} \quad B = (y_1, y_2, \dots, y_n) \quad \text{olmak üzere,}$$

$$\left[\sum_{i=1}^n (x_i - y_i)^2 \right]^{1/2}$$

Öklid uzaklığı, K-ortalama kümeleme algoritması, temel K-NN algoritması gibi sınıflandırma ve kümeleme algoritmalarında yakınlığın ölçülmesi için kullanılan temel uzaklık ölçütüdür.

2.1.3.2 Komşu Sayısı (k)

K-NN algoritmasında, komşu sayısı (k) parametresinin değeri sınıflandırmayı belirler. Sınıflandırma sürecinde, $k=1$ için, sadece en yakın komşunun bulunduğu sınıfa seçilir. Örneğin, k sayısı örnek sayısına (N) yaklaştıkça veri setindeki tüm veriler gözden geçirilir ve oylama yoluyla seçilir.

3.1.3.3 Ağırlıklandırma

Ağırlık değerleri sınıflandırılmak istenen noktaya daha yakın olan komşu noktaların, çoğunluk oylamasına daha fazla katkı koyması amaçlanır. Kullanılan ağırlık değeri atama yöntemleri, her bir komşunun ağırlığının, d , komşuların arasındaki uzaklık olmak üzere, $1/d$ ya da $1/d^2$ şeklinde alınmasıdır.

3.2 Genetik Algoritma

Genetik algoritma (GA), çevrelerindeki değişikliklere uyum sağlayabilen türlerin hayatta kalmayı ve üremeyi başarabildiği doğal seçim sürecini simüle eden bir algoritmadır. Basitçe söylemek gerekirse, bir problemin çözümü için ardışık nesillerin bireyleri arasındaki "en uygun olanın hayatta kalması" durumunu simüle ederler. Her nesil, bireylerin bir popülasyonunu içerir. Her birey arama uzayında bir noktayı ve olası bir çözümü temsil eder. Her birey, bir dize olarak temsil edilir. Bu dize, Kromozom ile benzerlik gösterir.

GA, popülasyon adı verilen bir çözüm setiyle başlatılmaktadır. Her bir jenerasyon boyunca popülasyon büyüklüğü korunmaktadır. Her bir jenerasyonda, her bir kromozomun uygunluk değeri değerlendirilmekte ve sonra uygunluk değerlerine göre bir sonraki jenerasyon için kromozomlar olasılıksal biçimde

seçilmektedir. Seçilen bazı kromozomlar rassal olarak eşleşmekte ve çaprazlama ve mutasyon rassal olarak gerçekleştirilerek genç bireyler üretilmektedir. Yüksek uygunluk değerine sahip kromozomlar yüksek olasılıkla seçildiğinden, yeni jenerasyona ait kromozomlar eski jenerasyona ait kromozomlardan daha iyi uygunluk değerine sahip olabilmektedir. Bu evrim süreci sonlandırma şartı sağlanana kadar tekrar etmektedir. GA'daki çözümler diziler veya kromozomlar olarak adlandırılmaktadır. Çoğu durumda, kromozomlar listeler veya diziler olarak gösterilmektedir. Dolayısıyla, GA ile yapılan pek çok işlem liste veya dizilerle yapılan işlemlerdir (Kumar vd., 2010).

GA, ilerleyiş olarak sırasıyla başlangıç popülasyonunun oluşturulması, uygunluk değerlerinin hesaplanması, seçim, çaprazlama, mutasyon, yeni nesil üretimi ve en iyi değerin döndürülmesi biçiminde yedi aşamada incelenebilmektedir.

i. Başlangıç Popülasyonu Oluşturulması:

Çözümler gösterimi oluşturulup, amaç tanımlandıktan sonra, değerlendirme sürecine ilerlenmektedir. İlk olarak, rassal olarak veya önceden sahip olunan bilgiyle kodlanmış çözümlerin bir başlangıç popülasyonu veya kromozomu yaratılmaktadır (Zhou, 2006)

ii. Uygunluk Değerlerinin Hesaplanması:

Başlangıç çözümlerinin oluşturulmasının ardından popülasyonun iteratif bir değerlendirme ve yeniden üretim sürecine girmesi sağlanmaktadır. Değerlendirme aşamasında amaç fonksiyonu doğrultusunda her bir kromozoma bir uygunluk değeri atanmaktadır. (Zhou, 2006)

iii. Seçim Operatörü:

Uygunluk değerleri hesaplandıktan sonra seçim aşaması ile devam edilmektedir. Her

bir jenerasyon sırasında var olan popülasyonun bir grubu yeni jenerasyonlar oluşturmak için korunmaktadır. Kromozomlar, uygunluk değeri daha iyi olan bireylerin daha büyük olasılıkla seçildiği uygunluk tabanlı bir süreçle seçilmektedir. Belirli seçim metotları ile her bir çözümün uygunluğu değerlendirilmekte ve tercihe bağlı olarak en iyi çözümler seçilmektedir. Diğer metotlarla sadece popülasyonun rassal bir örnekleme değerlendirilmekte, ancak bu süreç epey zaman alabilmektedir. Seçim metotlarına rulet çemberi seçimi, stokastik evrensel seçim, sıralama seçimi ve turnuva seçimi örnek gösterilebilmektedir. (Sastri, Goldberg ve Kendall, 2014; Zhou, 2006)

iv. Çaprazlama Operatörü:

Seçim aşamasından sonra belirlenen bir çaprazlama operatörü ve bir çaprazlama oranıyla ebeveynler üzerinde yeni bireylerin çaprazlaması gerçekleşmektedir. Çaprazlama operatörü, her bir ebeveynden seçilen genleri kopyalayarak iki yeni genç birey üretmektedir. Her bir genç bireydeki belirli bir pozisyonda bulunan gen, iki ebeveyninden birindeki aynı pozisyonundaki genden kopyalanmaktadır. Belirlenen pozisyonlardaki genler için hangi ebeveynden genlerin kopyalanacağına çaprazlama maskesi adı verilen ek bir dizi ile karar verilmektedir. Tek noktalı çaprazlama, iki noktalı çaprazlama ve uniform çaprazlama çaprazlama operatörlerine örnek olarak gösterilmiştir. Çaprazlama ile yeni kromozomlar oluşturulmaktadır. Böylece uygunluk değeri daha yüksek kromozomların ortaya çıkma olasılığı sağlanmaktadır. (Haldurai, Madhubala ve Rajalakshmi, 2016)

v. Mutasyon Operatörü:

Mutasyon operatörü genellikle çaprazlama operatöründen sonra devreye girmektedir. Mutasyon operatörüyle, tek bir rassal gen

seçilerek birey alelleri rassal kromozomların lokuslarında çok küçük bir olasılıkla değiştirilmektedir. Bir sonraki seçimde gelişecek veya yok olacak uygunluk değeri daha yüksek veya daha düşük bir kromozom yaratılabilmektedir. Amaç en iyiyi bulmak olduğundan eğer popülasyonda bir tek kötü çözüm varsa olumsuzluk yaratabilmektedir. Diğer bir taraftan mutasyon sürecinde iyi bir çözüm yaratılırsa en iyi çözüme ulaşmada çok yardımcı olmaktadır. Mutasyon operatörü için nokta mutasyonu, güç mutasyonu, küçültme mutasyonu, denetleme mutasyonu, eşsizlik mutasyonu, değişken-olasılık mutasyonu, iki noktalı mutasyon, gauss mutasyonu, değiş tokuş mutasyonu ve uniform mutasyon, non-uniform mutasyon örnek verilmiştir. Pek çok evrimsel algoritma yerel optimum bir noktada sıkışma eğilimi gösterebildiği için nihai sonuç evrensel optimum noktadan çok daha uzakta sonlandırılabilir. Bunun önüne geçilmesi ve algoritmanın yerel optimumlardan kurtarılabilmesi için mutasyon operatörü ayrıca önem kazanmaktadır. (Chauhan, Singh ve Aggarwal, 2021)

vi. Sonlandırma:

Uygunluk değerinin hesaplanması, seçim, çaprazlama, mutasyon adımları ve bir GA stratejisi döngüsü sonlandırma şartı sağlanana kadar uygulanmaktadır. onlandırma şartları, en küçük uygunluk değeri şartını sağlayan bir çözümün bulunması, belirli sayıda jenerasyona ulaşılması, ayrılan bütçeye ulaşılması (hesaplama zamanı/para), en yüksek sıralamadaki çözümün uygunluk değerine ulaşılması veya daha fazla iyileşme olmayacak kadar başarılı bir çözüme ulaşılması olarak düşünülebilmektedir. (Kumar vd., 2010)

vi.i. Yer Değiştirme (Replacement):

Her bir jenerasyonda genetik operatörleri gerçekleştirmek için bireyler seçilmektedir. Yaratılan her bir yeni birey popülasyondaki uygunluk değeri en kötü bireyle karşılaştırılmaktadır. Eğer yeni birey popülasyondaki uygunluk değeri en düşük bireyden daha iyi bir uygunluk değerine sahipse o bireyin yerine seçilmektedir.(Wu ve Ji, 2007)

vi.ii. Elitizm:

Bu süreçte, sonraki jenerasyonda yeniden üretim için en yüksek uygunluğa sahip bireyler korunmaktadır. Elitizm ile mevcut popülasyondan elde edilen en iyi çözümün sonraki popülasyon için değişmemiş biçimde kopyalanması ve GA tarafından elde edilen çözüm kalitesinin bir jenerasyondan diğerine geçerken her bir iterasyondan sonra artması garanti edilmektedir (Markandeswar vd., 2016)

3.3 Parçacık Sürü Optimizasyonu

Parçacık sürü optimizasyonu, balık ve kuş gibi canlıların doğadaki sosyal davranışlarından esinlenerek geliştirilen meta-sezgisel bir optimizasyon algoritmasıdır (Eberhart ve Kennedy, 1995). Parçacık sürü optimizasyonunun köklerini, evrimsel hesaplama ve yapay yaşam örüntüleri olmak üzere iki ana bileşene dayanır. Bu nedenle, parçacık sürü optimizasyonu genetik algoritmalarla ve evrimsel programlama ile doğrudan bir ilişkisi vardır.

Bir kuş sürüsü birbirine çarpmadan birlikte uçar, yön değiştirir ve bunları eşzamanlı olarak gerçekleştirir. Sürüdeki bir kuş yiyecek bulduğunda, diğer kuşlarla iletişim kurarak onların da yiyeceğe ulaşmasını sağlar. Aynı şekilde, balık sürülerinde de bir balık bir tehlike sezdiğinde sürüdeki diğer balıklarla iletişime geçerek onların tehlikeden haberdar olmasını sağlar. Özellikle kuş ve balık sürülerinin iş birliği

davranışlarından ilham alan parçacık sürü optimizasyonu, yapay parçacıkların arama uzayında kuş ve balık sürülerinin hız ve konum gibi fiziksel özelliklerini taklit etmesiyle optimal çözümü bulmayı amaçlar. Algoritmada, her parçacık muhtemel bir çözümü temsil etmektedir (Eberhart ve Kennedy, 1995; Kaewkamnerdpong ve Bentley, 2005; Krause vd., 2013; Vasuki, 2020).

Bir sürü modeli geliştirmek için belirli prensiplerin sağlanması gerekmektedir. Bunlar,

- Çoklu birimler arasında bilgi paylaşımı,
- Birimlerin kendi kendini organize etmesi ve evrim geçirmesi,
- Birimlerin eş öğrenmesi açısından etkili olması,
- Uygulamalı ve gerçek zamanlı problemlere kolayca uyarlanabilmesi

Parçacıkların arama uzayındaki konumlarının değişimi bireylerin sosyo-bilişsel eğilimine bağlıdır ve parçacıkların mevcut konumlarının güncellenmesi için hız vektörü kullanılır. Bir parçacığın en iyi konumu “pbest” olarak adlandırılır ve bu değer hız vektörünün güncellenmesi için hafızada tutulur. Bununla birlikte, hız vektörünün hesaplanmasında popülasyondaki herhangi bir parçacık tarafından elde edilen en iyi konum da kullanılmaktadır. Bu konum ise “gbest” olarak adlandırılmaktadır. “pbest” ve “gbest” değerleri bulunduktan sonra parçacıkların konumu ve hızı güncellenir (Eberhart ve Kennedy, 1995). Bir parçacığın hızı,

$$v_i^{t+1} = wv_i^t + c_1r_1(pbest_i - x_i^t) + c_2r_2(gbest - x_i^t)$$

biçiminde hesaplanır. Burada, v_i^{t+1} , i . parçacığın $(t + 1)$. iterasyondaki hızını; v_i^t , i . parçacığın t . iterasyondaki hızını ve x_i^t , i . parçacığın t . iterasyondaki konumunu ifade etmektedir. Ayrıca, w , atalet ağırlığıdır ve genellikle $[0,1.2]$ aralığından seçilir. Öğrenme katsayıları olarak adlandırılan c_1 ve c_2 genellikle $[0, 2]$ aralığında belirlenir.

c_1 , parçacığın kendi deneyimlerine, c_2 ise parçacığın sürüdeki diğer parçacıkların deneyimlerine göre hareket etmesini sağlar. r_1 ve r_2 ise $[0, 1]$ aralığında düzgün dağılımdan rastgele üretilen değerlerdir.

Parçacıkların konumu ise,

$$x_i^{t+1} = x_i^t + v_i^{t+1}$$

biçiminde güncellenir. Burada, x_i^{t+1} , i . parçacığın $(t + 1)$. iterasyondaki konumunu ifade etmektedir (Eberhart ve Kennedy, 1995; Shi ve Eberhart, 1998; Kashani vd., 2021).

4. Bulgular

Bu çalışmanın amacı, K-NN sınıflandırma algoritmasının parametrelerini zeki optimizasyon teknikleriyle iyileştirerek, göğüs kanseri teşhisinde daha yüksek doğruluk oranları elde etmektir. Bu amaç doğrultusunda, GridSearch, genetik algoritma ve parçacık sürü optimizasyonu gibi zeki optimizasyon teknikleri uygulanmış ve elde edilen sonuçlar analiz edilmiştir.

PSO'nun ilerle şekli farklı olduğu için eklenmemiştir.

	Uzaklık	Komşu Sayısı	Ağırlık
G.A.	Manhattan	13	Uniform
K-NN	Manhattan	7	Uniform

Tablo 1

Analizler sonucu K-NN uygulamasını geliştirmek adına makine öğrenmesinde mutlak doğruluğu yakalamak adına belirli parametreleri bulmak adına uygulanan GridSearch ile birlikte uyguladığımızda en iyi metrikler olarak bizlere {'distance': 'manhattan', 'n_neighbors': 13, 'weights': 'uniform'} Metriklerini önerdiğini görüyoruz bu metriklerde accuracy(doğruluk) değerimize baktığımızda 0.9340 olarak yakalamış oluyoruz.

Genetik Algoritma'da ise parametre optimizasyonunu GridSearch mantığıyla ilerleyerek en optimal sonucu bizlere veriyor. Parametre optimizasyonunda olası parametre değerlerinin verilmesinin yanı sıra genetik algoritmanın parametrelerini "jenerasyon = 20, popülasyon boyutu = 30, yavruların boyutu = 12" olarak girildi. "Tpot" kütüphanesini kullanarak uyguladığımızda 20. Jenerasyonda 0.95614035 Doğruluk oranını yakalamış olduk ve en iyi metrikler olarak {'distance': 'manhattan', 'n_neighbors': 7, 'weights': 'uniform'} belirlenmiştir.

Parçacık Sürü Optimizasyonu uygulamasında "pyswarm" kütüphanesi yardımıyla uygulandığında model doğruluk oranımız 0.9385964 olmuştur. Diğer modellerin aksine PSO'da en iyi metriğin hesaplanması yerine, sistemden çıkarılması gereken belirli özellikleri(feature) saptar ve o saptanan özelliklerin çıkartılması üzerine kurulmuş bir sistem olarak dizayn edilmiştir. Bu sisteme göre PSO'nun uygulama için belirlediği featurelar "mean texture, radius error, area error, smoothness error, compactness error, worst radius, worst texture, worst perimeter" olarak saptanmıştır.

	K-NN(G.S.)	Genetik A.	Parçacık Sürü O.
Doğruluk	0.9362	0.9561	0.9385

Tablo 2

Bulgularımız, K-NN sınıflandırma algoritmasının parametrelerini zeki optimizasyon teknikleriyle iyileştirmenin, sınıflandırma performansını artırılmasında önemli bir potansiyele sahip olduğunu göstermektedir. GridSearch, Genetik Algoritma ve Parçacık Sürü Optimizasyonu yöntemleriyle elde edilen sonuçlar, doğruluk oranlarında belirgin bir iyileşme sağladığını göstermektedir. Bu çalışma, K-NN sınıflandırma algoritmasının kullanıldığı diğer problemler için de benzer iyileştirmelerin elde edilebilir. Gelecekteki çalışmalarda, daha geniş veri setleri ve

farklı zeki optimizasyon tekniklerinin kullanılmasıyla elde edilen sonuçların daha da geliştirilmesi hedeflenebilir.

5. Sonuç

Bu çalışmada, K-NN sınıflandırma algoritmasının parametrelerini zeki optimizasyon teknikleriyle iyileştirme potansiyeli incelenmiştir. GridSearch, genetik algoritma ve parçacık sürü optimizasyonu yöntemleri kullanılarak algoritmanın performansı artırılmış ve doğruluk oranları elde edilmiştir.

Bulunan parametrelerle elde edilen K-NN sınıflandırma modelinde, GridSearch parametre optimizasyonunda göğüs kanseri veri setinde %93.62 doğruluk oranı elde edilmiştir. Genetik algoritma ile yapılan optimizasyonda ise jenerasyon sayısı 20 olarak belirlenmiş ve %95.61 doğruluk oranı elde edilmiştir. Parçacık sürü optimizasyonu kullanılarak elde edilen sonuçlarda ise %93.85 doğruluk oranı elde edilmiştir.

Sonuçlarımız, genetik algoritmanın K-NN sınıflandırma algoritmasının parametrelerini iyileştirmek için etkili bir yöntem olduğunu göstermektedir. Ayrıca, GridSearch ve parçacık sürü optimizasyonu yöntemlerinin de sınıflandırma alternatif olarak kullanılabileceği görülmüştür.

Bu bulgular, zeki optimizasyon tekniklerinin K-NN sınıflandırma algoritmasının performansını artırmada önemli bir rol oynayabileceğini göstermektedir. Bu çalışma, gelecekteki araştırmalara ve uygulamalara yönelik bir temel sağlamaktadır. Daha fazla veri seti ve parametre kombinasyonunun incelenmesi, algoritmanın performansının daha da iyileştirilmesine katkı sağlayabilir.

6. Referanslar

- Erdal Taşcı1, Ayтуğ Onan, “K-En Yakın Komşu Algoritması Parametrelerinin Sınıflandırma Performansı Üzerine Etkisinin İncelenmesi” *Akademik Bilişim Konferansı (AB 16)*, 30.01.2016 05.02.2016
- Anıl Yalçın, “Doktor Nöbet Çizelgeleme Problemi İçin Ağırlıklı Hedef Programlama Tabanlı Genetik Algoritma” *Kütahya Dumlupınar Üniversitesi Lisansüstü Eğitim Enstitüsü, Endüstri Mühendisliği Anabilim Dalı, Yüksek Lisans Tezi*
- Eda Özkul, “Yapay Çekirge Sürü Optimizasyonu” *Karadeniz Teknik Üniversitesi Fen Bilimleri Enstitüsü Trabzon, İstatistik ve Bilgisayar Bilimleri Anabilim Dalı, Doktora Tezi*
- Nurhayati, Fajar Agustian, Muhammad Dzil Ikram Lubis, “Particle Swarm Optimization Feature Selection for Breast Cancer Prediction” *2020 8th International Conference on Cyber and IT Service Management (CITSM)*
- M.Akhil jabbar, B.L Deekshatulu, Priti Chandra, “Classification of Heart Disease Using K- Nearest Neighbor and Genetic Algorithm” *Procedia Technology 10 (2013) 85 – 94*