

Identification of genes affecting growth traits in broiler chickens using linear regression and machine learning methods

ABSTRACT

Knowledge about the association between single nucleotide polymorphisms (SNPs) and important economic traits is one of the crucial tools in breeding programs within the poultry industry. Genome-wide studies for discovering SNP variations related to these traits are often conducted using simple linear models. However, due to certain assumptions of these models, some SNP markers may not be identified. This study aimed to evaluate the performance of random forest and gradient boosting methods compared to linear models in identifying SNP markers associated with body weight traits at 6 and 9 weeks of age in F2 broiler chickens resulting from crosses between the commercial Arian line and native Urmia birds. The results showed that the machine learning approaches were able to identify important markers, such as GGaluGA308573, GGaluGA255033, Gga_rs13614212, Gga_rs13743072, GGaluGA258772, Gga_rs14034395, and Gga_rs13858398, associated with body weight traits, which were related to genes *MAP2*, *ACSL1*, *CAMSAP2*, *FAM117B*, *SLC4A4*, *TIMP4*, and *LncRNA*, respectively. These genes are primarily involved in cellular division, growth control, regulation of cellular skeleton structure and microtubules, and transcription activity, constituting the most important biological processes. The identification of these novel genes using machine learning methods, which were not detected by linear models and previous studies in this population, could provide new insights into genetic control of growth traits in broiler chickens. Moreover, the discovered significant markers can be utilized in genetic improvement programs for birds.

Keywords: nucleotide polymorphisms, Genome-wide studies, broiler chickens, machine learning.



۱
۲
۳
۴
۵
۶
۷
۸
۹
۱۰
۱۱
۱۲
۱۳
۱۴
۱۵
۱۶
۱۷
۱۸
۱۹
۲۰
۲۱
۲۲
۲۳
۲۴
۲۵
۲۶
۲۷
۲۸
۲۹
۳۰
۳۱
۳۲
۳۳
۳۴
۳۵
۳۶
۳۷
۳۸
۳۹
۴۰
۴۱
۴۲
۴۳
۴۴

۸۱ الگوریتم‌های یادگیری داده محور هستند که می‌توانند با استفاده از داده‌های موجود، آموزش دیده و از الگوی کشف شده، در پیش‌بینی
 ۸۲ سناریوهای آینده استفاده کنند. یادگیری ماشین بسته به نوع آموزش، می‌تواند تحت نظارت یا بدون نظارت باشد. یادگیری نظارت
 ۸۳ شده از یک مجموعه‌ی آموزشی برای آموزش مدل‌ها و به دست آوردن خروجی مورد نظر استفاده می‌کند. این مجموعه‌ی آموزشی،
 ۸۴ شامل ورودی‌ها و خروجی‌های صحیح است که به مدل اجازه می‌دهد در طول زمان آموزش ببیند (Baştanlar and Özuysal, 2014).
 ۸۵ در میان روش‌های یادگیری نظارت شده می‌توان به روش جنگل تصادفی (Random Forest, RF) و گرادیان بوستینگ
 ۸۶ (Gradient Boosting, GB) اشاره نمود.

۸۷ مطالعات نشان داده است که روش‌های جنگل تصادفی و گرادیان بوستینگ در مطالعات پویش ژنومی، روش‌های مناسبی برای
 ۸۸ انتخاب نشان‌گرهای مهم و شناسایی ژن‌های تاثیرگذار بر صفات اقتصادی هستند. این روش‌ها بر خلاف مدل‌های آماری دیگر،
 ۸۹ نیازی به تعیین توزیع جمعیت مورد بررسی، ارزیابی برآوردگرها (فاصله اطمینان یا p-value) یا آزمون فرضیه‌های صفر ندارند و روش
 ۹۰ آزاد و بدون آمار هستند. علاوه بر این، این روش‌ها قادر هستند نشان‌گرهایی بدون تاثیر و یا دارای تاثیر منفی بر صفت را مشخص
 ۹۱ کنند و روابط غیرخطی بین نشان‌گرها و فنوتیپ را در نظر بگیرند و در داده‌هایی با ابعاد بزرگ به راحتی قابل استفاده باشند (Li et al., 2018; Sun et al., 2022).
 ۹۲ عواملی مانند عدم تعادل پیوستگی، اثرات غیرافزایشی، فراوانی آللی جزئی، روابط پیچیده و برهم‌کنش بین نشان‌گرها بر عملکرد
 ۹۳ روش‌های RF و GB موثر باشند (Li et al., 2018; Enoma et al., 2022). با این حال، مطالعات انجام شده برای بررسی الگوریتم-
 ۹۴ های مختلف یادگیری ماشین برای مطالعات پویش ژنومی در جوجه‌های گوشتی بسیار محدود است. بنابراین، هدف از مطالعه‌ی
 ۹۵ حاضر، شناسایی نشان‌گرهای مهم و تاثیرگذار بر صفات وزن بدن جوجه‌های نسل F₂ با استفاده از دو روش RF و GB و مقایسه آن-
 ۹۶ ها با روش مدل خطی (Linear Model, LM) است.
 ۹۷
 ۹۸
 ۹۹

روش شناسی پژوهش

۱۰۰ داده‌های فنوتیپی استفاده شده در پژوهش حاضر، از پرند‌های آمیخته نسل دوم (F₂) حاصل از تلاقی دو طرفه‌ی بین پرند‌های لاین
 ۱۰۱ تجاری آرین (A) و بومی استان آذربایجان غربی (U)، که در مرکز تحقیقات پرورش طیور دانشکده کشاورزی دانشگاه تربیت مدرس
 ۱۰۲ ایجاد شدند، بود. پرندگان F₁ از آمیزش $A \text{ ♂} \times U \text{ ♀}$ و $U \text{ ♂} \times A \text{ ♀}$ ایجاد شدند. به منظور ایجاد پرند‌های نسل دوم، نرهای F₁ از
 ۱۰۳ هر تلاقی متقابل یا معکوس، هر کدام با ۴ تا ۶ قطعه از ماده‌های خانواده‌های دیگر آمیزش داده شدند. تخم‌مرغ‌های جمع‌آوری شده،
 ۱۰۴ ابتدا شماره‌گذاری شده، سپس در دمای ۱۲ درجه سانتی‌گراد نگهداری شدند. پس از جداسازی تخم‌مرغ‌های کثیف، بدشکل و
 ۱۰۵ نامناسب، بقیه تخم‌مرغ‌ها جهت جوجه‌کشی در سینی‌های ستر قرار و سپس در روز ۱۸ به هچر انتقال داده شدند. جوجه‌ها پس از
 ۱۰۶ تولد تا ۷ روزگی به صورت گروهی نگهداری و به وسیله دان‌خوری سینی و آب‌خوری کله‌قندی تغذیه شدند و از ۸ روزگی تا پایان دوره
 ۱۰۷ به صورت انفرادی در قفس‌های تکی پرورش یافتند. این پرند‌ها با استفاده از سه جیره غذایی پلت شده شامل، آغازین، رشد و پایانی
 ۱۰۸ پرورش داده شدند. رکوردهای وزن بدن از تولد تا پایان دوره پرورش به صورت هفتگی ثبت گردید. برای مطالعه حاضر، از وزن بدن
 ۱۰۹ ثبت شده در سنین ۶ هفته (BW6) و ۹ هفته (BW9) استفاده شد. آماره‌های توصیفی صفات مورد مطالعه در جدول ۱ نشان داده شده
 ۱۱۰ است.

جدول ۱. آماره توصیفی صفات وزن بدن در ۶ هفته (BW6) و ۹ هفته (BW9)

صفت	تعداد	میانگین	انحراف معیار	حداقل	حداکثر
BW6	۳۰۰	۱۳۳۶	۲۵۷	۷۷۰	۲۰۵۰
BW9	۲۹۸	۲۲۹۴	۳۲۷	۱۷۵۵	۳۰۸۰

برای پرنده‌های تحقیق حاضر، نمونه‌های خون از جوجه‌های نسل دوم جمع‌آوری شد. استخراج DNA با استفاده از روش بهینه رسوب‌دهی نمکی انجام گرفت. جهت تعیین ژنوتیپ هر یک از افراد، نمونه‌ها به دانشگاه Aarhus کشور دانمارک انتقال داده شدند. در مجموع، تعداد ۳۱۲ نمونه DNA از نسل F₂ با همکاری این دانشگاه با استفاده از تراشه SNP تجاری ژنوم مرغ بنام Illumina 60K SNP chip (Illumina industry, USA)، که توسط شرکت تجاری کاب در اختیار قرار داده شده بود، تعیین ژنوتیپ شدند. برای هر نمونه ۵۴۳۴۰ نشان‌گر چندشکلی تک‌نوکلئوتیدی تعیین شد. این پرندگان نتاج حاصل از ۸ خانواده برادر – خواهر ناتنی بودند. مراحل کنترل کیفیت بر روی داده‌های اصلی با استفاده از نرم افزار PLINK 1.9 انجام شد (Purcell et al., 2007). نشان‌گرهای با نرخ فراخوانی کمتر از ۹۵ درصد یا حداقل فراوانی آلی (Minor allele frequency, MAF) کمتر از ۵ درصد، ژنوتیپ از دست رفته بیش از یک درصد و عدم تعادل هاردی واینبرگ در سطح احتمال کمتر از ۱۰^{-۶} حذف شدند. پس از فرآیند کنترل کیفیت، ۴۵۵۱۲ نشان‌گر و ۳۰۰ پرنده برای بررسی نهایی باقی ماندند. قبل از اجرای روش‌های یادگیری ماشین و رگرسیون خطی جهت شناسایی نشان‌گرهای مهم، داده‌های فنوتیپی مورد نظر برای اثرات جنس و نوبت جوجه‌کشی با روش تابعیت خطی تصحیح شدند.

روش‌های مورد استفاده جهت شناسایی نشان‌گرهای با اهمیت

مدل خطی

مدل خطی به صورت زیر است:

$$y = 1\mu + Zq + e$$

که در آن، y بردار مقادیر فنوتیپی تصحیح شده برای وزن بدن در سنین ۶ یا ۹ هفتگی، 1 بردار n بعدی از یک‌ها، μ میانگین جمعیت، q بردار نامعلوم اثر نشان‌گر، به عنوان تابعیت ثابت مشاهده از ژنوتیپ، Z بردار حاوی ژنوتیپ‌های نشان‌گر با مقادیر ۰ برای ژنوتیپ A_1A_1 ، ۱ برای ژنوتیپ A_1A_2 و ۲ برای ژنوتیپ A_2A_2 ، و e بردار مقادیر باقی‌مانده با فرض $e \sim N(0, I\sigma_e^2)$ است. فرض می‌شود که σ_e^2 واریانس باقیمانده و I ماتریس همبستگی است. آزمون ارتباط ژنتیکی با استفاده از دستور “--Linear” نرم افزار PLINK 1.9 انجام شد (Purcell et al., 2007). در این مدل، ۱۰ چندشکلی تک‌نوکلئوتیدی برتر براساس مقادیر p-value از نتایج مدل خطی انتخاب شدند. لازم به توضیح است که دو عامل ساختار جمعیت و مولفه‌های اصلی، به این دلیل که قابل برآزش در روش‌های یادگیری ماشین نبودند، در مدل فوق منظور نشدند.

جنگل تصادفی

روش جنگل تصادفی، یک الگوریتم ترکیبی از چندین درخت تصمیم است. در واقع، آن‌ها یک مجموعه‌ای از درختان تصمیم در یک جنگل ایجاد می‌کنند که می‌توانند پیش‌بینی بهتری نسبت به یک درخت فردی داشته باشند. در الگوریتم RF، برای رشد هر درخت از یک نمونه bootstrap (نمونه‌گیری با جایگزینی) از داده‌های آموزش اصلی استفاده می‌شود. داده‌های آموزشی شامل اطلاعات ژنوتیپی و فنوتیپ‌های تصحیح شده می‌باشد. برای تشکیل هر درخت براساس اطلاعات ژنوتیپی و فنوتیپی، یک نمونه با جایگزینی از داده‌های آموزشی گرفته و درخت روی نمونه گرفته شده رشد می‌کند و در هر گره از تعدادی از SNP ها که به صورت تصادفی انتخاب شده‌اند، یک SNP با توجه به قدرت تفکیک بهتر انتخاب می‌شود و گره براساس ژنوتیپ آن SNP منشعب و فنوتیپ‌ها بر اساس آن SNP به برگ یا گره بعدی تفکیک می‌شوند. در مراحل بعد و در گره‌های دیگر از نمونه‌های دیگر SNP، یک SNP انتخاب شده و انشعاب گره‌ها ادامه می‌یابد تا درخت کامل شود و همه‌ی فنوتیپ‌ها در برگ‌های درخت قرار بگیرند. به همین صورت درخت‌های بعدی شکل می‌گیرند. الگوریتم جنگل تصادفی، نتیجه را با میانگین‌گیری از خروجی‌های حاصل از همه درختان در جنگل پیش‌بینی می‌کند (Breiman, 2001). در هر بار bootstrap از اطلاعات، برخی اطلاعات هرگز نمونه‌گیری نمی‌شوند و برخی دیگر چند بار نمونه‌گیری می‌شوند. به عبارت دیگر، هر داده ورودی برای برخی درختان، داده‌های خارج از کیسه (out of Bag) خواهند بود.

۱۴۷ که در ایجاد برخی از درختان مشارکت ندارند. این داده‌ها به عنوان یک اعتبارسنج داخلی، که از طریق برآورد خطای خارج از کیسه
 ۱۴۸ (out of Bag error) انجام می‌شود، برای هر درخت عمل می‌کنند.

۱۴۹ برای محاسبه اهمیت هر چندشکلی تک‌نوکلئوتیدی، از روش جای‌گذاری تصادفی (Permutation) استفاده می‌شود. ابتدا خطای
 ۱۵۰ پیش‌بینی با استفاده از نمونه‌های خارج از کیسه در درخت مربوطه محاسبه می‌شود. سپس مقادیر ژنوتیپی هر چندشکلی
 ۱۵۱ تک‌نوکلئوتیدی به صورت تصادفی جای‌گذاری می‌شوند و خطای پیش‌بینی نمونه‌های خارج از کیسه مجدداً محاسبه می‌شود. تفاوت
 ۱۵۲ بین خطای نمونه‌های خارج از کیسه برای چندشکلی تک‌نوکلئوتیدی با مقادیر ژنوتیپی اصلی و مقادیر ژنوتیپی جای‌گذاری شده به
 ۱۵۳ شکل تصادفی (میانگین در تمام درختان در جنگل) نشان‌دهنده اهمیت یا قابلیت پیش‌بینی آن چندشکلی تک‌نوکلئوتیدی خاص است
 ۱۵۴ که به صورت درصد افزایش در میانگین مربعات خطا (percentage of increase in MSE) مشخص می‌شود (Nicodemus *et al.*,
 ۱۵۵ 2010; Li *et al.*, 2018). هر چقدر این مقدار بیشتر باشد نشان می‌دهد نبود چندشکلی تک‌نوکلئوتیدی مربوطه در نمونه، باعث
 ۱۵۶ افزایش خطای پیش‌بینی شده و بنابراین از اهمیت بیشتری برخوردار است.

۱۵۷ در مطالعه حاضر، ۳ پارامتر مهم روش RF برای تعیین مقدار مناسب آن‌ها در نظر گرفته شد که شامل تعداد درخت (Ntree) با
 ۱۵۸ مقادیر ۱۰۰۰ تا ۵۰۰۰ درخت (با فاصله ۱۰۰۰ درخت)، تعداد نمونه‌های انتخاب شده برای تشکیل هر گره (mtry) با مقادیر \sqrt{p} ،
 ۱۵۹ $p/3$ ، $0.5 \times (p/3)$ و $2 \times (p/3)$ ، که p نشان دهنده تعداد نشان‌گرها است و حداقل تعداد نشان‌گرها در گره‌های پایانی (node size) با
 ۱۶۰ مقادیر ۵، ۱۰ و ۱۵ بودند. از ریشه میانگین مربعات خطا (Root Mean Squared Error, RMSE) برای انتخاب بهترین ترکیب از
 ۱۶۱ پارامترها استفاده شد. برای اجرای روش RF از بسته نرم‌افزاری random Forest در نرم‌افزار R استفاده شد (Liaw and Wiener.,
 ۱۶۲ 2015).

گرادین بوستینگ

۱۶۵ توابع اصلی استفاده شده در الگوریتم گرادین بوستینگ (GB)، یادگیرنده‌های ضعیف مانند درخت‌های تصمیم (decision stump)
 ۱۶۶ هستند. در این روش سعی می‌شود تعدادی از یادگیرنده‌های ضعیف و کمکی تولید شود که به هم پیوسته بوده و با یک طرح آموزش
 ۱۶۷ مرحله‌ای برای کاهش خطا در پیش‌بینی کار کنند. هدف این الگوریتم، افزایش توان یک یادگیرنده قوی از مجموعه‌ای از
 ۱۶۸ یادگیرنده‌های ضعیف است. در این روش، یک یادگیرنده پایه به صورت ترتیبی به باقیمانده‌های درخت قبلی اضافه می‌شود و انتظار
 ۱۶۹ می‌رود با تمرکز بر روی داده‌های نادرست پیش‌بینی شده در درخت قبلی، نرخ خطا در درخت بعدی کاهش یابد و تا زمانی که نرخ
 ۱۷۰ خطا در حال کاهش باشد، الگوریتم بوستینگ ادامه خواهد داشت. در مطالعه حاضر، نشان‌گرهای مهم در روش GB با استفاده از تأثیر
 ۱۷۱ نسبی (relative influence) شناسایی می‌شوند که میانگین کاهش میانگین مربعات خطا (MSE) در همه درخت‌ها را محاسبه می‌کند
 ۱۷۲ زمانی که چندشکلی تک‌نوکلئوتیدی خاصی برای تقسیم داده‌ها استفاده می‌شود (Friedman, 2001).

۱۷۳ سه پارامتر مهم روش GB برای تعیین مقدار مناسب آن‌ها، شامل تعداد درخت (Ntree) با مقادیر ۱۰۰۰ تا ۵۰۰۰ درخت (با فاصله
 ۱۷۴ ۱۰۰۰ درخت)، حداقل تعداد نشان‌گرها در گره‌های پایانی (n.minobinnode) با مقادیر ۵، ۱۰ و ۱۵ و ضریب اختصاص داده شده به
 ۱۷۵ درختان (shrinkage) با مقادیر ۰/۱، ۰/۰۷، ۰/۰۴ و ۰/۰۱ بودند. از بسته نرم‌افزاری gbm در نرم‌افزار R برای اجرای روش GB
 ۱۷۶ استفاده شد (Ridgeway *et al.*, 2013). برای روش‌های RF و GB تنظیمات پارامترها از طریق جستجوی شبکه با استفاده از 3-fold
 ۱۷۷ cross-validation بر روی ۷۵ درصد نمونه‌های تصادفی انجام شد.

غربالگری گستره ژنوم برای نشان‌گرهای برتر

۱۸۰ تمام نشان‌گرها بر اساس اهمیت آن‌ها پس از اجرای مدل خطی و روش‌های یادگیری ماشین بر اساس بهترین ترکیب پارامترها
 ۱۸۱ رتبه‌بندی شدند. پس از رتبه‌بندی، ده نشان‌گر برتر از هر روش برای صفات وزن بدن مشخص شد. علاوه بر این، ژن‌های قرار گرفته

در فاصله ۱-Mb بالا و پایین از ۳ نشان گر برتر، که از هر روش در ناحیه ژنومی تشخیص داده شده بود، با استفاده از پایگاه داده‌های NCBI و Ensemble از ژنوم مرجع مرغ (Gallus-Gallus) مشخص شدند (Yates et al., 2016). جهت کاوش فرآیندهای بیولوژیکی و عملکردی ژن‌ها، هستی‌شناسی ژن با استفاده از نرم افزار آنلاین DAVID (<http://david.abcc.ncifcrf.gov/>) استفاده شد.

۱۸۵

۱۸۶

یافته های پژوهش و بحث

۱۸۷ نتایج ترکیب مختلف از پارامترهای مهم روش‌های یادگیری ماشین براساس ریشه میانگین مربعات خطا برای صفات وزن بدن در ۶ و
۱۸۸ ۹ هفتگی به ترتیب در شکل‌های ۱ و ۲ نشان داده شده است. در روش جنگل تصادفی، ترکیب تعداد نمونه‌های انتخاب شده برای
۱۸۹ تشکیل هر گره، تعداد درخت و حداقل تعداد نشان‌گرها در گره‌های پایانی به ترتیب برابر $p/3$ ، ۴۰۰۰ و ۵ برای صفت وزن بدن در ۶
۱۹۰ هفتگی و $0.5 \times p/3$ ، ۳۰۰۰ و ۵ برای صفت وزن بدن در ۹ هفتگی حداقل خطا را ایجاد کرد. در روش گرادیان بوستینگ، ترکیب تعداد
۱۹۱ درخت، حداقل تعداد نشان‌گرها در گره‌های پایانی و ضریب اختصاص داده شده به درختان به ترتیب برابر ۳۰۰۰، ۵ و ۰/۱ برای هر
۱۹۲ دو صفت وزن بدن انتخاب شدند. در مطالعات پویس ژنومی پیشنهاد شده است که مقدار پارامتر تعداد متغیر انتخاب شده در هر گره
۱۹۳ درخت (mtry) برای روش جنگل تصادفی باید از ۰/۱ تعداد SNPها بیشتر باشد (Goldstein et al., 2010) که در مطالعه حاضر برای
۱۹۴ این پارامتر، به ترتیب، ۱۵۱۷۱ و ۷۵۸۵ SNP برای صفت وزن ۶ و ۹ هفتگی در نظر گرفته شد. مقدار مناسب حداقل اندازه‌ی گره‌های
۱۹۵ پایانی در هر دو روش (node size در جنگل تصادفی و n.minobinnode در گرادیان بوستینگ) برای هر دو صفت برابر ۵ تعیین شد.
۱۹۶ هر اندازه مقدار این پارامتر بزرگتر باشد گره و برگ کم‌تری تشکیل خواهد شد و در نتیجه درختان کوچک‌تری تولید می‌شوند که
۱۹۷ ممکن است خیلی از نشان‌گرها در تشکیل درخت شرکت نکنند و در نتیجه کارایی مدل برای انتخاب نشان‌گرهای با اهمیت کاهش
۱۹۸ پیدا کند، بنابراین پیشنهاد شده است از مقادیر کوچک برای این پارامتر استفاده شود (Boulesteix et al., 2012). در هر دو روش با
۱۹۹ افزایش تعداد درخت، احتمال اینکه بیشتر نشان‌گرهای نمونه‌گیری و اهمیت آن‌ها سنجیده شود، افزایش می‌یابد. به عبارت دیگر، هر
۲۰۰ نشان‌گر فرصت آن را خواهد داشت که حداقل یک بار نمونه‌گیری شود. بنابراین، زمانی که تعداد نشان‌گرها زیاد است نباید از مقادیر
۲۰۱ کوچک یا پیش فرض برنامه برای تعداد درخت استفاده شود (Boulesteix et al., 2012). در روش گرادیان بوستینگ، ضریب ۰/۱۰
۲۰۲ اختصاص داده شده به درختان برای هر دو صفت، کمترین خطا را ایجاد کرد. این پارامتر، میزان مشارکت هر درخت در مدل را نشان
۲۰۳ می‌دهد و مقدار آن معمولاً با تعداد درخت رابطه عکس دارد، به طوری که با افزایش این پارامتر می‌توان تعداد درخت مورد نیاز جهت
۲۰۴ ایجاد مدل را کاهش داد. این پارامتر با مقداری کوچک‌تر از ۰/۱ و تعداد درخت زیاد جهت تحلیل دقیق‌تر در روش بوستینگ پیشنهاد
۲۰۵ شده است (Friedman, 2002).

۲۰۶ پلات منهن برای صفات وزن بدن در ۶ هفتگی و ۹ هفتگی با استفاده از روش‌های مدل خطی، جنگل تصادفی و گرادیان
۲۰۷ بوستینگ، به ترتیب در شکل‌های ۳ و ۴ نشان داده شده است. در روش جنگل تصادفی، اهمیت نشانگرها به صورت منفی، صفر و
۲۰۸ مثبت نشان داده می‌شود در حالی که در روش گرادیان بوستینگ اهمیت نشانگرها، صفر یا مثبت است (Li et al., 2018). برای
۲۰۹ صفت وزن بدن در ۶ و ۹ هفتگی به ترتیب ۴۷ درصد و ۴۸ درصد از نشانگرها در جنگل تصادفی و ۲۶ درصد و ۱۶ درصد از
۲۱۰ نشانگرها در روش گرادیان بوستینگ دارای اهمیت مثبت بودند. با استفاده از دو روش یادگیری ماشین در کروموزوم‌های ۴ و ۷ برای
۲۱۱ صفت وزن بدن در ۶ هفتگی و در کروموزوم‌های ۶ و ۸ برای صفت وزن بدن در ۹ هفتگی، نشانگرهای با اهمیتی کشف شدند.

۲۱۲ تعداد ده نشان گر برتر شناسایی شده از هر سه روش برای صفات وزن بدن در ۶ و ۹ هفتگی به ترتیب در جداول ۲ و ۳ نشان
۲۱۳ داده شده است. در مدل خطی، دامنه مقدار p-value برای ده نشان گر برتر از 10^{-5} تا $2/01 \times 10^{-4}$ و $1/53 \times 10^{-4}$ برای وزن ۶ هفتگی و از
۲۱۴ 10^{-5} تا $2/51 \times 10^{-4}$ برای وزن ۹ هفتگی بود. براساس این مقادیر، نشان‌گرهای GGAluGA141221 و
۲۱۵ GGAluGA142838 به ترتیب برای وزن ۶ و ۹ هفتگی مهم‌ترین نشان‌گرها بودند. در روش جنگل تصادفی، دامنه مقدار درصد
۲۱۶ افزایش خطای میانگین مربعات برای ده نشان گر برتر از ۰/۴۱ تا ۰/۹۶ برای وزن ۶ هفتگی و از ۰/۳۱ تا ۰/۶۶ برای وزن ۹ هفتگی

۲۱۷ بود. در این روش، برای وزن بدن در ۶ و ۹ هفتگی، مهم‌ترین نشان‌گرها، به ترتیب، Gga_rs15763229 و GGaluGA308573 بودند.
 ۲۱۸ دامنه مقدار تاثیر نسبی در روش گرادیان بوستینگ برای وزن ۶ هفتگی از ۵۲۶۲۳/۳۵ تا ۱۱۳۹۲۹/۴ و برای وزن ۹ هفتگی از
 ۲۱۹ ۴۰۴۷۹ تا ۱۶۳۳۵۲ بود. با در نظر گرفتن این مقادیر، دو نشان‌گر Gga_rs13743072 و Gga_rs13614212، به ترتیب، مهم‌ترین
 ۲۲۰ نشان‌گرهای کشف شده برای وزن‌های ۶ و ۹ هفتگی بودند.

۲۲۱
 ۲۲۲

۲۲۳
 ۲۲۴

۲۲۵
 ۲۲۶

۲۲۷
 ۲۲۸

۲۲۹
 ۲۳۰

۲۳۱
 ۲۳۲

۲۳۳
 ۲۳۴

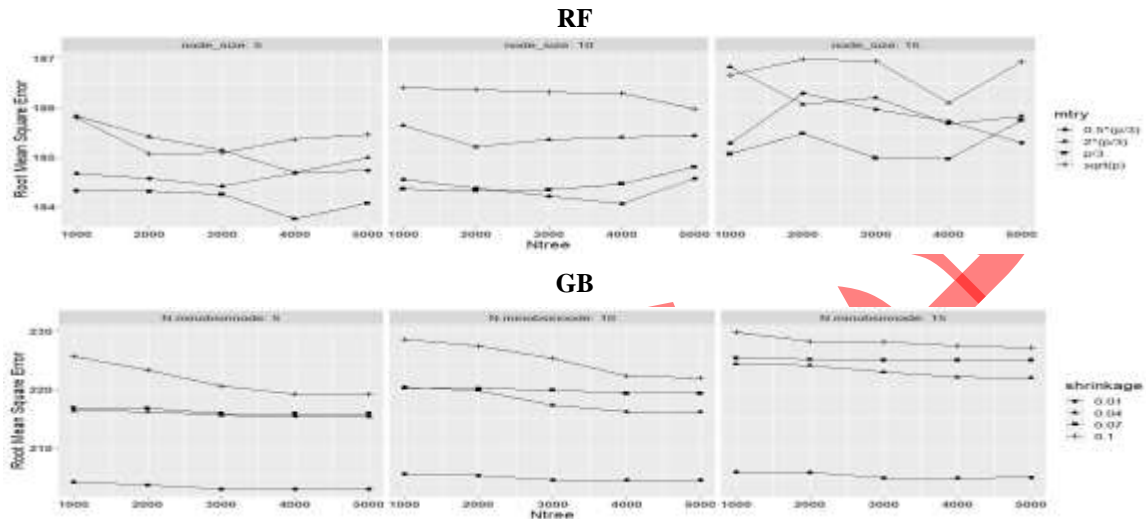
۲۳۵
 ۲۳۶

۲۳۷
 ۲۳۸

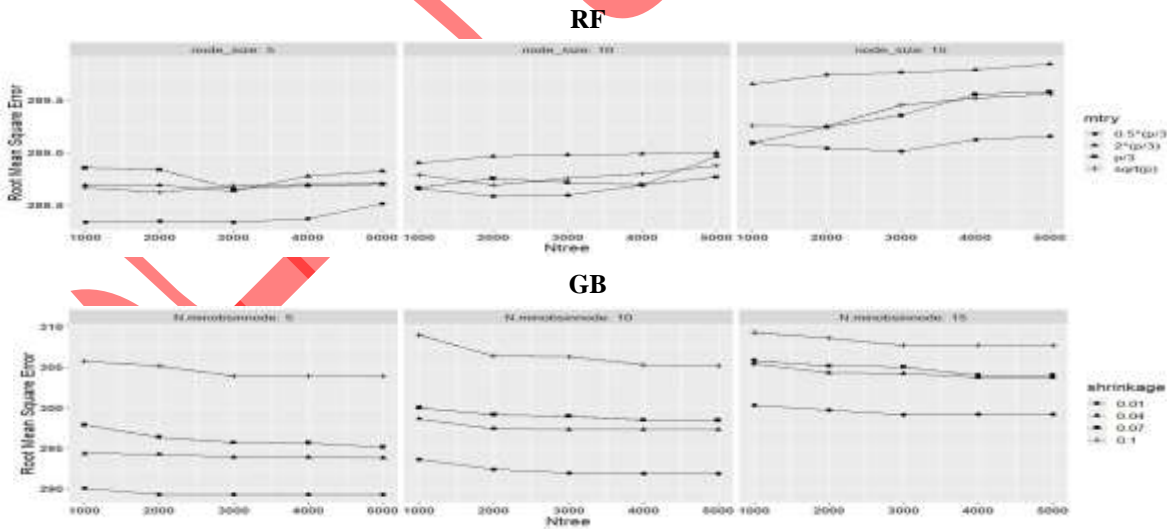
۲۳۹
 ۲۴۰

۲۴۱
 ۲۴۲

۲۴۳
 ۲۴۴

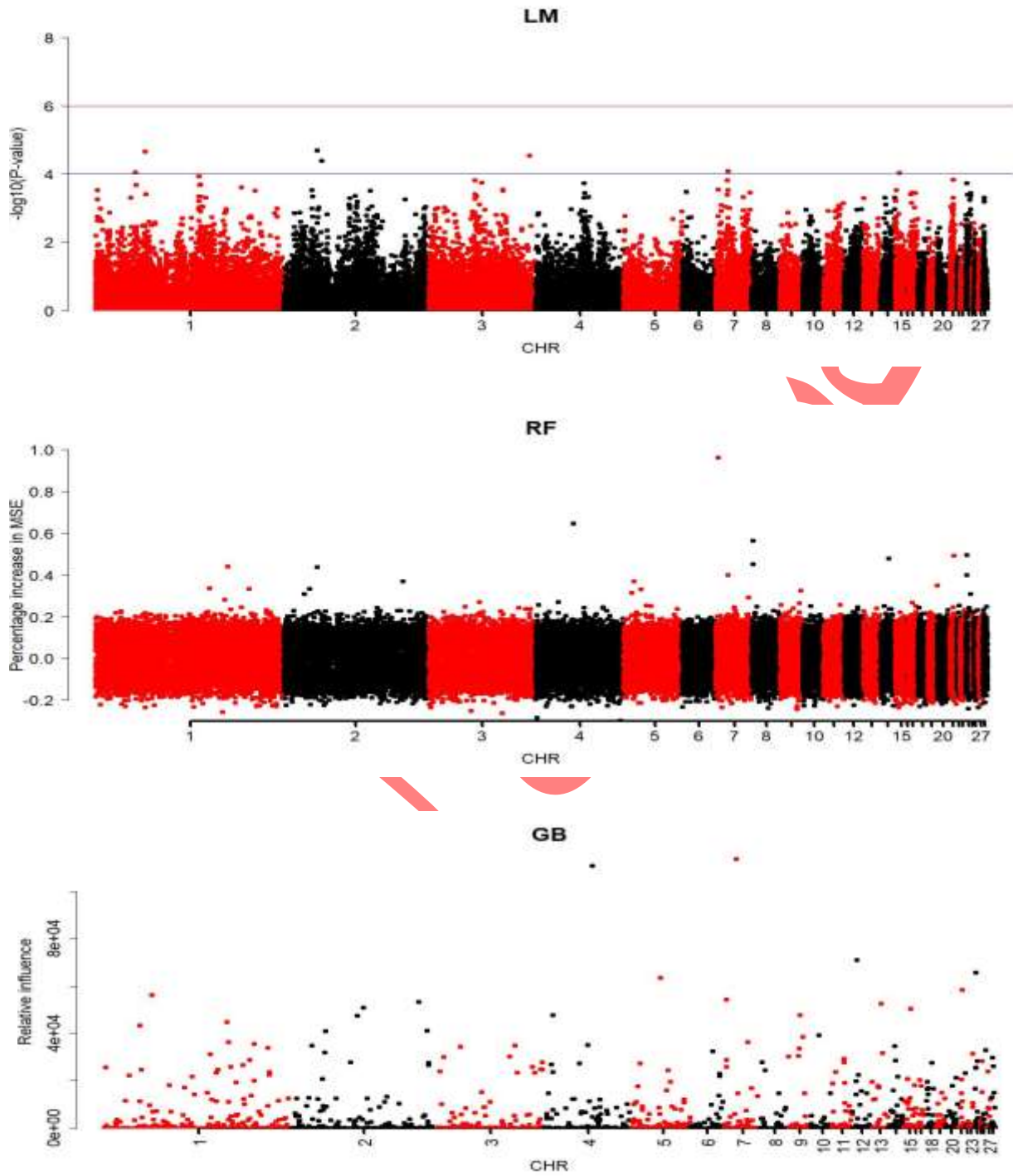


شکل ۱. ارتباط بین ترکیب‌های مختلف پارامترها و ریشه میانگین مربعات خطا برای روش‌های جنگل تصادفی (RF) و گرادیان بوستینگ (GB) در صفت وزن بدن در ۶ هفتگی. mtry: تعداد نشان‌گرهای انتخاب شده برای ساخت درخت، P: تعداد کل SNP ها، محور X به اندازه درخت جنگل (NTree) است.



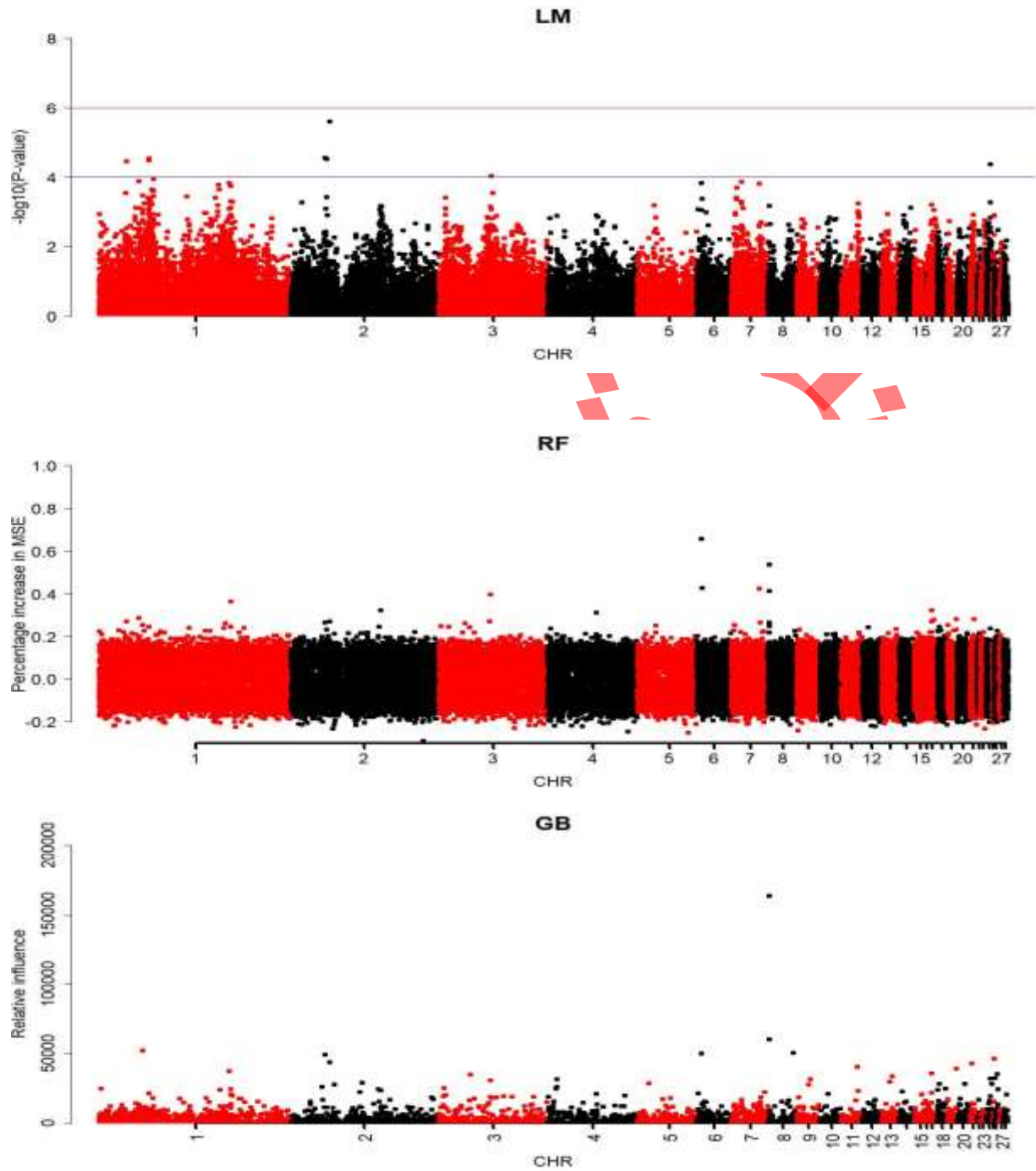
شکل ۲. ارتباط بین ترکیب‌های مختلف پارامترها و ریشه میانگین مربعات خطا برای روش‌های جنگل تصادفی (RF) و گرادیان بوستینگ (GB) در صفت وزن بدن در ۹ هفتگی. mtry: تعداد نشان‌گرهای انتخاب شده برای ساخت درخت، P: تعداد کل SNP ها، محور X به اندازه درخت جنگل (NTree) است.

۲۴۵
۲۴۶
۲۴۷
۲۴۸
۲۴۹
۲۵۰
۲۵۱
۲۵۲
۲۵۳
۲۵۴
۲۵۵
۲۵۶
۲۵۷
۲۵۸
۲۵۹
۲۶۰
۲۶۱
۲۶۲
۲۶۳
۲۶۴
۲۶۵
۲۶۶
۲۶۷
۲۶۸
۲۶۹
۲۷۰
۲۷۱
۲۷۲
۲۷۳
۲۷۴
۲۷۵
۲۷۶
۲۷۷
۲۷۸
۲۷۹



شکل ۳. پلات منهن برای صفت وزن بدن در ۶ هفتهگی با استفاده از روش‌های مدل خطی (LM)، جنگل تصادفی (RF) و گرادیان بوستینگ (GB)

۲۸۰
۲۸۱
۲۸۲
۲۸۳
۲۸۴
۲۸۵
۲۸۶
۲۸۷
۲۸۸
۲۸۹
۲۹۰
۲۹۱
۲۹۲
۲۹۳
۲۹۴
۲۹۵
۲۹۶
۲۹۷
۲۹۸
۲۹۹
۳۰۰
۳۰۱
۳۰۲
۳۰۳
۳۰۴
۳۰۵
۳۰۶
۳۰۷
۳۰۸
۳۰۹
۳۱۰
۳۱۱
۳۱۲
۳۱۳
۳۱۴
۳۱۵



شکل ۴. پلات منهن برای صفت وزن بدن در ۹ هفتگی با استفاده از روش‌های مدل خطی (LM)، جنگل تصادفی (RF) و گرادیان بوستینگ (GB)

جدول ۲. ده نشان‌گرهای برتر از مدل خطی (LM)، جنگل‌های تصادفی (RF) و گرادیان بوستینگ (GB) برای صفت وزن بدن در ۶ هفته‌گی

Method	Rank	Chr.	SNP	Position (bp)	p-value	MAF
LM	۱	۲	GGaluGA141221	۳۵۹۱۹۵۲۹	۲/۰۱E-۰۵	۰/۱۷
	۲	۱	Gga_rs13866016	۵۱۴۲۵۹۷۸	۲/۱۸E-۰۵	۰/۴۰
	۳	۳	GGaluGA237887	۱۰۷۵۰۸۲۱۹	۲/۸۶E-۰۵	۰/۱۴
	۴	۲	GGaluGA142838	۴۰۵۲۸۶۶۶	۴/۰۲E-۰۵	۰/۳۳
	۵	۷	Gga_rs13743072	۱۱۹۳۲۶۸۰	۸/۲۱E-۰۵	۰/۳۴
	۶	۱	GGaluGA014205	۴۲۴۹۸۳۳۲	۸/۹۲E-۰۵	۰/۰۷
	۷	۱۵	Gga_rs13629367	۴۳۲۰۰۲۰	۹/۳۹E-۰۵	۰/۰۶
	۸	۱	Gga_rs13918443	۱۱۰۶۵۴۵۱۲	۰/۰۰۰۱۱۸۷	۰/۱۹
	۹	۲۱	Gga_rs14284004	۳۵۰۲۵۶۳	۰/۰۰۰۱۴۹۵	۰/۱۶
	۱۰	۷	Gga_rs14607739	۱۲۷۳۲۸۴۶	۰/۰۰۰۱۵۲۹	۰/۱۶
Method	Rank	Chr.	SNP	Position (bp)	Percentage increase in MSE	
RF	۱	۷	GGaluGA308573	۲۵۶۴۱۳۲	۰/۹۶	۰/۲۸
	۲	۴	GGaluGA255033	۳۹۰۳۷۸۶۸	۰/۶۵	۰/۰۷
	۳	۸	Gga_rs13614212	۱۳۹۴۳۶۹	۰/۵۶	۰/۲۳
	۴	۲۴	Gga_rs14292823	۱۱۵۶۸۴۰	۰/۵۰	۰/۴۱
	۵	۲۱	GGaluGA184364	۴۳۰۹۲۶۳	۰/۴۹	۰/۲۶
	۶	۱۴	Gga_rs14076673	۸۷۹۸۳۷۸	۰/۴۹	۰/۰۹
	۷	۸	GGaluGA322130	۱۵۴۱۱۵۲	۰/۴۸	۰/۲۱
	۸	۱	Gga_rs10727436	۱۴۱۰۸۱۷۸۸	۰/۴۵	۰/۲۹
	۹	۲	GGaluGA141221	۳۵۹۱۹۵۲۹	۰/۴۴	۰/۱۷
	۱۰	۲۴	Gga_rs14292795	۱۱۳۰۳۷۶	۰/۴۱	۰/۴۱
Method	Rank	Chr.	SNP	Position (bp)	Relative influence	
GB	۱	۷	Gga_rs13743072	۱۱۹۳۲۶۸۰	۱۱۳۹۲۹	۰/۳۴
	۲	۴	GGaluGA258772	۴۹۸۶۳۷۶۱	۱۱۱۰۰۵	۰/۱۸
	۳	۱۲	Gga_rs14034395	۴۹۲۴۱۰۹	۷۰۹۹۳	۰/۳۶
	۴	۲۴	Gga_rs15214346	۱۸۹۰۱۴۰	۶۵۷۳۹	۰/۱۵
	۵	۵	Gga_rs14527434	۳۰۱۸۰۳۹۳	۶۳۶۵۶	۰/۱۲
	۶	۲۱	GGaluGA184364	۴۳۰۹۲۶۳	۵۸۵۵۶	۰/۲۶
	۷	۱	GGaluGA016965	۵۱۴۸۱۷۵۲	۵۶۳۴۵	۰/۲۸
	۸	۷	GGaluGA308573	۲۵۶۴۱۳۲	۵۴۳۶۸	۰/۲۹
	۹	۲	Gga_rs14251348	۱۳۴۹۰۹۲۹۸	۵۳۳۰۴	۰/۴۷
	۱۰	۱۳	GGaluGA094688	۱۰۸۲۸۸۶۵	۵۲۶۲۳	۰/۳۰

Chr: شماره کروموزوم، Rank: رتبه نشان‌گر، SNP: نام نشان‌گر، MAF: فراوانی آلل کمیاب.

جدول ۳. ده نشان‌گرهای برتر از مدل خطی (LM)، جنگل‌های تصادفی (RF) و گرادیان بوستینگ (GB) برای صفت وزن بدن در ۹ هفتگی

Method	Rank	Chr.	SNP	Position (bp)	p-value	MAF
LM	۱	۲	GGaluGA142838	۴۱۰۹۳۸۶۸	۲/۵۱E-۰۵	۰/۳۳
	۲	۲	GGaluGA141221	۳۵۹۱۹۵۲۹	۲/۷۹E-۰۵	۰/۱۷
	۳	۱	Gga_rs13865536	۵۰۶۸۵۶۲۲	۲/۹۱E-۰۵	۰/۴۶
	۴	۲	GGaluGA141644	۳۷۴۱۵۱۸۲	۳/۰۲E-۰۵	۰/۲۸
	۵	۱	GGaluGA017356	۵۲۳۵۸۹۷۰	۳/۲۱E-۰۵	۰/۴۳
	۶	۱	Gga_rs13747634	۲۸۴۹۷۱۹۲	۳/۵۴E-۰۵	۰/۳۰
	۷	۲۴	Gga_rs14295712	۳۹۱۴۶۹۹	۴/۲۱E-۰۵	۰/۲۹
	۸	۳	Gga_rs16277926	۵۴۳۳۰۷۰۶	۹/۲۴E-۰۵	۰/۴۵
	۹	۱	Gga_rs13654297	۵۷۳۵۵۲۰۶	۰/۰۰۰۱۱۲	۰/۲۵
	۱۰	۱	Gga_rs14814275	۴۱۴۴۷۹۳۷	۰/۰۰۰۱۳۱	۰/۴۲
Method	Rank	Chr.	SNP	Position (bp)	Percentage increase in MSE	
RF	۱	۶	Gga_rs15763229	۶۰۵۹۴۱۰	۰/۶۶	۰/۰۸
	۲	۸	Gga_rs13614212	۱۳۹۴۳۶۹	۰/۵۴	۰/۳۳
	۳	۶	GGaluGA295472	۶۳۰۹۳۶۳	۰/۴۲	۰/۰۶
	۴	۷	GGaluGA318133	۲۹۲۵۵۶۱۶	۰/۴۲	۰/۴۴
	۵	۸	GGaluGA322130	۱۵۴۱۱۵۳	۰/۴۱	۰/۲۱
	۶	۳	Gga_rs14359646	۵۴۰۷۶۶۸۹	۰/۳۹	۰/۳۱
	۷	۱	Gga_rs13625427	۱۳۸۰۴۴۸۴۷	۰/۳۶	۰/۲۹
	۸	۲	GGaluGA158217	۹۳۸۲۵۶۲۱	۰/۳۲	۰/۴۶
	۹	۱۷	GGaluGA114401	۵۰۰۶۷۳۹	۰/۳۲	۰/۱۷
	۱۰	۴	GGaluGA258464	۵۱۲۲۶۹۵۱	۰/۳۱	۰/۲۴
Method	Rank	Chr.	SNP	Position (bp)	Relative influence	
GB	۱	۸	Gga_rs13614212	۱۳۹۴۳۶۹	۱۶۳۵۶۲	۰/۲۳
	۲	۸	GGaluGA322130	۱۵۴۱۱۵۳	۶۰۱۱۲	۰/۲۱
	۳	۱	Gga_rs13858398	۴۳۹۵۹۹۰۲	۵۲۲۹۱	۰/۳۰
	۴	۸	Gga_rs14654881	۲۶۳۶۶۷۷۵	۵۰۵۱۶	۰/۲۳
	۵	۶	Gga_rs15763229	۶۰۵۹۴۱۰	۵۰۱۴۶	۰/۰۸
	۶	۲	GGaluGA141221	۳۵۹۱۹۵۲۹	۴۹۴۲۴	۰/۱۷
	۷	۲۵	Gga_rs16723884	۱۹۴۴۰۹۶	۴۶۱۲۵	۰/۱۶
	۸	۲	GGaluGA142838	۴۱۰۹۳۸۶۸	۴۳۵۶۲	۰/۳۳
	۹	۲۱	Gga_rs16178040	۱۶۱۵۸۴۶	۴۳۰۷۵	۰/۰۸
	۱۰	۱۱	Gga_rs14695411	۱۶۶۹۸۶۴۴	۴۰۴۷۹	۰/۲۸

Chr: شماره کروموزوم، Rank: رتبه نشان‌گر، SNP: نام نشان‌گر، MAF: فراوانی آلل مینور.

در روش‌های مورد استفاده در مطالعه حاضر، نشان‌گرهای متفاوتی از لحاظ اهمیت رتبه‌بندی شدند ولی برخی از نشان‌گرها حداقل در دو روش، جزء ده نشان‌گرهای برتر بودند. در وزن ۶ هفتگی، نشان‌گرهای GGaluGA184364 و GGaluGA308573 در دو روش جنگل تصادفی و گرادیان بوستینگ، نشان‌گر GGaluGA141221 در دو روش جنگل تصادفی و مدل خطی، و نشان‌گر Gga_rs13743072 در دو روش مدل خطی و گرادیان بوستینگ با رتبه‌های متفاوت، مشترک و جز ده نشان‌گر برتر بودند. در وزن ۹ هفتگی، سه نشان‌گر Gga_rs15763229، Gga_rs13614212 و GGaluGA322130 در بین دو روش جنگل تصادفی و گرادیان

۳۳۴ بوستینگ و نشان گر GGAluGA142838 در بین دو روش مدل خطی و گرادیان بوستینگ با رتبه‌های متفاوت جزو ده نشان‌گر برتر
 ۳۳۵ بودند. در هر دو صفت، نشان‌گرهای مشترک بین دو روش یادگیری ماشین بیشتر بود که می‌تواند اهمیت آن‌ها بر صفات مورد نظر را
 ۳۳۶ نشان دهد. ممکن است در مدل‌های خطی به دلیل فرضیات و سخت‌گیری‌های موجود در معنی‌دار شدن نشان‌گرها، امکان شناسایی
 ۳۳۷ آن‌ها میسر نشده باشد.

۳۳۸ ژن‌های مرتبط با ۳ نشان‌گر مهم هر روش برای صفات وزن بدن در ۶ و ۹ هفتگی در جدول ۴ نشان داده شده است. نشان‌گر
 ۳۳۹ GGAluGA141221 توسط مدل خطی در دو صفت وزن بدن جز ۳ نشان‌گر برتر و برای وزن بدن در ۶ و ۹ هفتگی، به ترتیب، دارای
 ۳۴۰ رتبه ۱ و ۲ بود. این نشان‌گر، که توسط الگوریتم‌های یادگیری ماشین هم جز نشان‌گرهای برتر و با اهمیت برای هر دو صفت وزن
 ۳۴۱ بدن شناخته شده است، بر روی کروموزوم ۲ قرار دارد و با ژن *SGOLI* در ارتباط است. با وجود این که در مطالعه حاضر، داده‌ها برای
 ۳۴۲ ساختار جمعیت تصحیح نشدند، اما برخی نتایج حاصل از آن مثل شناسایی نشان‌گر مرتبط با ژن *SGOLI*، با یافته‌های Emrani و
 ۳۴۳ همکاران (۲۰۱۷) مطابقت داشت. وجود تفاوت در سایر نشان‌گرهای مطالعه حاضر با یافته‌های Emrani و همکاران (۲۰۱۷)، علی‌رغم
 ۳۴۴ استفاده از داده‌های مشابه، به این دلیل بود که این محققین پوشش ژنومی داده‌ها را با مدل خطی انجام دادند که در آن اثرات ساختار
 ۳۴۵ جمعیت، وزن بدن جوجه‌های یک‌روزه و پلی‌ژنیک منظور شده بود. حذف این عوامل در مدل خطی مطالعه حاضر به این دلیل بود که
 ۳۴۶ نتایج حاصل از این مدل با نتایج دو روش یادگیری ماشین (جنگل تصادفی و گرادیان بوستینگ)، که امکان برازش این عوامل در آن
 ۳۴۷ روش‌ها وجود نداشت (Arabnejad et al., 2020)، قابل مقایسه باشد.

۳۴۸ نشان‌گر Gga_rs13614212 یکی دیگر از نشان‌گرهای مهم برای هر دو صفت وزن بود که توسط الگوریتم‌های یادگیری ماشین
 ۳۴۹ شناسایی شد. این نشان‌گر، که در وزن ۶ هفتگی توسط روش جنگل تصادفی و در وزن ۹ هفتگی توسط هر دو روش یادگیری ماشین
 ۳۵۰ نیز به‌عنوان ۳ نشان‌گر برتر شناسایی شده بود، بر روی کروموزوم ۸ و در درون ژن *CAMSAP2* قرار داشت. پروتئین حاصل از این
 ۳۵۱ ژن به سر آنتی‌سانتروم‌های میکروتوبول‌ها متصل می‌شود و این امکان را به آن می‌دهد تا دینامیک، ساختار و پلیمریزاسیون
 ۳۵۲ میکروتوبول‌ها را کنترل کند (Jiang et al., 2014). این ژن از جمله ژن‌های معنی‌دار برای صفات لاشه در گاوهای گوشتی
 ۳۵۳ (Srikanth et al., 2020) و ضریب تبدیل خوراک در جوجه‌های گوشتی (Wang et al., 2019) گزارش شده است.

۳۵۴ در صفت وزن بدن در ۶ هفتگی، دو نشان‌گر GGAluGA308573 و Gga_rs13743072 در هر ۳ روش اهمیت زیادی داشتند.
 ۳۵۵ نشان‌گر GGAluGA308573 در روش‌های مدل خطی، جنگل تصادفی و گرادیان بوستینگ به ترتیب رتبه ۱۸، ۱ و ۸ و نشان‌گر
 ۳۵۶ Gga_rs13743072 به ترتیب رتبه ۵، ۱۱ و ۱ را داشتند. نشان‌گر GGAluGA308573 بر روی کروموزوم ۷ و در درون ژن *MAP2* قرار
 ۳۵۷ داشت. این ژن یک بخش اصلی از سیستم عصبی است و بیشتر در دندریت‌ها یافت شده و به‌عنوان فعال‌کننده و استحکام‌دهنده
 ۳۵۸ میکروتوبول و کنترل‌کننده سیستم‌های میکروتوبول در دندریت‌های نورون‌ها شناخته شده است که منجر به کشش دندریت می‌شود
 ۳۵۹ (Harada et al., 2002). در مطالعه Fehm و همکاران (۲۰۰۴) گزارش شده است که وزن بدن در انسان از طریق سیستم عصبی مرکزی
 ۳۶۰ تنظیم می‌شود، زیرا گیرنده‌های گلوکوکورتیکوئید و مینرالوکورتیکوئید در نورون‌ها، تعیین‌کننده تعادل بین فرآیندهای تخصیص قند و
 ۳۶۱ مصرف غذا هستند. ژن *MAP2* به‌عنوان یکی از ژن‌های مؤثر در تمایز بین جمعیت گوسفند غیر تجاری در آفریقای جنوبی (Nguni
 ۳۶۲ sheep)، که برای ویژگی تولیدی خاصی انتخاب نشده، و نژادهای تجاری (Mutton Merino and Dorset Horn)، که برای رشد و تولید
 ۳۶۳ گوشت انتخاب شده‌اند، معرفی شد (Dzomba et al., 2020). نشان‌گر Gga_rs13743072 در کروموزوم ۷ و در بالا دست ژن
 ۳۶۴ *FAM117B* قرار داشت. در مطالعات گزارش شده است که این ژن در انسان بر سطوح کلسترول بدن (Willer et al., 2013) و در
 ۳۶۵ جوجه‌های گوشتی، بر توسعه و رشد جنین مؤثر است (Kanachari et al., 2021). نوع دیگری از خانواده این ژن (*FAM184B*) با
 ۳۶۶ همان توالی، به‌عنوان ژن مؤثر بر وزن بدن در سنین ۴، ۶، ۱۰ و ۱۲ هفته در جوجه‌های گوشتی گزارش شده است (Jin et al., 2015).

۳۶۷
 ۳۶۸

۳۶۹
۳۷۰
۳۷۱
۳۷۲
۳۷۳

جدول ۴. سه نشان گر مهم شناسایی شده با روش های مدل خطی، جنگل تصادفی و گرادیان بوستینگ و ژن های مرتبط با آن ها برای صفات وزن بدن در ۶ هفته گی (BW6) و ۹ هفته گی (BW9)

SNP	Chr.	Position (bp)	Rank in LM	Rank in RF	Rank in GB	P-value	PIncmSE	RI	Gene	Location
BWT6										
GGaluGA141221	۲	۳۵۹۱۹۵۲۹	۱	۹	۱۹	۲/۰۱E-۰۵	۰/۴۴	۴۰۸۹۳	SGOL1	Intron
Gga_rs13866016	۱	۵۱۴۲۵۹۷۸	۲	۹۰	۵۱۹	۲/۱۸E-۰۵	۰/۲۲	۱۱۲۲	CYTH4	Intron
GGaluGA237887	۳	۱۰۷۵۰۸۲۱۹	۳	۱۱۴	۷۰۴	۲/۸۶E-۰۵	۰/۲۱	۳۴۴	DTNB	Intron
GGaluGA308573	۷	۲۵۶۴۱۳۲	۱۸	۱	۸	-/۰۰۰۳	۰/۹۶	۵۴۳۶۸	MAP2	Intron
GGaluGA255033	۴	۳۹۰۳۷۸۶۸	۲۴۳۷	۲	۹۸۱	-/۰۳۵	۰/۶۵	۶۵	ACSL1	Intron
Gga_rs13614212	۸	۱۳۹۴۳۶۹	۵۰۰۰	۳	۹۴۳	-/۰۰۸	۰/۵۶	۸۳	CAMSAP2	Intron
Gga_rs13743072	۷	۱۱۹۳۲۶۸۰	۵	۱۱	۱	۸/۲۱E-۰۵	۰/۴۰	۱۱۳۹۳۹	FAM117B	Upstream
GGaluGA258772	۴	۴۹۸۶۳۷۶۱	۱۷۷	۶۵۵	۲	-/۰۰۱	۰/۱۷	۱۱۱۰۰۵	SLC4A4	Intron
Gga_rs14034395	۱۲	۴۹۲۴۱۰۹	۴۸۹۱	۹۹۴	۳	-/۰۰۷	۰/۱۶	۷۰۹۹۳	TIMP4	Intron
BWT9										
GGaluGA142838	۲	۴۱۰۹۳۸۶۸	۱	۱۶	۸	۲/۵۱E-۰۵	۰/۲۸	۴۳۵۶۳	CNOT10	Intron
GGaluGA141221	۲	۳۵۹۱۹۵۲۹	۲	۲۲	۶	۲/۷۹E-۰۵	۰/۲۶	۴۹۴۲۴	SGOL1	Intron
Gga_rs13865536	۱	۵۰۶۸۵۶۲۲	۳	۹۸۵	۷۷۱	۲/۹۱E-۰۵	۰/۱۶	۵۱۷۰	MGAT3	Intron
Gga_rs15763229	۶	۶۰۵۹۴۱۰	۱۳	۱	۵	-/۰۰۰۱	۰/۶۶	۵۰۱۴۶	-	Intergenic
Gga_rs13614212	۸	۱۳۹۴۳۶۹	۵۱	۲	۱	-/۰۰۰۷	۰/۵۴	۱۶۳۵۶۲	CAMSAP2	Intron
GGaluGA295472	۶	۶۳۰۹۳۶۳	۳۵	۳	۱۲۸	-/۰۰۰۴	۰/۴۲	۱۳۷۳۳	LncRNA	Intron
Gga_rs13614212	۸	۱۳۹۴۳۶۹	۵۱	۲	۱	-/۰۰۰۷	۰/۵۴	۱۶۳۵۶۲	CAMSAP2	Intron
GGaluGA322130	۸	۱۵۴۱۱۵۳	۲۱۷۲	۵	۲	-/۰۰۳	۰/۴۱	۶۰۱۱۲	-	Intergenic
Gga_rs13858398	۱	۴۳۹۵۹۹۰۲	۱۳۴	۶۸۴	۳	-/۰۰۱	۰/۱۷	۵۳۲۹۱	LncRNA	Upstream

۳۷۴ Chr: تعداد کروموزوم، Position: موقعیت نشان گر، LM: مدل خطی، RF: جنگل تصادفی، GB: گرادیان بوستینگ، PIncmSE: درصد افزایش میانگین
۳۷۵ مربعات خطا، RI: تاثیر نسبی.

۳۷۶
۳۷۷ برای صفت وزن ۶ هفته گی، دو ژن *CYTH4* و *DTNB* در روش مدل خطی، ژن *ACSL1* در روش جنگل تصادفی و دو ژن
۳۷۸ *SLC4A4* و *TIMP4* در روش گرادیان بوستینگ در ارتباط با نشان گرهای با اهمیت شناسایی شدند. ژن *CYTH4*، ژن کدگذار برای
۳۷۹ تولید پروتئینی به نام Cytohesin4 است که عضوی از خانواده ی پروتئین های Cytohesine می باشد و در فعال سازی سیگنال های
۳۸۰ داخلی سلولی نقش دارد. این گروه پروتئینی ممکن است در عملکردهای متعددی از جمله تنظیم سوخت و ساز، تنظیم رشد و اندام
۳۸۱ سلولی و فعالیت های نورونی نقش داشته باشند (Teuliere et al., 2014; Cremonesi et al., 2012). ژن *DTNB*، که توسط Emrani
۳۸۲ و همکاران (۲۰۱۷) برای صفات رشد در جوجه های گوشتی گزارش شده است، در تولید کمپلکس پروتئین مرتبط با دیستروفین
۳۸۳ (DPC) نقش دارد. اختلال در این کمپلکس پروتئینی، که در سارکولم مشاهده می شود، با اشکال مختلف دیستروفی عضلانی همراه
۳۸۴ است. ژن *ACSL1* تولید پروتئینی به همین نام در جوجه های گوشتی را کد می کند و نقش مهمی در متابولیسم چربی ها و اسیدهای
۳۸۵ چرب زنجیره بلند دارد. در جوجه های گوشتی، استفاده از ژن *ACSL1* می تواند به عنوان یکی از راهکارهای بهبود ترکیبات چربی در
۳۸۶ بافت عضلانی و سیستم ایمنی جوجه ها مورد استفاده قرار گیرد. این ژن می تواند بر رشد و توسعه عضلات جوجه های گوشتی تأثیر
۳۸۷ داشته و باعث بهبود کیفیت گوشت جوجه ها شود (Liu et al 2019., Tian et al., 2021). ژن *SLC4A4* پروتئینی به نام NHE4 را کد
۳۸۸ می کند که نقش اصلی این پروتئین در تنظیم غلظت یون هیدروژن و بی کربنات در سلول ها است. در جوجه های گوشتی، ژن

۳۸۹ *SLC4A4* نقش مهمی در تنظیم تعادل اسیدی-بازی و موازنه آب و الکترولیت دارد. این ژن در بافت‌های مخاطی روده‌ها و کلیه‌ها
 ۳۹۰ فعالیت می‌کند و در حفظ تعادل pH خون بدن نقش دارد. ژن *SLC4A4* از جمله ژن‌های تاثیرگذار در کنترل استرس گرمایی، چربی-
 ۳۹۱ های شکمی و رشد در پرندگان است (Resnyk et al., 2017; Bahadoran et al., 2018). ژن *TIMP4* کد کننده پروتین و عضوی از
 ۳۹۲ خانواده تیمپس (TIMPs) است که به عنوان مهارکننده طبیعی آنزیم‌های فعالیت متالوپروتئیناز (Metalloproteinases) عمل می‌کند.
 ۳۹۳ آنزیم‌های متالوپروتئیناز انواعی از آنزیم‌ها هستند که قابلیت تخریب و تفکیک پروتین‌ها را دارند. این آنزیم‌ها در فرآیندهایی مانند
 ۳۹۴ رشد و توسعه سلولی، ترمیم بافت‌های ملتهب و آسیب‌دیده نقش مهمی داشته و از جمله ژن‌های تاثیرگذار در رشد عضله در انسان و
 ۳۹۵ توسعه بافت روده جوجه‌های در حال رشد است (Brew et al., 2000; Hrabia et al., 2022). در مطالعه Yue و همکاران (۲۰۲۲)
 ۳۹۶ گزارش شده است که این ژن از جمله ژن‌های تاثیرگذار بر استحکام استخوان در پرندگان تخم‌گذار است.
 ۳۹۷ برای صفت وزن بدن در ۹ هفته‌گی، نشان‌گرهای Gga_rs15763229 و GGaluGA322130، که به ترتیب بر روی کروموزوم
 ۳۹۸ های ۶ و ۸ قرار داشتند، توسط دو روش یادگیری ماشین جزو نشان‌گرهای با اهمیت شناخته شدند که هر دو از نشان‌گرهای بین ژنی
 ۳۹۹ بودند. ژن‌های *CNOT10* و *MGAT3* توسط مدل خطی و ژن *LncRNA* توسط دو روش یادگیری ماشین برای صفت وزن بدن در ۹
 ۴۰۰ هفته‌گی شناسایی شدند. ژن *CNOT10* توسط Emrani و همکاران (۲۰۱۷) برای صفات رشد در جوجه‌های گوشتی گزارش شده است.
 ۴۰۱ ژن *MGAT3* (مانوزیل گلیکوزامینیل ترانسفراز ۳) ژنی است که مسئول تولید آنزیم *MGAT3* است. این آنزیم در حیوانات و انسان در
 ۴۰۲ فرایند گلیکوزیلاسیون، یعنی اضافه کردن گروه‌های گلیکوز به پروتین‌ها، نقش دارد. این فرآیند به پروتین‌ها کمک می‌کند تا ساختار
 ۴۰۳ و عملکرد صحیح خود را به دست آورند و در عملکردهای مختلف در سطح سلول و بافت شرکت کنند. این ژن در مطالعات به عنوان
 ۴۰۴ ژن تاثیرگذار بر صفات رشد و محدود کننده سرعت جذب چربی در رژیم غذایی در گوساله‌ها گزارش شده است (Lyu et al., 2021;
 ۴۰۵ Sun et al 2012). ژن *LncRNA* در جریان فعالیت ژنتیکی و تنظیم بیان ژن‌ها نقش دارد و چون در کدگذاری پروتین نقش ندارد به-
 ۴۰۶ عنوان "غیرکدگذارنده" نامیده می‌شوند. این ژن معمولاً بیشتر از ۲۰۰ نوکلئوتید طول دارد و به صورت تراکم بالا در نواحی ژنوم
 ۴۰۷ یافت می‌شود. بطور خاص، برخی از *LncRNA* ها در رشد و تکثیر عضله نقش دارند و می‌توانند در تنظیم بیان ژن‌ها و عملکرد
 ۴۰۸ پروتین‌های مرتبط با رشد و تکثیر عضلات مؤثر باشند (Ren et al., 2018). نتایج هستی‌شناسی (Gene Ontology) در سطح
 ۴۰۹ فرآیندهای بیولوژیکی با استفاده از ۱۰ نشانگر برتر انتخاب شده توسط ۳ روش برای صفات وزن بدن در جدول ۵ گزارش شده است.
 ۴۱۰

جدول ۵. هستی‌شناسی و فرآیند بیولوژی با استفاده از ۱۰ نشانگر برتر تعیین شده با روش‌های رگرسیون خطی، جنگل تصادفی و گرادیان بوستینگ

نام ژن	p-value	توصیف هستی‌شناسی	شناسه هستی‌شناسی
رگرسیون خطی			
CACNA1I, PPP1R12A, GPR85, RARB, IGF1	0.008	signal transduction	GO:0007165
THRB, RARB	0.012	retinoic acid receptor signaling pathway	GO:0048384
THRB, RARB	0.029	hormone-mediated signaling pathway	GO:0009755
جنگل تصادفی			
ARHGAP32, TAGAP, KALRN	0.005	regulation of small GTPase mediated signal transduction	GO:0051056
MAP2, CAMSAP2, CDK5RAP2	0.007	microtubule cytoskeleton organization	GO:0000226
MAP2, CAMSAP2	0.011	regulation of microtubule polymerization	GO:0031113
BRINP1, MAP2	0.018	central nervous system neuron development	GO:0021954
MAP2, CDK5RAP2	0.030	microtubule bundle formation	GO:0001578
KCNJ5, KCNJ1	0.040	potassium ion import across plasma membrane	GO:1990573
BRINP1, FZD7	0.064	cellular response to retinoic acid	GO:0071300
KCNJ5, KCNJ1	0.069	potassium ion transport	GO:0006813
SGO1, CDK5RAP2	0.085	chromosome segregation	GO:0007059
گرادیان بوستینگ			
RER1, PLPP3	0.038	retrograde vesicle-mediated transport, Golgi to ER	GO:0006890
FZD7, PLPP3	0.075	canonical Wnt signaling pathway	GO:0060070
SPEN, TIMP4	0.087	Notch signaling pathway	GO:0007219

۴۱۱
۴۱۲ ژن‌های گزارش شده توسط روش جنگل تصادفی فرآیندهای بیولوژیکی بیش‌تری را نسبت به دو روش دیگر نشان دادند. مسیرهای
۴۱۳ بیولوژیکی ژن‌های کشف شده توسط این روش، نقش‌های بسیار مهمی در سیستم عصبی، استحکام اسکلت سلولی، تقسیم سلولی و
۴۱۴ کنترل‌کننده سیستم‌های میکروتوبولی دارند که این فرآیندها می‌توانند وزن بدن را تحت تاثیر قرار دهند (Fehm et al 2014). از جمله
۴۱۵ این ژن‌ها می‌توان به دو ژن *MAP2* و *CAMSAP2* اشاره کرد که با ۳ نشانگر برتر شناسایی شده توسط روش‌های یادگیری ماشین
۴۱۶ مطالعه حاضر مرتبط بوده و در بیشتر فرآیندهای بیولوژیکی بدن شرکت دارند. مشاهده نشان‌گرهای با اهمیت و ژن‌های مرتبط با
۴۱۷ آن‌ها توسط روش‌های یادگیری ماشین در مقابل با مدل خطی نشان می‌دهد که ممکن است برخی از نشان‌گرها و ژن‌های تاثیرگذار
۴۱۸ بر صفات به دلیل فرضیات موجود در مدل‌های خطی کشف نشوند. در مطالعات به تاثیر تعداد کم نمونه‌ها بر افزایش نرخ مثبت کاذب
۴۱۹ نتایج حاصل از مدل‌های خطی اشاره شده است (Hong and Park, 2012). بنابراین استفاده از رویکردهای ترکیبی و مقایسه مدل‌های
۴۲۰ مختلف می‌تواند راه‌کار مناسبی برای پیدا کردن نشان‌گرهای مرتبط با صفات اقتصادی باشد.

۴۲۱ ۴۲۲ نتیجه گیری و پیشنهادها

۴۲۳ در مطالعه حاضر، پویس ژنومی صفات وزن بدن در جوجه‌های گوشتی در سنین ۶ و ۹ هفتگی با روش‌های جنگل تصادفی، گرادیان
۴۲۴ بوستینگ، و مدل خطی ساده، جهت شناسایی نشان‌گرهای مهم با استفاده از تراشه 60K چندشکلی‌های تک‌نوکلئوتیدی انجام شد.
۴۲۵ نشان‌گرهای مهمی از جمله *GGaluGA141221* و *GGaluGA142838* توسط مدل خطی و نشانگرهای *GGaluGA308573*.
۴۲۶ *Gga_rs13743072*، *Gga_rs13614212*، *GGaluGA322130* و *Gga_rs15763229* توسط روش‌های یادگیری ماشین برای صفات
۴۲۷ وزن بدن در ۶ و ۹ هفتگی شناسایی شدند. نشان‌گرهای شناسایی شده در روش‌های یادگیری ماشین با ژن‌های جدیدی از جمله
۴۲۸ *MAP2*، *ACSL1*، *CAMSAP2*، *FAM117B*، *SLC4A4*، *TIMP4* و *LncRNA* مرتبط بودند. مطالعات هستی‌شناسی ژن نشان داد که
۴۲۹ ژن‌های کشف شده توسط روش جنگل تصادفی در عملکردهای بیولوژیکی متفاوت و مهمی نقش داشتند که قبلاً در جمعیت مورد
۴۳۰ مطالعه توسط مدل خطی گزارش نشده بودند. بنابراین به نظر می‌رسد به جای استفاده از مدل‌های آماری رایج مانند رگرسیون خطی
۴۳۱ به تنهایی، استفاده از روش‌های یادگیری ماشین و حتی ادغام روش‌ها با یکدیگر می‌تواند در تجزیه و تحلیل دقیق‌تر پویس ژنومی
۴۳۲ کمک‌کننده باشد.

۴۳۳ ۴۳۴ سپاسگزاری

۴۳۵ بدین‌وسیله، نگارندگان مقاله، از پروفسور Just Jensen از دانشگاه Aarhus دانمارک برای تعیین ژنوتیپ‌ها و تامین هزینه‌های مالی
۴۳۶ آن، و از شرکت کاب و نترس برای در اختیار قرار دادن تراشه ایلومینا ۶۰k قدردانی می‌کنند.

۴۳۷ ۴۳۸ منابع

۴۳۹ Arabnejad, M., Montgomery, C. G., Gaffney, P. M., McKinney, B. A. (2020). Nearest-neighbor projected distance
۴۴۰ regression for epistasis detection in GWAS with population structure correction. *Frontier Genetics*, 11:784.
۴۴۱ Bahadoran, S., Dehghani Samani, A., & Hassanpour, H. (2018). Effect of heat stress on the gene expression of ion
۴۴۲ transporters/channels in the uterus of laying hens during eggshell formation. *Stress*, 21(1), 51-58.
۴۴۳ Baştanlar, Y., & Özuysal, M. (2014). Introduction to machine learning. *Methods in Molecular Biology*, 1107: 105-
۴۴۴ 128.
۴۴۵ Boulesteix, A. L., Janitza, S., rupp, J. K., & König I. R. (2012). Overview of random forest methodology and
۴۴۶ practical guidance with emphasis on computational biology and bioinformatics. Technical Report. *Department of*
۴۴۷ *Statistics*, University of Munich.
۴۴۸ Breiman, L. (2001). Random forests. *Machine. Learning*, 45, 5-32.

- 449 Brew, K., Dinakarpanthian, D., & Nagase, H. (2000). Tissue inhibitors of metalloproteinases: evolution, structure
450 and function. *Biochimica et Biophysica Acta (BBA)-Protein Structure and Molecular Enzymology*, 1477(1-2), 267-
451 283.
- 452 Cha, J., Choo, H., Srikanth, K., Lee, S. H., Son, J. W., Park, M. R., Kim, N., Jang, G. W., & Park, J. E. (2021).
453 Genome-wide association study identifies 12 loci associated with body weight at age 8 weeks in Korean native
454 chickens. *Genes*, 12(8), 1170.
- 455 Cremonesi, P., Capoferri, R., Pisoni, G., Del Corvo, M., Strozzi, F., Rupp, R., Caillat, H., Modesto, P., Moroni, P.,
456 Williams, J. L., & Castiglioni, B. (2012). Response of the goat mammary gland to infection with *Staphylococcus*
457 *aureus* revealed by gene expression profiling in milk somatic and white blood cells. *BMC Genomics*, 13(1), 1-17.
- 458 Dadousis, C., Somavilla, A., Ilska, J. J., Johnsson, M., Batista, L., Mellanby, R. J., Headon, D., Gottardo, P.,
459 Whalen, A., Wilson, D., & Dunn, I. C. (2021). A genome-wide association analysis for body weight at 35 days
460 measured on 137,343 broiler chickens. *Genetics Selection Evolution*, 53, 1-14.
- 461 Dzomba, E. F., Chimonyo, M., Snyman, M. A., & Muchadeyi, F. C. (2020). The genomic architecture of South
462 African mutton, pelt, dual- purpose and nondescript sheep breeds relative to global sheep populations. *Animal*
463 *Genetics*, 51 (6), 910-923.
- 464 Emrani, H., Vaez Torshizi, R., Masoudi, A. A., & Ehsani, A. (2017). Identification of new loci for body weight
465 traits in F2 chicken population using genome-wide association study. *Livestock Science*, 206, 125–131.
- 466 Enoma, D. O., Bishung, J., Abiodun, T., Ogunlana, O., & Osamor, V. C. (2022). Machine learning approaches to
467 genome-wide association studies. *Journal of King Saud University-Science*, 34(4), 101847.
- 468 Fehm, L., Kern, W., & Peters, A. (2004). Body weight regulation through the central nervous system. The
469 development of a pathogenetically based adiposity therapy. *Medizinische Klinik*, 99 (11), 674-679.
- 470 Friedman, J. (2002). Stochastic gradient boosting. *Computational Statistics & Data Analysis*, 38,
471 367–378
- 472 Friedman, J. H. (2001). Greedy function approximation: a gradient boosting machine. *Annals of Statistics*, 29 (5),
473 1189-1232.
- 474 Goddard M. E., & Hayes, B. J. (2009) Mapping genes for complex traits in domestic
475 animals and their use in breeding programmes. *Nature Reviews Genetics*, 10, 381–
476 391.
- 477 Goldstein, B. A., Hubbard, A. E., Cutler A., & Barcellos L. F. (2010). An application of random forests to a
478 genome-wide association dataset: Methodological considerations and new findings. *BMC Genetics*, 11, 49.
- 479 Harada, A., Teng, J., Takei, Y., Oguchi, K., & Hirokawa, N. (2002). MAP2 is required for dendrite elongation, PKA
480 anchoring in dendrites, and proper PKA signal transduction. *The Journal of Cell Biology*, 158 (3), 541-549.
- 481 Hayes B. (2013). Overview of Statistical Methods for Genome-Wide Association Studies (GWAS). *Methods in*
482 *Molecular Biology*, 1019, 149-69.
- 483 Hong E. P. & Park J. W. (2012). Sample size and statistical power calculation in genetic association studies.
484 *Genomics Inform*, 10(2), 117-22.
- 485 Hrabia, A., Miska, K. B., Schreier, L. L., Proszkowiec-Weglarz, M. (2022) Altered gene expression of selected
486 matrix metalloproteinase system proteins in the broiler chicken gastrointestinal tract during post-hatch development
487 and coccidia infection. *Poultry Science*, 101(8):101915.
- 488 Hu, T., Darabos, C., & Urbanowicz, R. (2020). Machine learning in genome-wide association studies. *Frontiers in*
489 *Genetics*, 11, 593958.
- 490 Jin, C. F., Chen, Y. J., Yang, Z. Q., Shi, K., & Chen, C. K. (2015). A genome-wide association study of growth trait-
491 related single nucleotide polymorphisms in Chinese Yancheng chickens. *Genetics and Molecular Research*, 14 (4),
492 15783-15792.
- 493 Kanakachari, M., Ashwini, R., Chatterjee, R. N., & Bhattacharya, T. K. (2021). Transcriptome analysis reveals
494 potential mechanisms and pathways underlying embryonic development with respect to muscle growth and egg
495 production in slow and fast growing chickens. *BMC Genomics*, 13 (1), 58.
- 496 Li, B., Zhang, N., Wang, Y. G., George, A. W., Reverter, A., & Li, Y. (2018). Genomic prediction of breeding
497 values using a subset of SNPs identified by three machine learning methods. *Frontiers in Genetics*, 9, 237.
- 498 Liaw, A., & Wiener, M. (2015). randomForest: Breiman and Cutler's random forests for classification and
499 regression. *R package version*, 4, 14.
- 500 Liu, L., Cui, H., Xing, S., Zhao, G., & Wen, J. (2019). Effect of divergent selection for intramuscular fat content on
501 muscle lipid metabolism in chickens. *Animals*, 10(1), 4.
- 502 Lyu, S., Yang, P., Liu, Y., Song, T., Zhang, Z., Shi, Q., Chen, F., Liu, X., Li, Z., Ru, B., & Cai, C. (2021). Genetic
503 effects of MOGAT1 gene SNP in growth traits of Chinese cattle. *Gene*, 769, 145201.

004 Mebratie, W., Madsen, P., Hawken, R. Rome, H., Marois, D., Henshall, J., Bovenhuis, H., & Jensen J. (2019).
005 Genetic parameters for body weight and different definitions of residual feed intake in broiler chickens. *Genetics*
006 *Selection Evolution*, 51, 53.
007
008 Nicodemus, K. K., Malley, J. D., Strobl, C., & Ziegler, A. (2010). The behaviour of random forest permutation-
009 based variable importance measures under predictor correlation. *BMC Bioinformatics*. 11,110.
010 Purcell, S., Neale, B., Todd-Brown, K., Thomas, L., Ferreira, M. A., Bender, D., Maller, J., Sklar, P., De Bakker, P.
011 I., Daly, M. J., & Sham, P. C. (2007). PLINK: a tool set for whole-genome association and population-based
012 linkage analyses. *American Journal of Human Genetics*, 81 (3), 559-575.
013 Ren, T., Li, Z., Zhou, Y., Liu, X., Han, R., Wang, Y., Yan, F., Sun, G., Li, H., & Kang, X. (2018). Sequencing and
014 characterization of lncRNAs in the breast muscle of Gushi and Arbor Acres chickens. *Genome*, 61(5), 337-347.
015 Resnyk, C. W., Carré, W., Wang, X., Porter, T. E., Simon, J., Le Bihan-Duval, E., Duclos, M. J., Aggrey, S. E., &
016 Cogburn, L.A. (2017). Transcriptional analysis of abdominal fat in chickens divergently selected on bodyweight at
017 two ages reveals novel mechanisms controlling adiposity: validating visceral adipose tissue as a dynamic endocrine
018 and metabolic organ. *BMC Genomics*, 18, 1-31.
019 Ridgeway, G. (2013). Package 'GBM': Generalized Boosted Regression Models. *R Package version*, 2.
020 Srikanth, K., Lee, S. H., Chung, K. Y., Park, J. E., Jang, G. W., Park, M. R., Kim, N. Y., Kim, T. H., Chai, H. H.,
021 Sun, J., Zhang, C., Lan, X., Lei, C., & Chen, H. (2012). Exploring polymorphisms and associations of the bovine
022 MOGAT3 gene with growth traits. *Genome*, 55(1), 56-62.
023 Sun, S., Dong, B., & Zou, Q. (2021). Revisiting genome-wide association studies from statistical modelling to
024 machine learning. *Briefings in Bioinformatics*, 22(4), 263.
025 Teuliere, J., Cordes, S., Singhvi, A., Talavera, K., & Garriga, G. (2014). Asymmetric neuroblast divisions producing
026 apoptotic cells require the cytohesin GRP-1 in *Caenorhabditis Elegans*. *Genetics*, 198(1), 229-247.
027 Tian, W., Wang, D., Wang, Z., Jiang, K., Li, Z., Tian, Y., Kang, X., Liu, X., & Li, H. (2021). Evolution, expression
028 profile, and regulatory characteristics of ACSL gene family in chicken (*Gallus Gallus*). *Gene*, 764, 145094.
029 Wang, J., Yuan, X., Ye, S., Huang, S., He, Y., Zhang, H., Li, J., Zhang, X., & Zhang, Z. (2019). Genome wide
030 association study on feed conversion ratio using imputed sequence data in chickens. *Asian-Australasian Journal of*
031 *Animal Sciences*, 32 (4), 494-500.
032 Willer, C. J., Schmidt, E. M., Sengupta, S., Peloso, G. M., Gustafsson, S., Kanoni, S., Ganna, A., Chen, J.,
033 Buchkovich, M.L., Mora, S., & Beckmann, J.S. (2013). Discovery and refinement of loci associated with lipid
034 levels. *Nature Genetics*, 45 (11), 1274-1283.
035 Wray, N. R., Yang, J., Hayes, B. J., Price, A. L., Goddard, M. E., & Visscher, P. M. (2013). Pitfalls of predicting
036 complex traits from SNPs. *Nature Reviews Genetics*, 14 (7), 507-515.
037 Yates, A., Akanni, W., Amode, M. R., Barrell, D., Billis, K., Carvalho-Silva, D., Cummins, C., Clapham, P.,
038 Fitzgerald, S., Gil, L., & Girón, C. G. (2016). Ensembl 2016. *Nucleic Acids Research*, 44, 710-716.
039 Yue, Q., Chen, Y., Chen, H., & Zhou, R. (2022). Transcriptome profile reveals novel candidate genes associated
040 with bone strength in end-of-lay hens. *Animal Biotechnology*, 1-9.
041 Zhang, G. X., Fan, Q. C., Zhang, T., Wang, J. Y., Wang, W. H., Xue, Q., & Wang, Y. J. (2015). Genome-wide
042 association study of growth traits in the Jinghai Yellow chicken. *Genetics and Molecular Research*, 14(4), 15331-
043 15338.

044 Extended Abstract

045 Introduction

046 By employing Genome-Wide Association Studies (GWAS) and identifying Single Nucleotide Polymorphisms
047 (SNPs) and their associations with genes influencing traits, the necessary molecular information for marker and gene
048 selection for improving quantitative traits has been provided. However, in GWAS studies, the most common
049 approach used to examine the effects of markers is simple regression and P-value, which disregards issues such as
050 increasing the false positive rate, overestimation of marker effects, ignoring linkage disequilibrium between
051 markers, incorrect assumption of independence among markers, assumption of all genomic variables following a
052 normal distribution, and neglecting the interaction effects between the SNP markers. Consequently, a suitable
053 alternative to address these problems is conducting GWAS based on machine learning methods. The main objective
054 of this study is to identify important and influential markers for body weight traits measured at 6 and 9 weeks of age
055 in F₂ chickens population using the RF and GB methods and comparing them with the linear model (LM) approach.
056
057

008

009 **Material and methods**

060 For the current study, body weight data at ages 6 and 9 weeks of 312 F₂ chickens, resulting from two-way
061 crossbreeding between fast-growing commercial Arin line (R) and indigenous fowls from West Azerbaijan province
062 (O) were used. At age of 70 days, DNA from blood samples of chickens was extracted and stored at -20°C. These
063 DNAs were used to identify genotype of each bird, using the 60k Illumina Chicken SNP BeadChip, containing
064 54,340 SNP markers provided by Cobb Vantress with cooperation of Arhus University of Denmark. The phenotypic
065 data were adjusted for sex and hatching effects, and three methods including linear model, random forest, and
066 gradient boosting were used to identify the important markers. The top ten markers for body weight traits were
067 identified for each method. In addition, genes located within the Mb-1 region above and below the three top
068 markers, as identified by each method in the genomic region, were determined using the NCBI and Ensemble
069 databases from the reference genome of chicken (*Gallus Gallus*).

070

071 **Results and discussion**

072 Identification of important markers and the corresponding genes using machine learning methods showed that
073 some markers and genes influencing traits might not be identified by linear model, indicating that machine learning
074 methods including random forest and gradient boosting were suitable tools for selecting important markers. These
075 associated markers were GGaluGA308573, GGaluGA255033, Gga_rs13614212, Gga_rs13743072,
076 GGaluGA258772, Gga_rs14034395, and Gga_rs13858398. By examining only 3 of the most important markers in
077 each method, new related genes such as *MAP2*, *ACSL1*, *CAMSAP2*, *FAM117B*, *SLC4A4*, *TIMP4*, and *LncRNA* were
078 identified which were not detected by linear model in previous studies. Literature results reported that these genes
079 are regulating microtubule-stabilizing activity, cell's shape, intestinal tissues during post-hatch development,
080 augmenting adipogenesis, tissue growth and morphogenesis, fatty acid metabolism and abdominal fat deposition,
081 intramuscular fat content, meat tenderness, and flavor, axon development, dendrite development, organelle
082 organization, bicarbonate secretion and absorption and intracellular pH, and induction of apoptosis (Programmed
083 death of cells) in chickens.

084

085 **Conclusion**

086 In this study, random forest, gradient boosting, and simple linear models were used to detect important markers
087 association with body weight traits at 6 and 9 weeks of ages in broiler chickens using a 60k Chicken SNP BeadChip.
088 The most important markers that were identified by linear model were GGaluGA141221, GGaluGA142838. For
089 machine learning methods, the top markers were GGaluGA308573, Gga_rs13743072, Gga_rs13614212,
090 GGaluGA322130, and Gga_rs15763229. These markers were associated with genes that control several
091 biochemical, physiological and biological functions in chickens. Results indicate that machine learning algorithms
092 were able to identify new genes for body weight traits that were not previously identified by linear model.