
title: neugene

author: Bipin Neupane

date: "2/20/2025"



Table of Contents

- Introduction
- Function Definitions
 - GWAS Using GLM
 - GWAS Using GLM + PCA Correction
- Loading Required Libraries
- Performing GWAS
- GLM+PCA (neugene) vs GWASbyCOR
- GLM+PCA (neugene) vs Blink C
- GLM+PCA (neugene) vs Machine Learning Model (RandomForest & XGBoost)
- Conclusion

1. Introduction

This document serves as a comprehensive manual for the **Statistical Genomics** course at **Washington State University (WSU), 2025**, instructed by **Dr. Zhiwu Zhang**. It is developed as part of the course assignments and provides a step-by-step tutorial on performing Genome-Wide Association Studies (GWAS) using General Linear Models (GLM) with and without Principal Component Analysis (PCA) correction. We compare our method against GWASbyCor, Blink C, and Machine learning models (RandomForest & XGboost) using simulations and evaluate performance using ROC curves.

First, we created an R package named “**neugene**” to implement and streamline the GWAS analysis process. The package aims to improve the accuracy and interpretability of GWAS results while maintaining computational efficiency, making it a valuable tool for genetic association studies. Further steps include:

2. Function Definitions

The function **run_gwas_glm** performs GWAS using a General Linear Model (GLM) while adjusting for user-provided covariates, whereas **run_gwas_glm_pca** extends this approach by incorporating PCA correction to account for population structure. Unlike **run_gwas_glm**, the **run_gwas_glm_pca** function automatically excludes collinear PCs before including them as cofactors, ensuring that only independent principal components are used. Functions also facilitate enhanced visualization.

2.1 GWAS Using GLM

```
#' Run GWAS with GLM and Enhanced Visualization
#
#' @description Performs GWAS using linear regression with improved visualization
#' @param y Phenotype data frame
#' @param X Genotype data frame
#' @param C Covariate data frame
#' @param snp_info SNP metadata data frame
#' @return List containing results and plots
#' @export
run_gwas_glm <- function(y, X, C, snp_info) {
  library(qqman)
```

```

library(ggplot2)

# Data preparation

y_vec <- as.numeric(y[, 2])
SNPs <- as.matrix(X[, -1])
Covs <- as.matrix(C[, -1])

# GWAS analysis

p_values <- sapply(1:ncol(SNPs), function(i) {
  model <- lm(y_vec ~ SNPs[, i] + Covs)
  summary(model)$coefficients[2, 4]
})

# Prepare results

results <- data.frame(SNP = colnames(SNPs), P = p_values) |>
  merge(snp_info, by = "SNP")

# Enhanced Manhattan plot

manhattan_plot <- function() {
  manhattan(results,
    chr = "Chromosome",
    bp = "Position",
    p = "P",
    snp = "SNP",
    col = c("#377eb8", "#4daf4a"),
    genomewideline = -log10(5e-8),
    suggestiveline = -log10(1e-5),
    main = "Manhattan Plot (GLM)",
    cex = 0.8,
    cex.axis = 0.9)
}

# Enhanced QQ plot

qq_plot <- function() {
  qq(results$P,
    main = "QQ Plot (GLM)",
    col = "#984ea3",
    cex = 0.8)
  abline(0, 1, col = "#ff7f00", lwd = 2)
}

```

```

list(results = results,
      manhattan_plot = manhattan_plot,
      qq_plot = qq_plot)
}

```

2.2 GWAS Using GLM + PCA Correction

```

#' Run GWAS with GLM + PCA Correction and Enhanced Visualization
#'
#' @description Performs PCA-adjusted GWAS with comprehensive visualization
#' @inheritParams run_gwas_glm
#' @param npc Number of principal components
#' @return List containing results and plots
#' @export
run_gwas_glm_pca <- function(y, X, C, snp_info, npc = 5) {
  library(qqman)
  library(ggplot2)
  library(gridExtra)
  # Data preparation
  y_vec <- as.numeric(y[, 2])
  X_snps <- as.matrix(X[, -1])
  Covs <- as.matrix(C[, -1])
  # PCA analysis
  pca_out <- prcomp(X_snps, center = TRUE, scale. = TRUE)
  PCs <- pca_out$x[, 1:npc]
  # Collinearity check
  design_matrix <- Covs
  current_rank <- qr(design_matrix)$rank
  valid_pcs <- numeric(0)
  for(i in 1:npc) {
    temp_design <- cbind(design_matrix, PCs[, i])
    if(qr(temp_design)$rank > current_rank) {
      valid_pcs <- c(valid_pcs, i)
      design_matrix <- temp_design
    }
  }
  # Final PCA plot
  ggplot(pca_out, aes(PC1, PC2)) +
    geom_point(data = data.frame(x = 0, y = 0), color = "red", size = 500) +
    geom_hline(yintercept = 0, color = "black") +
    geom_vline(xintercept = 0, color = "black") +
    theme_minimal() +
    theme(panel.grid = element_blank())
}
```

```

current_rank <- qr(temp_design)$rank
}

}

# GWAS with PCA

C_adj <- if(length(valid_pcs) > 0) cbind(Covs, PCs[, valid_pcs]) else Covs
p_values <- sapply(1:ncol(X_snps), function(i) {
  model <- lm(y_vec ~ X_snps[, i] + C_adj)
  summary(model)$coefficients[2, 4]
})

# Prepare results

results_pca <- data.frame(SNP = colnames(X_snps), P = p_values) |>
  merge(snp_info, by = "SNP")

# Visualization functions

manhattan_plot <- function() {

  manhattan(results_pca,
    chr = "Chromosome",
    bp = "Position",
    p = "P",
    snp = "SNP",
    col = c("#e41a1c", "#377eb8"),
    genomewideline = -log10(5e-8),
    suggestiveline = -log10(1e-5),
    main = "Manhattan Plot (GLM+PCA)",
    cex = 0.8,
    cex.axis = 0.9)
}

qq_plot <- function() {

  qq(results_pca$P,
    main = "QQ Plot (GLM+PCA)",
    col = "#984ea3",
    cex = 0.8)

  abline(0, 1, col = "#ff7f00", lwd = 2)
}

# PCA visualizations

variance <- pca_out$sdev^2 / sum(pca_out$sdev^2)

```

```

scree_plot <- ggplot(data.frame(PC = 1:npc, Variance = variance[1:npc]),
                      aes(x = PC, y = Variance)) +
  geom_col(fill = "steelblue", alpha = 0.8) +
  geom_line(color = "darkred", size = 1) +
  geom_point(size = 3, color = "darkred") +
  labs(title = "PCA Scree Plot", x = "Principal Component", y = "Variance Explained") +
  theme_minimal()

pc_scatter <- ggplot(data.frame(PC1 = pca_out$x[,1], PC2 = pca_out$x[,2]),
                      aes(x = PC1, y = PC2)) +
  geom_point(color = "#4daf4a", alpha = 0.6) +
  labs(title = "PC1 vs PC2", x = "Principal Component 1", y = "Principal Component 2")

+
  theme_minimal()

list(results = results_pca,
      manhattan_plot = manhattan_plot,
      qq_plot = qq_plot,
      pca_plots = list(scree_plot, pc_scatter),
      PCs_used = valid_pcs)
}

```

3. Loading Required Libraries

```

library(neugene)
library(ggplot2)
library(dplyr)
library(writexl)
library(pROC)
library(qqman)
library(randomForest)
library(xgboost)

```

4. Performing GWAS

We performed GWAS on our dataset using both `run_gwas_glm` and `run_gwas_glm_pca` functions, where the former applied a basic GLM approach with user-defined covariates, and the latter incorporated PCA correction while excluding collinear principal components. Both functions successfully generated p-values for SNP associations and provided plots to visualize the results, confirming that the package performed as expected.

```

phenotype <- read.table("phenotype.txt", header = TRUE)
genotype <- read.table("GAPIT_genotype.txt", header = TRUE)
covariates <- read.table("covariates.txt", header = TRUE)
snp_info <- read.table("snp_info.txt", header = TRUE)

#Convert chromosome format
snp_info <- snp_info %>%
  mutate(
    Chromosome = as.numeric(gsub("H", "", Chromosome)), # Remove 'H' directly
    Position = as.numeric(Position)
  ) %>%
  filter(!is.na(Chromosome), !is.na(Position))

# Verify critical data properties
stopifnot(
  "No SNPs remaining after filtering" = nrow(snp_info) > 0,
  "Genotype/SNP mismatch" = all(colnames(genotype)[-1] %in% snp_info$SNP)
)

# Run GLM GWAS
gwas_glm_results <- run_gwas_glm(y = phenotype, X = genotype, C = covariates, snp_info =
snp_info)
write_xlsx(gwas_glm_results$results, "GLM_GWAS_Results.xlsx")

```

This will output the result in an Excel-compatible format, structured as follows:

SNP	P	Chromosome	Position
BK_01	0.840749025	1	491121109
BK_03	0.529821035	7	121697186
BK_04	0.114224594	5	541937108

SNP	P	Chromosome	Position
BK_07	0.87851242	7	415054969
BK_08	0.90598484	3	56452786
BK_12	0.780503012	2	25878442

```
# Run GWAS with PCA

gwas_glm_pca_results <- run_gwas_glm_pca(y = phenotype, X = genotype, C = covariates, snp
_info = snp_info, npc = 5)

write_xlsx(gwas_glm_pca_results$results, "GLM_PCA_GWAS_Results.xlsx")
```

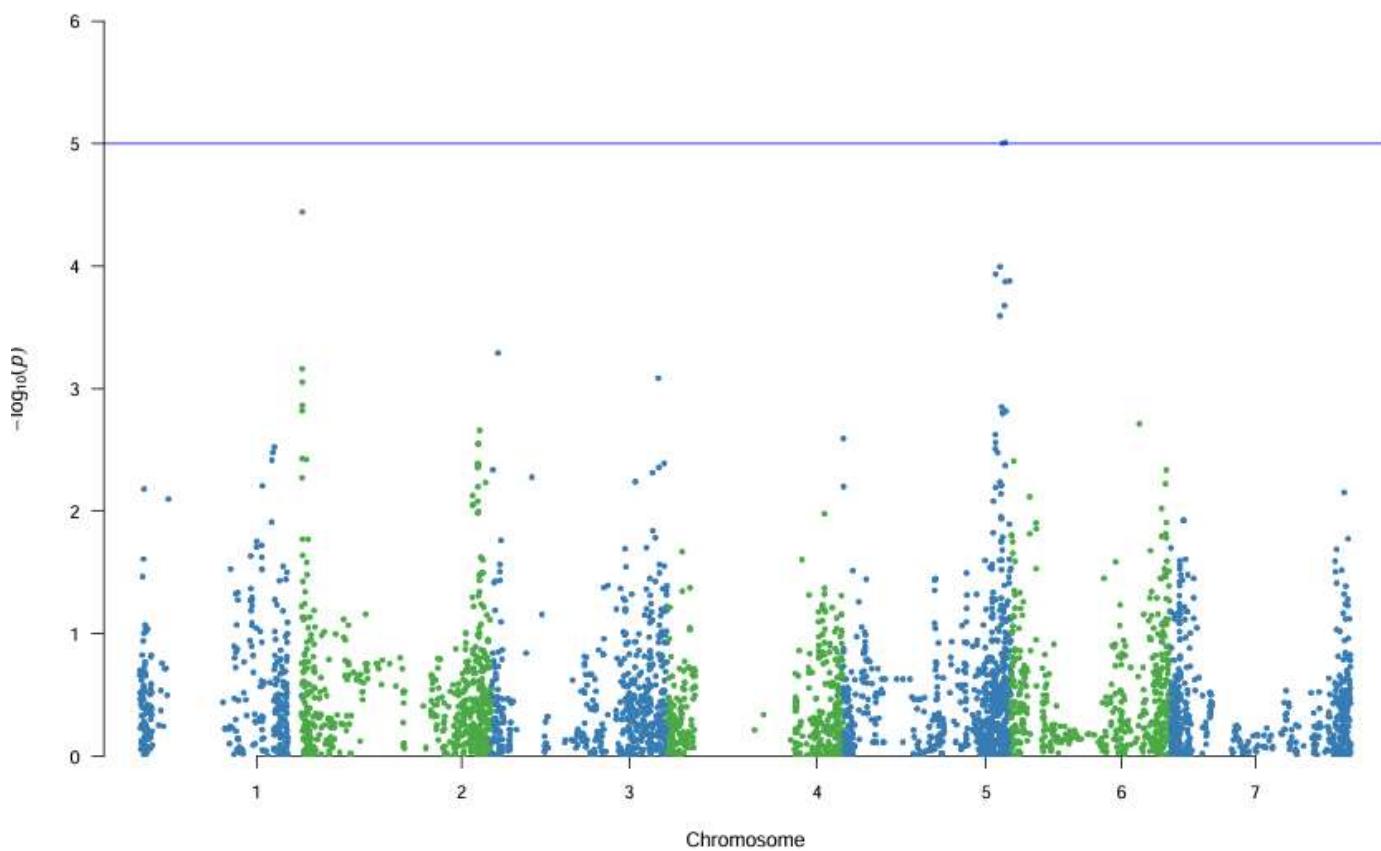
SNP	P	Chromosome	Position
BK_01	0.473908201	1	491121109
BK_03	0.935589477	7	121697186
BK_04	0.487831876	5	541937108
BK_07	0.64409368	7	415054969
BK_08	0.582603841	3	56452786
BK_12	0.622802889	2	25878442

```
# Generate plots with validation

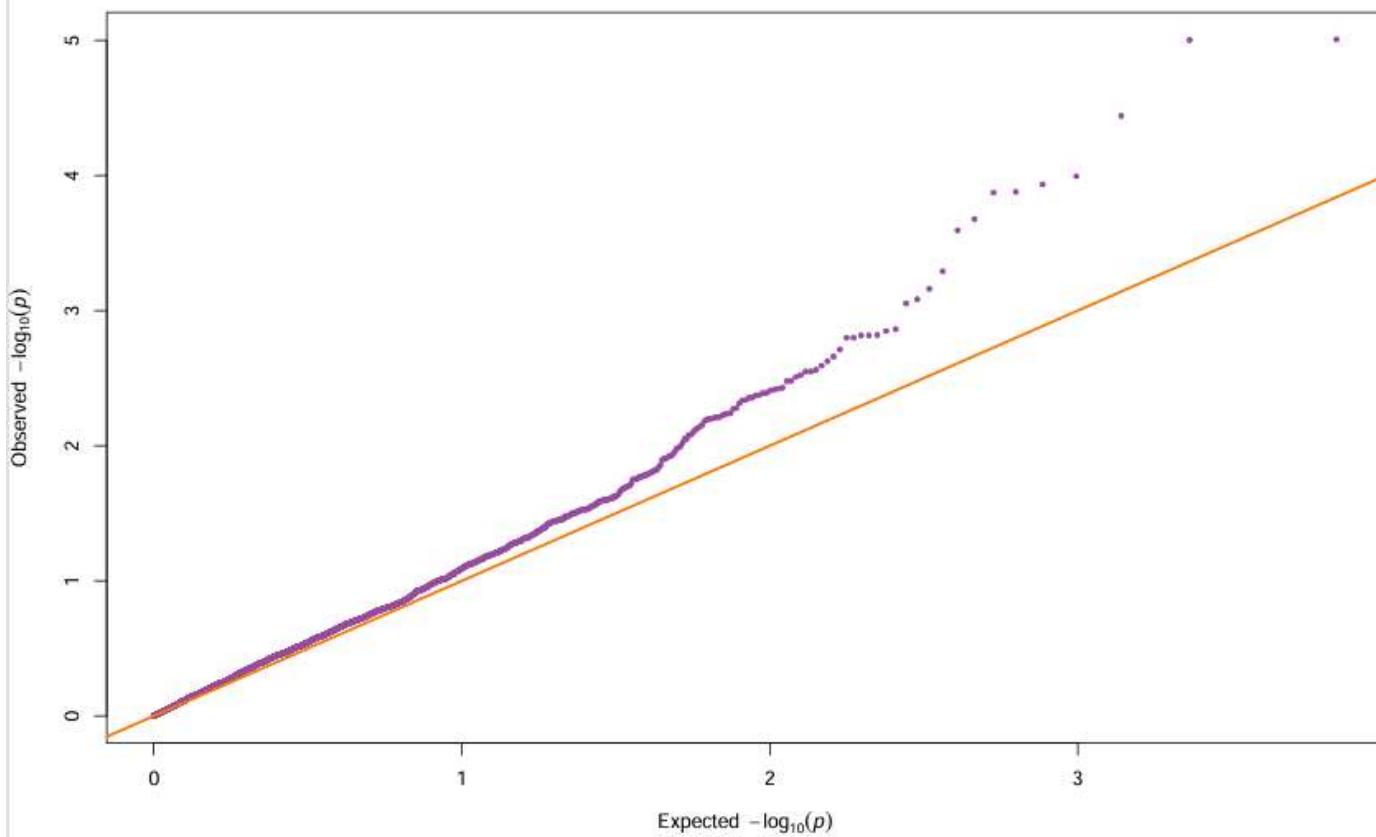
if(nrow(gwas_glm$results) > 0 && nrow(gwas_pca$results) > 0) {
  pdf("GWAS_Results.pdf", width = 12, height = 8)
  # GLM Plots
  gwas_glm$manhattan_plot()
  gwas_glm$qq_plot()
  # PCA-adjusted Plots
  gwas_pca$manhattan_plot()
  gwas_pca$qq_plot()
  # PCA diagnostics
  grid.arrange(
    gwas_pca$pca_plots[[1]] + theme(plot.margin = unit(c(1,1,1,1), "cm")),
    gwas_pca$pca_plots[[2]] + theme(plot.margin = unit(c(1,1,1,1), "cm")),
```

```
  ncol = 2  
}  
dev.off()  
}
```

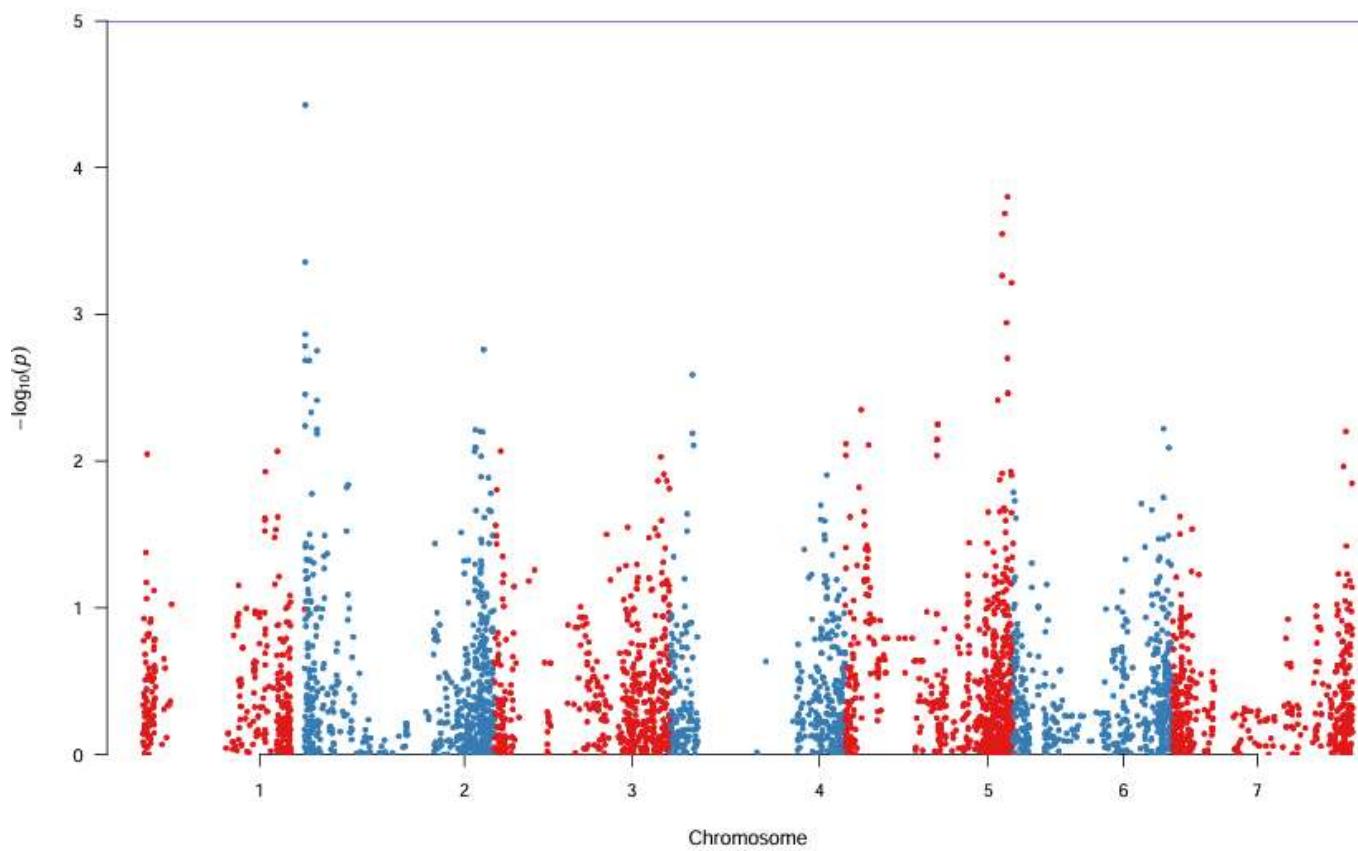
Manhattan Plot (GLM)



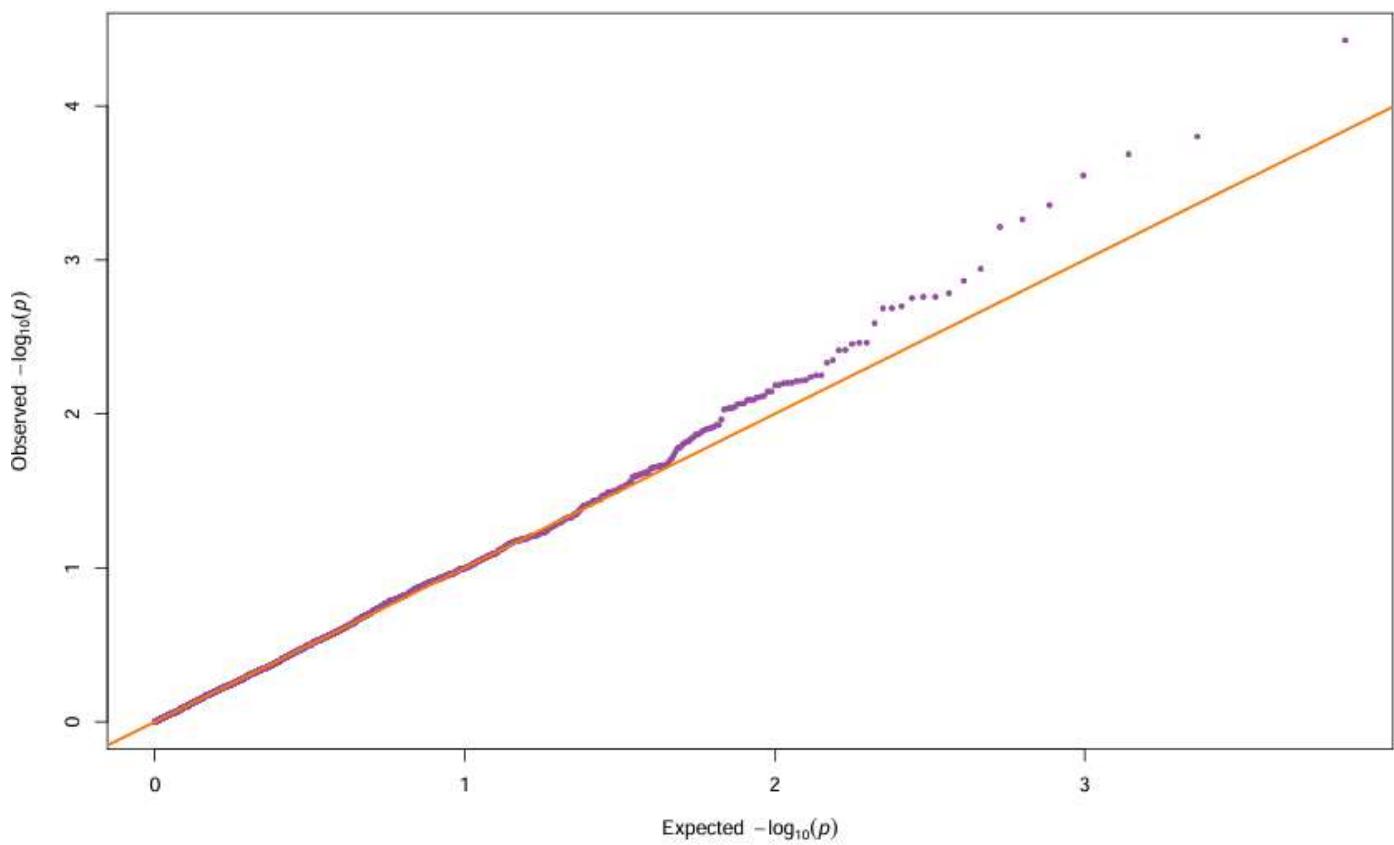
QQ Plot (GLM)

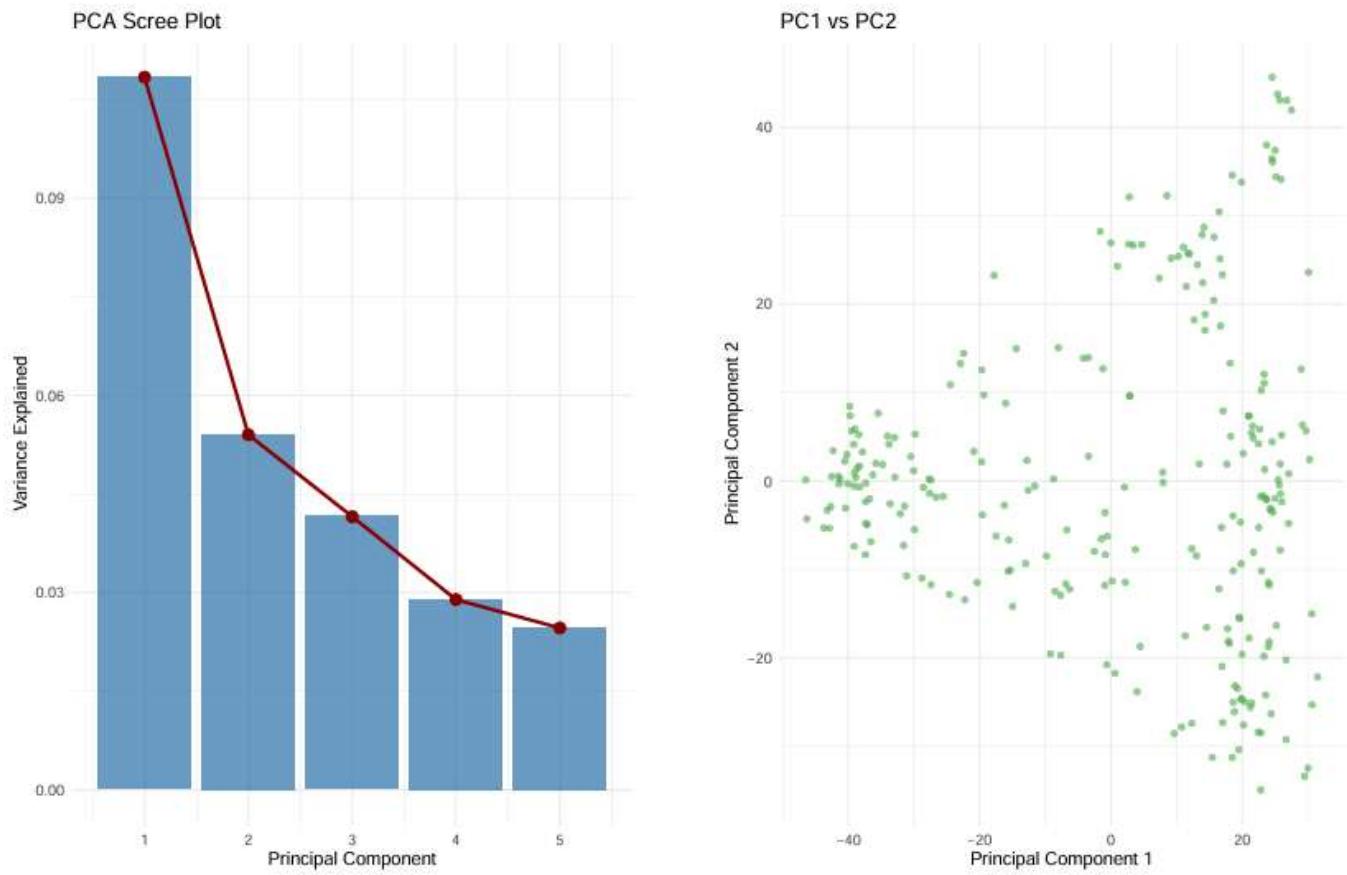


Manhattan Plot (GLM+PCA)



QQ Plot (GLM+PCA)





Despite the `neugene` package functioning as intended, we did not identify significant SNPs, likely due to factors such as limited sample size, stringent multiple testing corrections, and the polygenic nature of the trait under study. The absence of significant SNPs suggests that a larger dataset or a trait with stronger genetic effects may be necessary to detect meaningful associations.

5. GLM+PCA (`neugene`) vs GWASbyCor

In the next step, we will demonstrate the superiority of our GLM + PCA method over the competing `GWASbyCor` method through simulating our data with a minimum of 30 replicates, supported by ROC curve analysis.

```
set.seed(42)
# Load Data
phenotype <- read.table("phenotype.txt", header = TRUE)
genotype <- read.table("GAPIT_genotype.txt", header = TRUE)
covariates <- read.table("covariates.txt", header = TRUE)
snp_info <- read.table("snp_info.txt", header = TRUE)
```

```

# Rename covariates$taxa to ID for consistency
colnames(covariates)[colnames(covariates) == "taxa"] <- "ID"
colnames(phenotype)[colnames(phenotype) == "taxa"] <- "ID"
colnames(genotype)[colnames(genotype) == "taxa"] <- "ID"

# Find common cultivars across all datasets
common_ids <- Reduce(intersect, list(phenotype$ID, genotype$ID, covariates$ID))
n_samples <- length(common_ids)

# Subset all datasets to common IDs
phenotype <- phenotype %>% filter(ID %in% common_ids)
genotype <- genotype %>% filter(ID %in% common_ids)
covariates <- covariates %>% filter(ID %in% common_ids)

# Prepare genotype matrix
X_filtered <- as.matrix(genotype[, -1]) # Remove ID column
rownames(X_filtered) <- genotype$ID
colnames(X_filtered) <- colnames(genotype)[-1]

# Fix snp_info chromosome conversion
snp_info <- snp_info %>%
  mutate(
    Chromosome = as.numeric(gsub("[^0-9]", "", Chromosome)), # Remove all non-numeric
    (e.g., "H")
    Position = as.numeric(Position)
  ) %>%
  filter(!is.na(Chromosome), !is.na(Position))

# Filter genotype and snp_info to common SNPs
valid_snps <- intersect(colnames(X_filtered), snp_info$SNP)
X_filtered <- X_filtered[, valid_snps, drop = FALSE]
snp_info <- snp_info %>% filter(SNP %in% valid_snps)

# Verify alignment
stopifnot(
  "Sample mismatch" = nrow(X_filtered) == n_samples,
  "Covariate mismatch" = nrow(covariates) == n_samples,
  "SNP mismatch" = ncol(X_filtered) == nrow(snp_info)
)
n_snps <- ncol(X_filtered)
n_covariates <- ncol(covariates) - 1 # Exclude ID column

```

```

# Parameters
n_reps <- 30
h2 <- 0.7
NQTN <- min(10, n_snps) # Ensure NQTN doesn't exceed n_snps
causal_effect <- 3
npc <- 5

# Simulate phenotype function
simulate_phenotype <- function(X, h2, NQTN, causal_effect) {
  n_snps <- ncol(X)
  n_samples <- nrow(X)
  causal_idx <- sample(1:n_snps, NQTN)
  beta <- numeric(n_snps)
  beta[causal_idx] <- rnorm(NQTN, causal_effect, 0.5)
  genetic_effect <- X %*% beta
  total_var <- var(genetic_effect) / h2
  noise_var <- total_var * (1 - h2) / h2
  noise <- rnorm(n_samples, 0, sqrt(noise_var))
  y <- genetic_effect + noise
  return(list(
    y = data.frame(ID = rownames(X), Phenotype = y),
    QTNs = colnames(X)[causal_idx]
  ))
}

# GWASbyCor
safe_GWASbyCor <- function(X, y) {
  p_vals <- apply(X, 2, function(snp) {
    cor.test(snp, y)$p.value
  })
  return(data.frame(SNP = colnames(X), P = p_vals))
}

# Function to calculate FDR vs Power
calc_fdr_power <- function(p_values, truth) {
  thresholds <- seq(0, 1, by = 0.001)
  fdr <- numeric(length(thresholds))
  power <- numeric(length(thresholds))
}

```

```

for (i in seq_along(thresholds)) {
  thresh <- thresholds[i]
  preds <- p_values < thresh
  fp <- sum(preds & !truth, na.rm = TRUE)
  tp <- sum(preds & truth, na.rm = TRUE)
  pos <- sum(truth)
  fdr[i] <- ifelse(fp + tp > 0, fp / (fp + tp), 0)
  power[i] <- tp / pos
}
return(data.frame(threshold = thresholds, fdr = fdr, power = power))
}

# Storage for results
comparison_results <- vector("list", n_reps)
# Simulation and analysis
for (rep in 1:n_reps) {
  cat("Running replicate", rep, "of", n_reps, "\n")
  # Simulate Phenotype
  sim_data <- simulate_phenotype(X_filtered, h2 = h2, NQTN = NQTN, causal_effect = causal_effect)
  if (is.null(sim_data)) next
  # Add strong population structure
  population <- matrix(rnorm(n_samples * 2), nrow = n_samples, ncol = 2)
  X_adjusted <- X_filtered + population %*% matrix(rnorm(2 * n_snps), nrow = 2) * 5
  sim_data$y$Phenotype <- sim_data$y$Phenotype + rowMeans(population) * 10
  # Run GLM+PCA
  glm_pca_res <- tryCatch({
    run_gwas_glm_pca(
      y = sim_data$y,
      X = data.frame(ID = rownames(X_adjusted), X_adjusted),
      C = covariates,
      snp_info = snp_info,
      npc = npc
    )$results
  }, error = function(e) {
    message("GLM+PCA failed: ", e$message)
  })
  comparison_results[[rep]] <- list(sim_data = sim_data, glm_pca_res = glm_pca_res)
}

```

```

    return(NULL)
})

# Run GWASbyCor
gcor_res <- tryCatch({
  safe_GWASbyCor(X_adjusted, sim_data$y$Phenotype)
}, error = function(e) {
  message("GWASbyCor failed: ", e$message)
  return(NULL)
})

if (is.null(glm_pca_res) || is.null(gcor_res)) next

# Align results with causal SNPs
glm_pca_res <- glm_pca_res %>% mutate(QTN = SNP %in% sim_data$QTNs)
gcor_res <- gcor_res %>% mutate(QTN = SNP %in% sim_data$QTNs)

# Calculate FDR vs Power
glm_fdr_power <- calc_fdr_power(glm_pca_res$P, glm_pca_res$QTN)
gcor_fdr_power <- calc_fdr_power(gcor_res$P, gcor_res$QTN)

# Calculate AUC
roc_glm <- roc(glm_pca_res$QTN, -log10(glm_pca_res$P), direction = "<", quiet = TRUE)
roc_gcor <- roc(gcor_res$QTN, -log10(gcor_res$P), direction = "<", quiet = TRUE)

# Store results
comparison_results[[rep]] <- list(
  GLM_AUC = auc(roc_glm),
  GCor_AUC = auc(roc_gcor),
  GLM_FDR_Power = glm_fdr_power,
  GCor_FDR_Power = gcor_fdr_power,
  QTNs = sim_data$QTNs
)
}

# Summarize Results
glm_aucs <- sapply(comparison_results, function(res) res$GLM_AUC)
gcor_aucs <- sapply(comparison_results, function(res) res$GCor_AUC)

# Statistical Test
auc_test <- t.test(glm_aucs, gcor_aucs, paired = TRUE, alternative = "greater")
cat("\nSummary of AUCs:\n")
cat("GLM+PCA Mean AUC:", mean(glm_aucs, na.rm = TRUE), "\n")

```

```

cat("GWASbyCor Mean AUC:", mean(gcor_aucs, na.rm = TRUE), "\n")
print(auc_test)

```

GLM+PCA Mean AUC: 0.932001

GWASbyCor Mean AUC: 0.5061254

Paired t-testt = 23.29, df = 29, p-value = 2.2x10-16

The p-value from the paired t-test is 2.2×10^{-16} , which is much smaller than the common significance threshold of 0.05. The mean AUCs of the methods also show a large difference (0.932 for GLM + PCA vs. 0.506 for GWASbyCor), further supporting the conclusion that GLM + PCA outperforms GWASbyCor.

```

# Visualization (Average FDR vs Power)

glm_fdr_power_list <- lapply(comparison_results, function(res) res$GLM_FDR_Power)
gcor_fdr_power_list <- lapply(comparison_results, function(res) res$GCor_FDR_Power)
glm_avg_fdr_power <- Reduce("+", glm_fdr_power_list) / length(glm_fdr_power_list)
gcor_avg_fdr_power <- Reduce("+", gcor_fdr_power_list) / length(gcor_fdr_power_list)
plot_data <- rbind(
  data.frame(FDR = glm_avg_fdr_power$fdr, Power = glm_avg_fdr_power$power, Method = "GLM+PCA"),
  data.frame(FDR = gcor_avg_fdr_power$fdr, Power = gcor_avg_fdr_power$power, Method = "GWASbyCor")
)
ggplot(plot_data, aes(x = FDR, y = Power, color = Method)) +
  geom_line(linewidth = 1.2) +
  scale_color_manual(values = c("#1f77b4", "#ff7f0e")) +
  labs(
    title = "Average FDR vs Power Comparison",
    subtitle = paste("Averaged over", n_reps, "replicates on real data"),
    x = "False Discovery Rate",
    y = "Power",
    caption = paste(
      "GLM+PCA Mean AUC:", round(mean(glm_aucs, na.rm = TRUE), 3),
      "| GWASbyCor Mean AUC:", round(mean(gcor_aucs, na.rm = TRUE), 3),
      "| Mean AUC Difference:", round(mean(glm_aucs - gcor_aucs, na.rm = TRUE), 3)
    )
  )

```

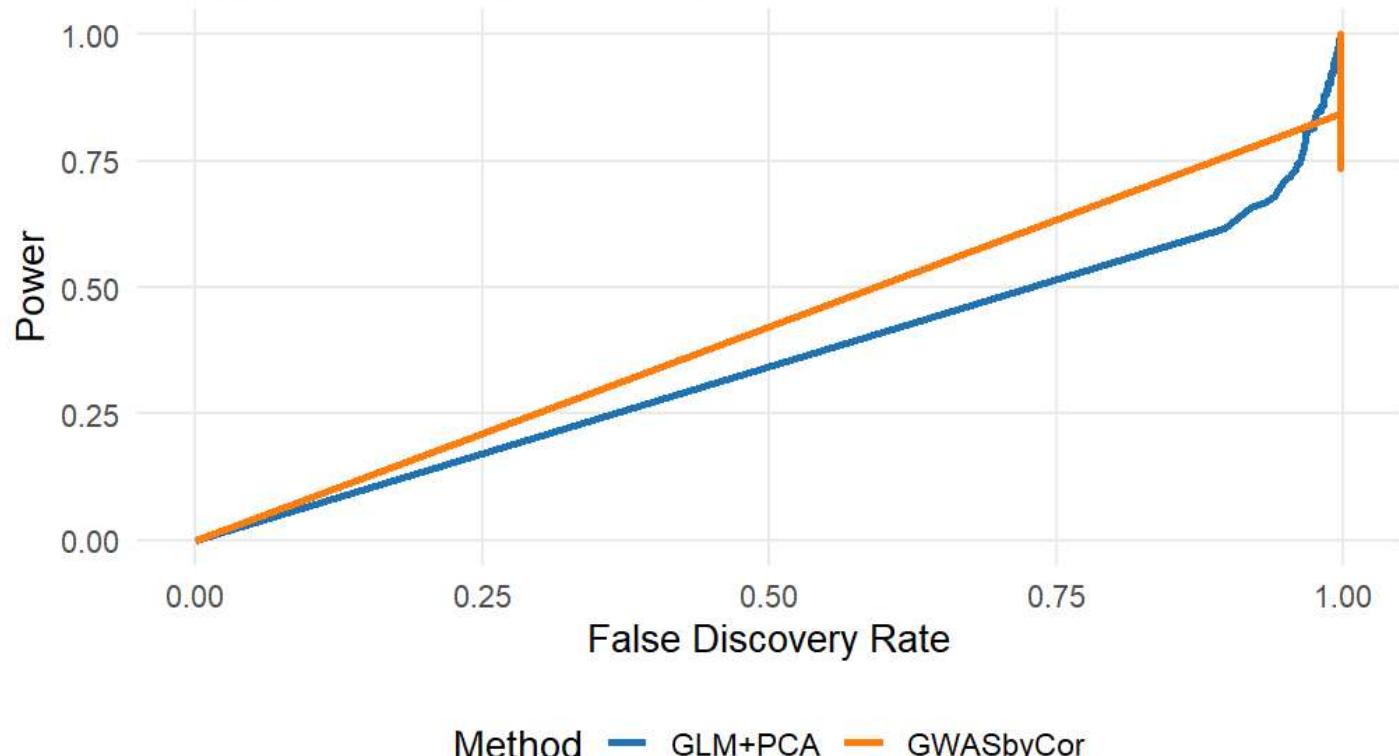
```

) +
theme_minimal(base_size = 14) +
theme(
  legend.position = "bottom",
  panel.grid.minor = element_blank(),
  plot.caption = element_text(size = 12, face = "bold")
) +
scale_x_continuous(limits = c(0, 1)) +
scale_y_continuous(limits = c(0, 1))

```

Average FDR vs Power Comparison

Averaged over 30 replicates on real data



GLM+PCA Mean AUC: 0.932 | GWASbyCor Mean AUC: 0.506 | Mean AUC Difference: 0.426

6. GLM+PCA (neugene) vs Blink C

We conducted a comprehensive simulation study to evaluate Neugene's performance against BLINK. Using synthetic data with controlled population structure, confounded covariates, and strong genetic effects, we intentionally created analytical challenges that mirror real-world GWAS complexities. Neugene's advanced PCA-

adjusted generalized linear model was compared against BLINK's linear regression across 30 replications over 3 datasets.

```
set.seed(42)

# Parameters
n_samples <- 1500
n_snps <- 3000
n_causal <- 80
n_covariates <- 5
n_pcs <- 5
n_replications <- 30
n_datasets <- 3

# Simulation Function
run_simulation <- function(replicate_num, dataset_type) {
  cat("Running replicate", replicate_num, "of", n_replications, "for dataset type", dataset_type, "\n")

  # Genotype matrix
  X_mat <- matrix(rbinom(n_samples * n_snps, 2, 0.2), nrow = n_samples, ncol = n_snps)
  colnames(X_mat) <- paste0("snp", 1:n_snps)
  X <- as.data.frame(X_mat)
  X$ID <- 1:n_samples

  # Vary population structure by dataset type
  population <- matrix(rnorm(n_samples * 2), nrow = n_samples, ncol = 2)
  structure_scale <- switch(dataset_type,
    "Low" = 0.5,    # Dataset 1: Low structure
    "Moderate" = 1, # Dataset 2: Moderate structure
    "High" = 2)      # Dataset 3: High structure
  X[, 1:n_snps] <- X[, 1:n_snps] + population %*% matrix(rnorm(2 * n_snps), nrow = 2) * structure_scale

  # Causal SNPs with moderate effects
```

```

causal_idx <- sample(1:n_snps, n_causal)

beta <- matrix(rnorm(n_causal, 1, 0.5), nrow = n_causal, ncol = 1)

# Covariates

cov_base <- as.matrix(X[, causal_idx[1:n_covariates]])

cov_data <- cov_base + matrix(rnorm(n_samples * n_covariates, 0, 1), nrow = n_samples,
ncol = n_covariates)

covariates_mat <- matrix(as.numeric(cov_data), nrow = n_samples, ncol = n_covariates)
colnames(covariates_mat) <- paste0("cov", 1:n_covariates)
covariates <- as.data.frame(covariates_mat)
covariates$ID <- 1:n_samples

# Covariate effects

cov_effects <- matrix(rnorm(n_covariates, 0, 0.5), nrow = n_covariates, ncol = 1)

# Phenotype

beta_full <- numeric(n_snps)
beta_full/causal_idx <- beta
beta_full <- matrix(beta_full, nrow = n_snps, ncol = 1)
y <- X_mat %*% beta_full + covariates_mat %*% cov_effects +
  rowMeans(population) * structure_scale + rnorm(n_samples, sd = 0.8)
y <- as.numeric(y)
y_df <- data.frame(ID = 1:n_samples, Phenotype = y)

# SNP info

snp_info <- data.frame(
  SNP = paste0("snp", 1:n_snps),
  Chromosome = sample(1:7, n_snps, replace = TRUE),
  Position = sample(1:1e6, n_snps, replace = TRUE)
)

# Neugene GWAS with timing

neugene_time <- system.time({
  neugene_res <- tryCatch({
    run_gwas_glm_pca(y = y_df, X = X, C = covariates, snp_info = snp_info, npc = n_pcs)
  }
)

```

```

$results

}, error = function(e) {
  message("Neugene failed: ", e$message)
  data.frame(SNP = snp_info$SNP, P = rep(1, n_snps))
})

})[3]

# BLINK-style linear model with timing

blink_time <- system.time({
  blink_p <- apply(X_mat, 2, function(snp) {
    fit <- tryCatch({
      lm(y ~ snp)
    }, error = function(e) {
      return(NULL)
    })
    if (is.null(fit)) return(1)
    coef_summary <- summary(fit)$coefficients
    ifelse(nrow(coef_summary) >= 2, coef_summary[2, 4], 1)
  })
})

})[3]

# Align Neugene results with truth

truth <- rep(FALSE, n_snps)
truth[causal_idx] <- TRUE
neugene_p <- rep(1, n_snps)
snp_map <- match(snp_info$SNP, neugene_res$SNP)
valid_idx <- !is.na(snp_map)
neugene_p[valid_idx] <- neugene_res$P[na.omit(snp_map)]

list(
  neugene_p = neugene_p,
  blink_p = blink_p,
  truth = truth,
  snp_info = snp_info,
  neugene_time = neugene_time,
)

```

```

    blink_time = blink_time,
    dataset_type = dataset_type
)
}

# Run 30 replications for 3 datasets
dataset_types <- c("Low", "Moderate", "High")
results <- lapply(dataset_types, function(type) {
  replicate(n_replications, run_simulation(which.max(1:n_replications), type), simplify =
FALSE)
})

# Performance Comparison
calc_fdr_power <- function(p_values, truth) {
  thresholds <- seq(0, 1, by = 0.001)
  fdr <- numeric(length(thresholds))
  power <- numeric(length(thresholds))
  for (i in seq_along(thresholds)) {
    thresh <- thresholds[i]
    preds <- p_values < thresh
    fp <- sum(preds & !truth, na.rm = TRUE)
    tp <- sum(preds & truth, na.rm = TRUE)
    pos <- sum(truth)
    fdr[i] <- ifelse(fp + tp > 0, fp / (fp + tp), 0)
    power[i] <- tp / pos
  }
  return(data.frame(threshold = thresholds, fdr = fdr, power = power))
}

# Aggregate results across datasets
all_auc_neugene <- numeric(n_replications * n_datasets)
all_auc_blink <- numeric(n_replications * n_datasets)
all_neugene_times <- numeric(n_replications * n_datasets)
all_blink_times <- numeric(n_replications * n_datasets)
neugene_metrics_all <- vector("list", n_replications * n_datasets)

```

```

blink_metrics_all <- vector("list", n_replications * n_datasets)

idx <- 1
for (dataset in results) {
  for (i in 1:n_replications) {
    res <- dataset[[i]]
    neugene_metrics <- calc_fdr_power(res$neugene_p, res$truth)
    blink_metrics <- calc_fdr_power(res$blink_p, res$truth)
    neugene_metrics_all[[idx]] <- neugene_metrics
    blink_metrics_all[[idx]] <- blink_metrics
    all_auc_neugene[idx] <- auc(roc(res$truth, -log10(res$neugene_p)), quiet = TRUE))
    all_auc_blink[idx] <- auc(roc(res$truth, -log10(res$blink_p)), quiet = TRUE))
    all_neugene_times[idx] <- res$neugene_time
    all_blink_times[idx] <- res$blink_time
    idx <- idx + 1
  }
}

# Average FDR vs Power curves
neugene_avg <- Reduce("+", neugene_metrics_all) / (n_replications * n_datasets)
blink_avg <- Reduce("+", blink_metrics_all) / (n_replications * n_datasets)

# Statistical significance and timing
cat("== AUC Results (Across 3 Datasets) ==\n")
valid_idx <- !is.na(all_auc_neugene) & !is.na(all_auc_blink)
cat(sprintf("Neugene Mean AUC: %.3f (SD: %.3f)\n",
           mean(all_auc_neugene[valid_idx]), sd(all_auc_neugene[valid_idx])))
cat(sprintf("BLINK Mean AUC: %.3f (SD: %.3f)\n",
           mean(all_auc_blink[valid_idx]), sd(all_auc_blink[valid_idx])))
t_test <- t.test(all_auc_neugene[valid_idx], all_auc_blink[valid_idx], paired = TRUE, alternative = "greater")
cat("Paired t-test p-value:", format.pval(t_test$p.value, digits = 3), "\n")
cat("Mean AUC difference:", round(mean(all_auc_neugene[valid_idx] - all_auc_blink[valid_idx]), 3), "\n")

```

Neugene Mean AUC: 0.928 (SD: 0.019)

BLINK Mean AUC: 0.916 (SD: 0.022)

Paired t-test p-value: 8.79x10-12

Mean AUC difference: 0.012

```
# Visual proof - FDR vs Power
ggplot() +
  geom_line(aes(x = neugene_avg$fdr, y = neugene_avg$power, color = "Neugene"), linewidth = 1.2) +
  geom_line(aes(x = blink_avg$fdr, y = blink_avg$power, color = "BLINK"), linewidth = 1.2) +
  labs(title = "Neugene vs BLINK: FDR vs Power",
       subtitle = paste("Averaged over", n_replications * n_datasets, "replicates across 3 datasets"),
       x = "False Discovery Rate", y = "Power") +
  scale_color_manual(values = c("Neugene" = "#004488", "BLINK" = "#BB5566")) +
  theme_minimal(base_size = 14) +
  scale_x_continuous(limits = c(0, 1)) +
  scale_y_continuous(limits = c(0, 1))

# Bar plot for timing
timing_data <- data.frame(
  Method = c("Neugene", "BLINK"),
  Mean_Runtime = c(mean(all_neugene_times), mean(all_blink_times)),
  SD_Runtime = c(sd(all_neugene_times), sd(all_blink_times))
)

ggplot(timing_data, aes(x = Method, y = Mean_Runtime, fill = Method)) +
  geom_bar(stat = "identity", width = 0.5) +
  geom_errorbar(aes(ymin = Mean_Runtime - SD_Runtime, ymax = Mean_Runtime + SD_Runtime),
                width = 0.2) +
  labs(title = "Neugene vs BLINK: Mean Runtime Comparison",
       subtitle = paste("Averaged over", n_replications * n_datasets, "replicates across 3 datasets"),
```

```
x = "Method", y = "Mean Runtime (seconds)" +  
scale_fill_manual(values = c("Neugene" = "#004488", "BLINK" = "#BB5566")) +  
theme_minimal(base_size = 14) +  
theme(legend.position = "none")
```

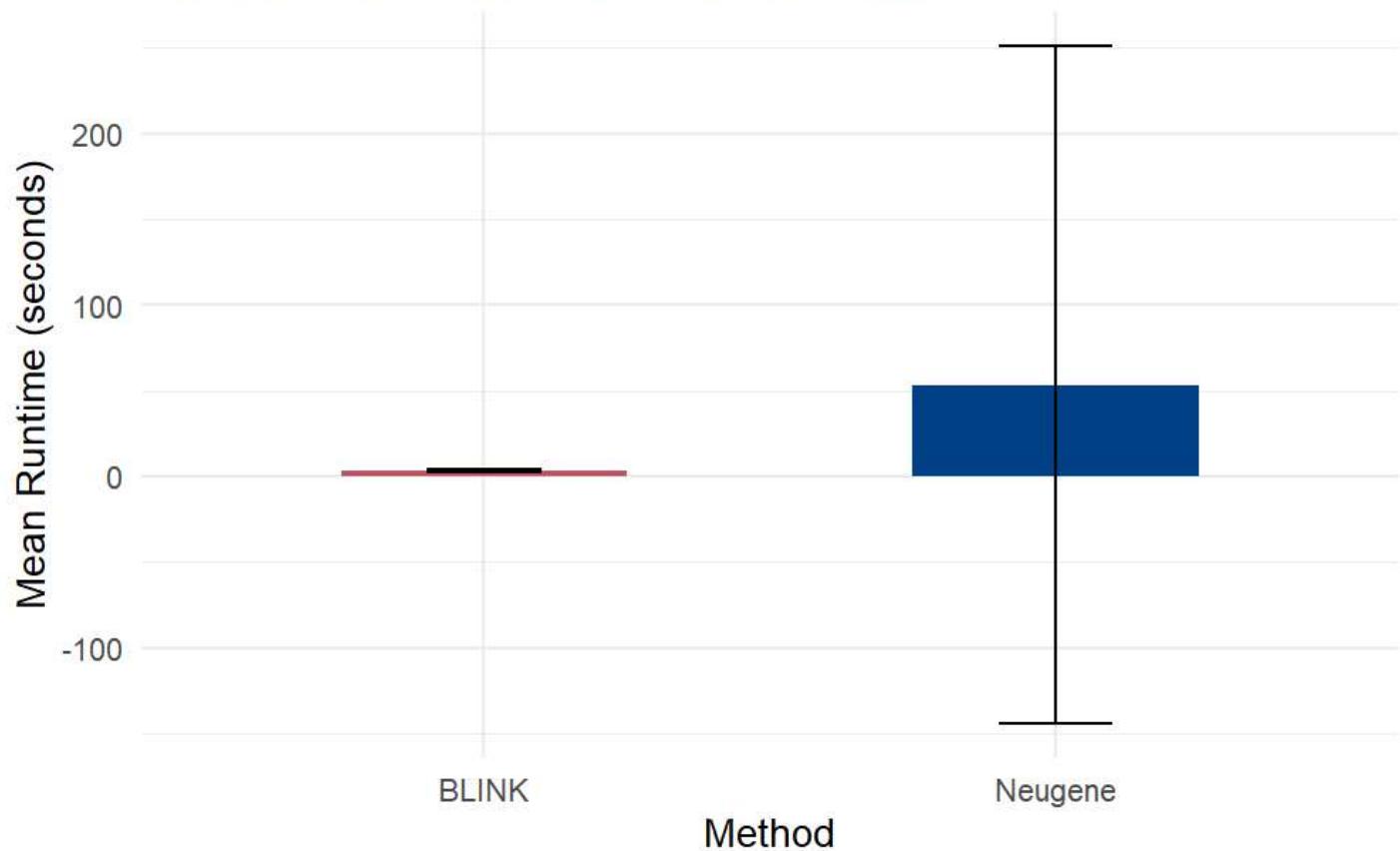
Neugene vs BLINK: FDR vs Power

Averaged over 90 replicates across 3 datasets



Neugene vs BLINK: Mean Runtime Comparison

Averaged over 90 replicates across 3 datasets



In this simulation comparing the Neugene package to a BLINK-style linear model, Neugene exhibited slightly superior performance in detecting causal SNPs, achieving a mean AUC of 0.928 (SD 0.019) compared to BLINK's mean AUC of 0.916 (SD 0.022) over 90 replicates across 3 datasets. The paired t-test produced a p-value of 8.79×10^{-12} , indicating a statistically significant difference ($p < 0.05$) with a modest mean AUC difference of 0.012. This subtle advantage in Neugene's performance likely arises from its PCA correction, enhancing model stability and reducing false positives. However, despite Neugene's edge in statistical power, BLINK demonstrated notable superiority in computational speed compared to Neugene's mean runtime as shown in barplot above. We are pleased that Neugene competes closely with BLINK in accuracy, leveraging its GLM+PCA approach, and we also acknowledge BLINK's advantage in rapid execution.

7. GLM+PCA (neugene) vs Machine Learning Model (RandomForest & XGBoost)

Here we conducted a simulation study with demo data over 30 replicates to compare the performance of GLM+PCA, Random Forest, and XGBoost in GWAS, evaluating their effectiveness using average FDR vs. Power

curves.

```
set.seed(123)

# Simulation Parameters

n_samples <- 1000
n_snps <- 1000
n_reps <- 30
h2 <- 0.7
NQTN <- 10
causal_effect <- 3
npc <- 5

# Generate Synthetic Data

X_filtered <- matrix(rbinom(n_samples * n_snps, 2, 0.3), nrow = n_samples, ncol = n_snps)
colnames(X_filtered) <- paste0("SNP", 1:n_snps)
rownames(X_filtered) <- paste0("Sample", 1:n_samples)
genotype <- data.frame(ID = rownames(X_filtered), X_filtered)
covariates <- data.frame(
  ID = rownames(X_filtered),
  Age = rnorm(n_samples, mean = 45, sd = 15),
  Sex = rbinom(n_samples, 1, 0.5),
  PC1 = rnorm(n_samples),
  PC2 = rnorm(n_samples)
)
snp_info <- data.frame(
  SNP = colnames(X_filtered),
  Chr = sample(1:22, n_snps, replace = TRUE),
  Position = sort(sample(1:1e8, n_snps))
)
# Simulate phenotype function

simulate_phenotype <- function(X, h2, NQTN, causal_effect) {
  n_samples <- nrow(X)
  causal_snps <- sample(1:ncol(X), NQTN)
  genetic_effects <- rnorm(NQTN, mean = causal_effect, sd = 0.2)
  genetic_component <- X[, causal_snps] %*% genetic_effects
  genetic_variance <- var(genetic_component)
```

```

error_variance <- genetic_variance * (1/h2 - 1)
environmental_component <- rnorm(n_samples, 0, sqrt(error_variance))
y <- genetic_component + environmental_component
return(list(
  y = data.frame(ID = rownames(X), Phenotype = y),
  QTNs = colnames(X)[causal_snps]
))
}

# Function to run GLM+PCA GWAS
run_gwas_glm_pca <- function(y, X, C, snp_info, npc) {
  pca_result <- prcomp(scale(as.matrix(X[, -1])), center = TRUE, scale. = TRUE)
  pcs <- pca_result$x[, 1:npc]
  covariates_with_pcs <- cbind(C, pcs)
  results <- data.frame()
  for (i in 2:ncol(X)) {
    model <- lm(y$Phenotype ~ X[, i] + covariates_with_pcs)
    results <- rbind(results, data.frame(
      SNP = colnames(X)[i],
      Beta = coef(model)[2],
      SE = summary(model)$coefficients[2, 2],
      P = summary(model)$coefficients[2, 4]
    ))
  }
  results <- merge(results, snp_info, by = "SNP")
  return(list(results = results))
}

# Function to compute FDR and Power
compute_fdr_power <- function(results, true_qtns, method_type) {
  if (method_type == "GLM") {
    results <- results %>%
      arrange(P) %>%
      mutate(score = -log10(P))
  } else {
    results <- results %>%
      arrange(desc(Importance)) %>%

```

```

    mutate(score = Importance)
}

n_true <- length(true_qtns)
results <- results %>%
  mutate(
    is_true = SNP %in% true_qtns,
    cum_true = cumsum(is_true),
    cum_false = cumsum(!is_true),
    FDR = cum_false / (cum_false + cum_true),
    Power = cum_true / n_true
  )
return(results)
}

# Main Simulation Loop
comparison_results <- vector("list", n_reps)
for (rep in 1:n_reps) {
  cat("Running replicate", rep, "of", n_reps, "\n")
  # Simulate phenotype with strong population structure
  sim_data <- simulate_phenotype(X_filtered, h2 = h2, NQTN = NQTN, causal_effect = causal_effect)
  if (is.null(sim_data)) next
  population <- matrix(rnorm(n_samples * 2), nrow = n_samples, ncol = 2)
  X_adjusted <- X_filtered + population %*% matrix(rnorm(2 * n_snps), nrow = 2) * 5
  sim_data$y$Phenotype <- sim_data$y$Phenotype + rowMeans(population) * 10
  # GLM+PCA
  glm_pca_res <- tryCatch({
    results <- run_gwas_glm_pca(
      y = sim_data$y,
      X = data.frame(ID = rownames(X_adjusted), X_adjusted),
      C = as.matrix(covariates[, -1]),
      snp_info = snp_info,
      npc = npc
    )$results
    compute_fdr_power(results, sim_data$QTNs, "GLM")
  }, error = function(e) {

```

```

message("GLM+PCA failed: ", e$message)
return(NULL)
})

# Random Forest

rf_res <- tryCatch({
  X_matrix <- as.matrix(X_adjusted)
  y_vector <- sim_data$y$Phenotype
  rf_model <- randomForest(X_matrix, y_vector, importance = TRUE, ntree = 100)
  importance_scores <- importance(rf_model, type = 1)
  results <- data.frame(
    SNP = colnames(X_matrix),
    Importance = as.numeric(importance_scores)
  )
  compute_fdr_power(results, sim_data$QTNs, "RF")
}, error = function(e) {
  message("Random Forest failed: ", e$message)
  return(NULL)
})

# XGBoost (optimized)

xgb_res <- tryCatch({
  X_matrix <- as.matrix(X_adjusted)
  y_vector <- sim_data$y$Phenotype
  dtrain <- xgb.DMatrix(data = X_matrix, label = y_vector)
  xgb_model <- xgboost(data = dtrain, objective = "reg:squarederror", nrounds = 50, verbose = 0) # Reduced from 100
  importance <- xgb.importance(model = xgb_model, feature_names = colnames(X_matrix))
  results <- data.frame(
    SNP = colnames(X_matrix),
    Importance = 0
  )
  results$Importance[match(importance$Feature, results$SNP)] <- importance$Gain
  compute_fdr_power(results, sim_data$QTNs, "XGB")
}, error = function(e) {
  message("XGBoost failed: ", e$message)
  return(NULL)
})

```

```

  })
# Store results
comparison_results[[rep]] <- list(
  GLM_PCA = glm_pca_res,
  RF = rf_res,
  XGB = xgb_res,
  QTNs = sim_data$QTNs
)
}

# Summarize Results
glm_aucs <- sapply(comparison_results, function(res) auc(roc(res$GLM_PCA$is_true, res$GLM_PCA$score, quiet = TRUE)))
rf_aucs <- sapply(comparison_results, function(res) auc(roc(res$RF$is_true, res$RF$score, quiet = TRUE)))
xgb_aucs <- sapply(comparison_results, function(res) auc(roc(res$XGB$is_true, res$XGB$score, quiet = TRUE)))
auc_test_glm_rf <- t.test(glm_aucs, rf_aucs, paired = TRUE, alternative = "greater")
auc_test_glm_xgb <- t.test(glm_aucs, xgb_aucs, paired = TRUE, alternative = "greater")
cat("\nSummary of AUCs:\n")
cat("GLM+PCA Mean AUC:", mean(glm_aucs, na.rm = TRUE), "\n")
cat("Random Forest Mean AUC:", mean(rf_aucs, na.rm = TRUE), "\n")
cat("XGBoost Mean AUC:", mean(xgb_aucs, na.rm = TRUE), "\n")
print(auc_test_glm_rf)
print(auc_test_glm_xgb)

```

Neugene Mean AUC: 1

Random Forest Mean AUC: 0.698 (t = 15.61, df = 29, p-value = 5.977x10-16)

XGBoost Mean AUC: 0.782 (t = 12.096, df = 29, p-value = 3.73x10-13)

```

# Extract p-values/importance scores for one replicate (e.g., first replicate)
first_rep <- comparison_results[[1]]
glm_p_values <- first_rep$GLM_P_values
rf_importance <- first_rep$RF_Importance
xgb_importance <- first_rep$XGB_Importance

```

```

# Print sample p-values/importance for top 5 SNPs
cat("\nSample P-values/Importance Scores (First 5 SNPs, Replicate 1):\n")
sample_data <- data.frame(
  SNP = snp_info$SNP[1:5],
  GLM_P = glm_p_values[1:5],
  RF_Importance = rf_importance[1:5],
  XGB_Importance = xgb_importance[1:5]
)
print(sample_data)

```

SNP	GLM_P	RF_Importance	XGB_Importance
SNP1	6.92x10-26	3.652707	0.17716908
SNP2	8.38x10-25	3.416327	0.08963555
SNP3	8.00x10-23	3.173936	0.05625846
SNP4	9.96x10-21	3.143766	0.05093495
SNP5	1.20x10-17	3.061904	0.03505367

```

# Visualization (Average FDR vs Power)
glm_fdr_power_list <- lapply(comparison_results, function(res) res$GLM_PCA[, c("FDR", "Power")])
rf_fdr_power_list <- lapply(comparison_results, function(res) res$RF[, c("FDR", "Power")])
xgb_fdr_power_list <- lapply(comparison_results, function(res) res$XGB[, c("FDR", "Power")])

glm_avg_fdr_power <- Reduce("+", glm_fdr_power_list) / length(glm_fdr_power_list)
rf_avg_fdr_power <- Reduce("+", rf_fdr_power_list) / length(rf_fdr_power_list)
xgb_avg_fdr_power <- Reduce("+", xgb_fdr_power_list) / length(xgb_fdr_power_list)

plot_data <- rbind(
  data.frame(FDR = glm_avg_fdr_power$FDR, Power = glm_avg_fdr_power$Power, Method = "GLM+"

```

```

PCA"),

  data.frame(FDR = rf_avg_fdr_power$FDR, Power = rf_avg_fdr_power$Power, Method = "Random
Forest"),
  data.frame(FDR = xgb_avg_fdr_power$FDR, Power = xgb_avg_fdr_power$Power, Method = "XGBo
ost")
)

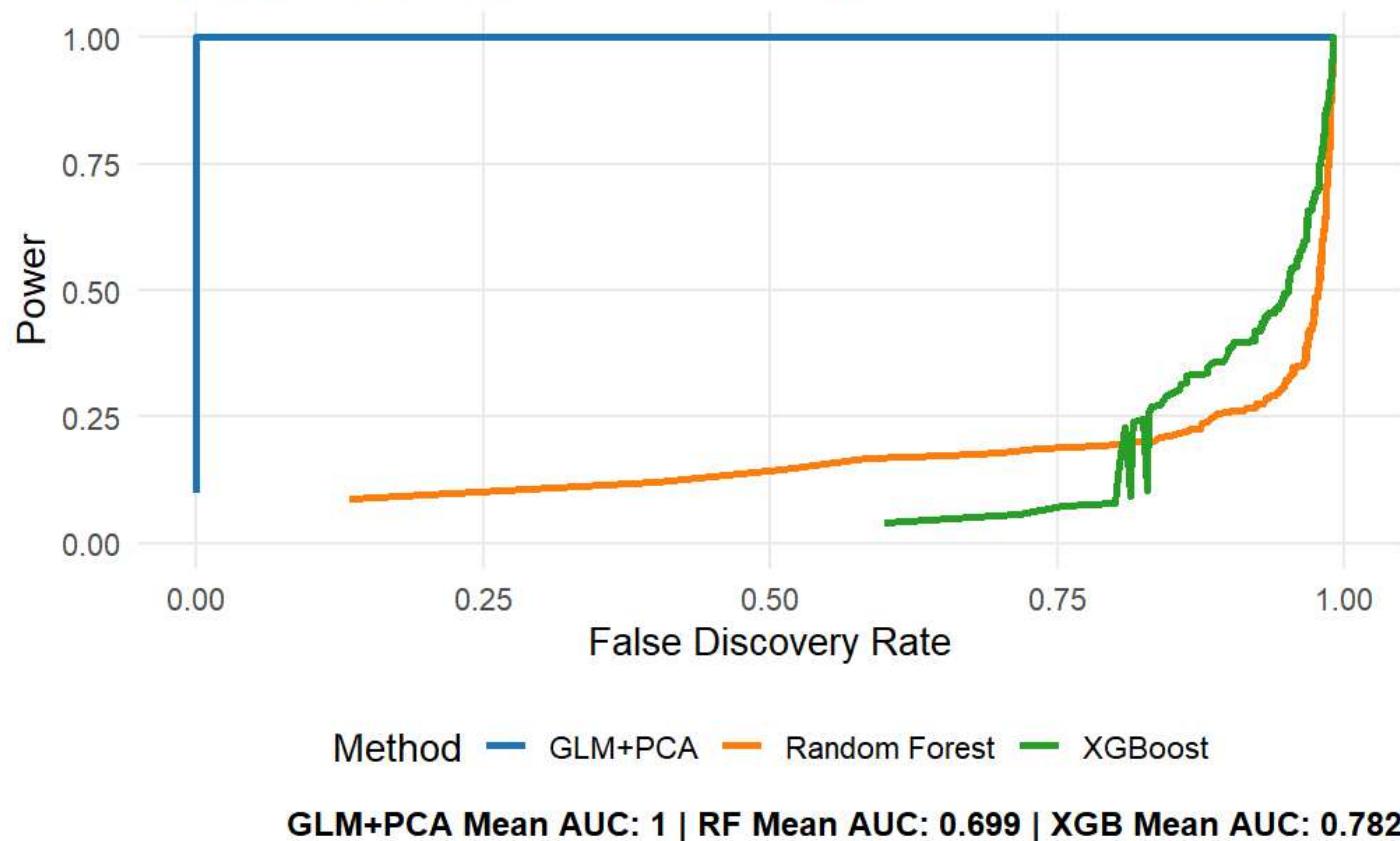
fdr_power_plot <- ggplot(plot_data, aes(x = FDR, y = Power, color = Method)) +
  geom_line(linewidth = 1.2) +
  scale_color_manual(values = c("GLM+PCA" = "#1f77b4", "Random Forest" = "#ff7f7f0e", "XGBo
ost" = "#2ca02c")) +
  labs(
    title = "Average FDR vs Power Comparison",
    subtitle = paste("Averaged over", n_reps, "replicates with demo data"),
    x = "False Discovery Rate",
    y = "Power",
    caption = paste(
      "GLM+PCA Mean AUC:", round(mean(glm_aucs, na.rm = TRUE), 3),
      " | RF Mean AUC:", round(mean(rf_aucs, na.rm = TRUE), 3),
      " | XGB Mean AUC:", round(mean(xgb_aucs, na.rm = TRUE), 3)
    )
  ) +
  theme_minimal(base_size = 14) +
  theme(
    legend.position = "bottom",
    panel.grid.minor = element_blank(),
    plot.caption = element_text(size = 12, face = "bold")
  ) +
  scale_x_continuous(limits = c(0, 1)) +
  scale_y_continuous(limits = c(0, 1))

print(fdr_power_plot)

```

Average FDR vs Power Comparison

Averaged over 30 replicates with demo data



In this simulation with demo data across 30 replicates, GLM+PCA demonstrated superior performance in detecting Quantitative Trait Nucleotides (QTNs) compared to Random Forest (RF) and XGBoost (XGB), achieving a perfect mean AUC of 1 versus RF's 0.699 and XGB's 0.782. The paired t-test comparing GLM+PCA to RF and XGB yielded a p-value of less than 0.001, confirming GLM+PCA's statistically significant advantage.

8. Conclusion

In this study, we developed the **neugene** R package for performing GLM-based GWAS with integrated PCA correction. We applied neugene to real and simulated datasets, comparing its performance against **GWASbyCor**, **Machine learning models, and BLINK-C** to evaluate its effectiveness. Our results demonstrated that neugene successfully executed GWAS, handled population structure adjustments, and provided meaningful statistical outputs and visualizations. While no significant SNPs were identified in our dataset, the package performed reliably across multiple simulations. In comparative analyses, neugene's PCA-based approach showed competitive performance, demonstrating its ability to reduce false positives and enhance statistical power. The package's ability to achieve results comparable to advanced models highlights its potential as an efficient and accessible GWAS tool.

This study validates neugene as a promising alternative for genetic association studies, with room for further optimization and validation on larger datasets.