

20 Kasım 2023

Açıklanabilir Yapay Zeka

7. Hafta: Global düzeyde açıklayıcılar | Değişken önemi yöntemi

Mustafa Cavus, Ph.D.

 Eskişehir Teknik Üniversitesi - İstatistik Bölümü

 mustafacavus@eskisehir.edu.tr

 linktr.ee/mustafacavus

Giriş

Değişken önemliliğinin kullanımına ihtiyaç duyulan durumlar:

- **Model basitleştirme:** Bir modelin tahminlerini etkilemeyen değişkenler modelin dışında tutulabilir.
- **Model araştırması:** Değişkenlerin farklı modellerdeki öneminin karşılaştırılması, değişkenler arasındaki ilişkilerin keşfedilmesine yardımcı olabilir.
- **Alan bilgisine dayalı model doğrulama:** En önemli değişkenlerin belirlenmesi, alan bilgisine dayalı modelin geçerliliğinin değerlendirilmesinde yardımcı olabilir.
- **Bilgi üretimi:** En önemli değişkenlerin tanımlanması, belirli bir mekanizmaya dahil olan yeni faktörlerin keşfedilmesine yol açabilir.

Giriş

Doğrusal modeller ve diğer birçok model türü için, modelin yapısının belirli unsurlarından yararlanan açıklayıcı değişkenin önemini değerlendirmeye yönelik modele özgü yöntemler bulunmaktadır.

Örneğin;

- doğrusal modeller için normalleştirilmiş regresyon katsayısının değeri veya buna karşılık gelen p değeri,
- ağaç tabanlı modeller için değişkenlerin dallanmalarda kullanım sayıları olabilir.

Permütasyonel Değişken Önemliliği

Permütasyonel Değişken Önemliliği

Fisher, Rudin ve Dominici (2019) seçilen bir açıklayıcı değişkenin veya bir grup değişkenin etkisi ortadan kaldırıldığında modelin performansının ne kadar değişeceğini ölçmeye dayalı modelden bağımsız bir değişken önemliliği yaklaşımı geliştirmişlerdir.

Etkiyi ortadan kaldırmak için ampirik bir dağılımdan yeniden örnekleme veya değişkenin değerlerinin permütasyonu gibi pertürbasyonlar kullanılabilir.

Fisher, Aaron, Cynthia Rudin, and Francesca Dominici. 2019. “All Models Are Wrong, but Many Are Useful: Learning a Variable’s Importance by Studying an Entire Class of Prediction Models Simultaneously.” *Journal of Machine Learning Research* 20 (177): 1–81. <http://jmlr.org/papers/v20/18-760.html>.

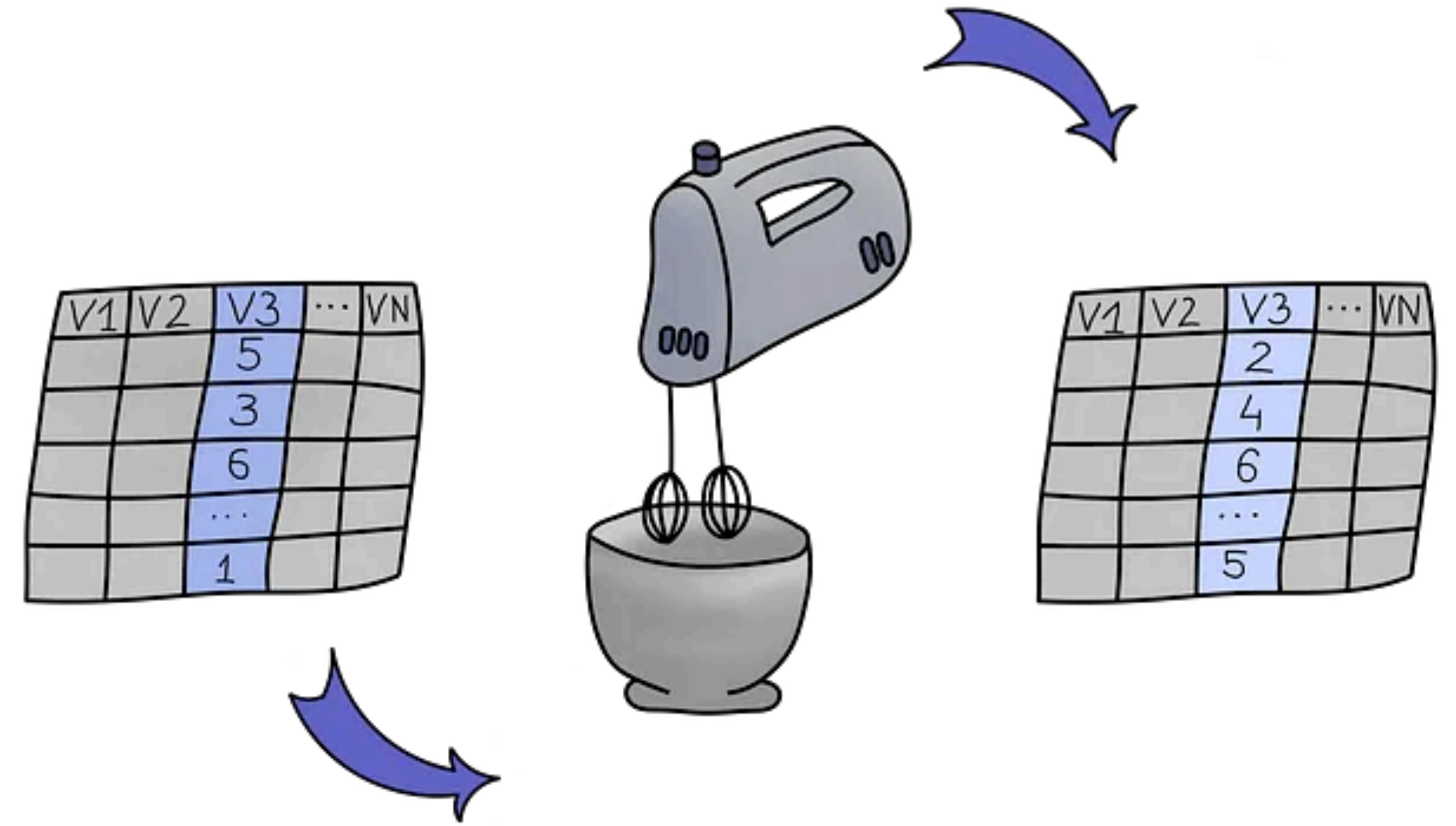
Yöntem

Bu fikir, Leo Breiman (2001) tarafından rastgele orman için önerilen değişken önem ölçüsünden alınmıştır. Eğer bir değişken önemliyse, o zaman değişkenin değerlerine değişiklik yapıldıktan sonra modelin performansının kötüleşmesini bekleriz. Performanstaki değişiklik ne kadar büyükse değişken de o kadar önemlidir.

Breiman, Leo. 2001. “Random Forests.” *Machine Learning* 45: 5–32. <https://doi.org/10.1023/a:1010933404324>.

Permütasyonel Değişken Önemliliği

Etkiyi ortadan kaldırmak için ampirik bir dağılımdan yeniden örnekleme veya değişkenin değerlerinin permütasyonu gibi pertürbasyonlar kullanılabilir.



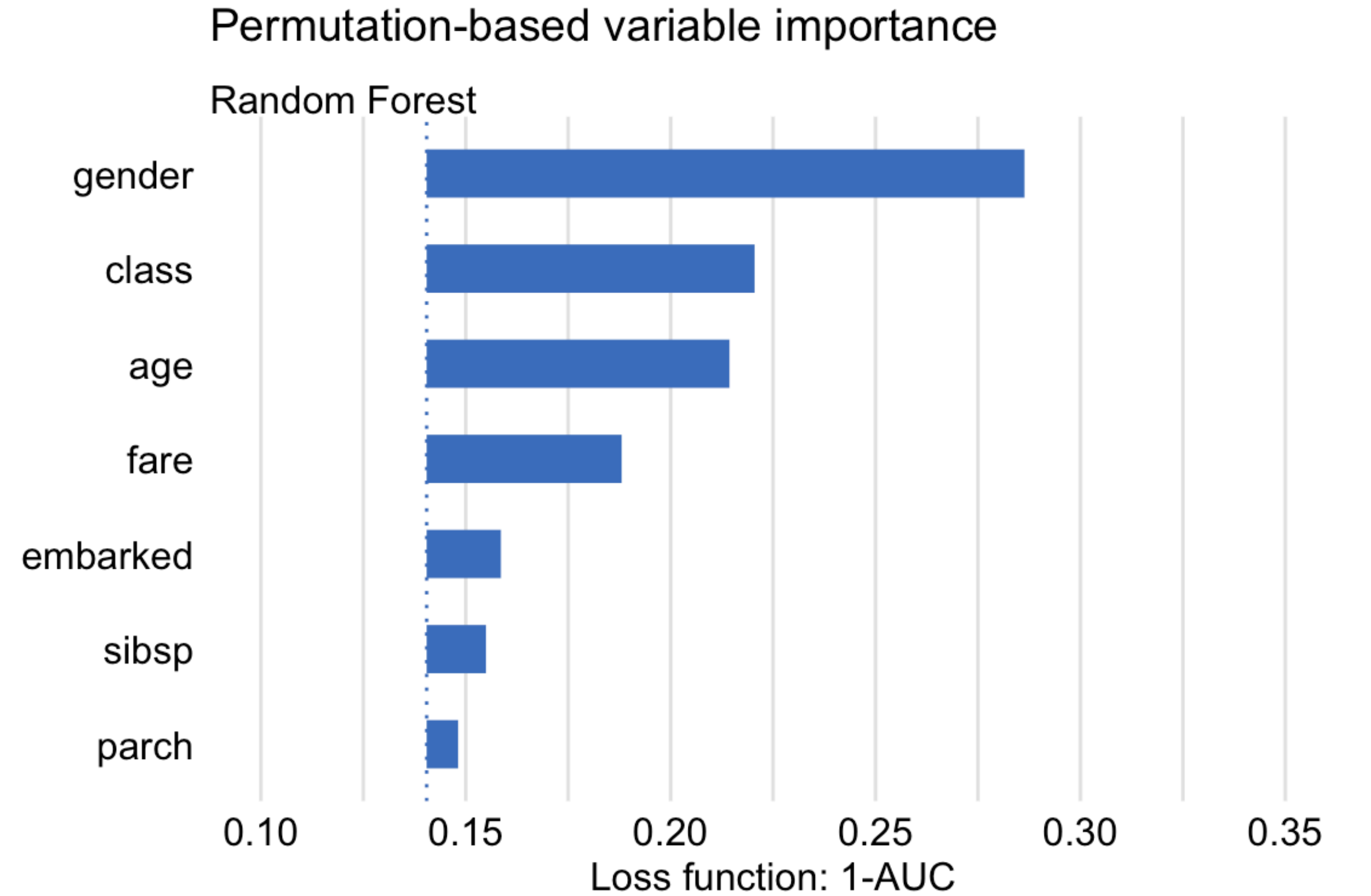
Kaynak: <https://medium.com/responsibleml/basic-xai-with-dalex-part-2-permutation-based-variable-importance-1516c2924a14>

Uygulama

Uygulama

Modeldeki en önemli değişkenin **gender** olduğunu görülmektedir.

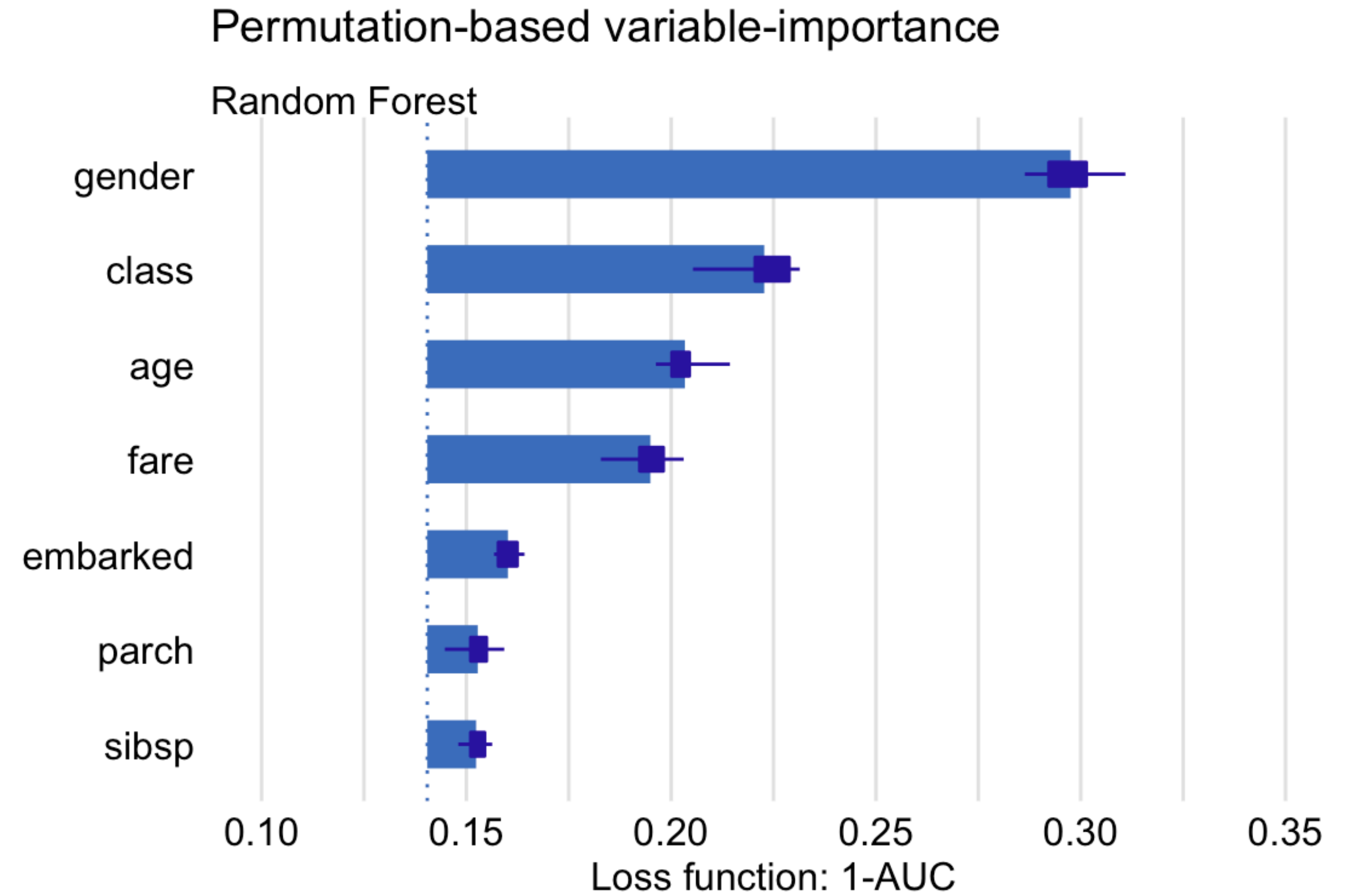
Sonraki en önemli üç değişken **class** (birinci sınıfta seyahat eden yolcuların hayatta kalma şansının daha yüksek olduğu), **age** (çocukların hayatta kalma şansının daha yüksek olduğu) ve **fare** (daha pahalı bilet sahiplerinin hayatta kalma şansı daha yüksektir).



Kayıp fonksiyonu olarak 1-AUC kullanılarak Titanic verileri için **Random Forest** modelinde yer alan açıklayıcı değişkenler için tek permütasyona dayalı değişken önem ölçümleri.

Uygulama

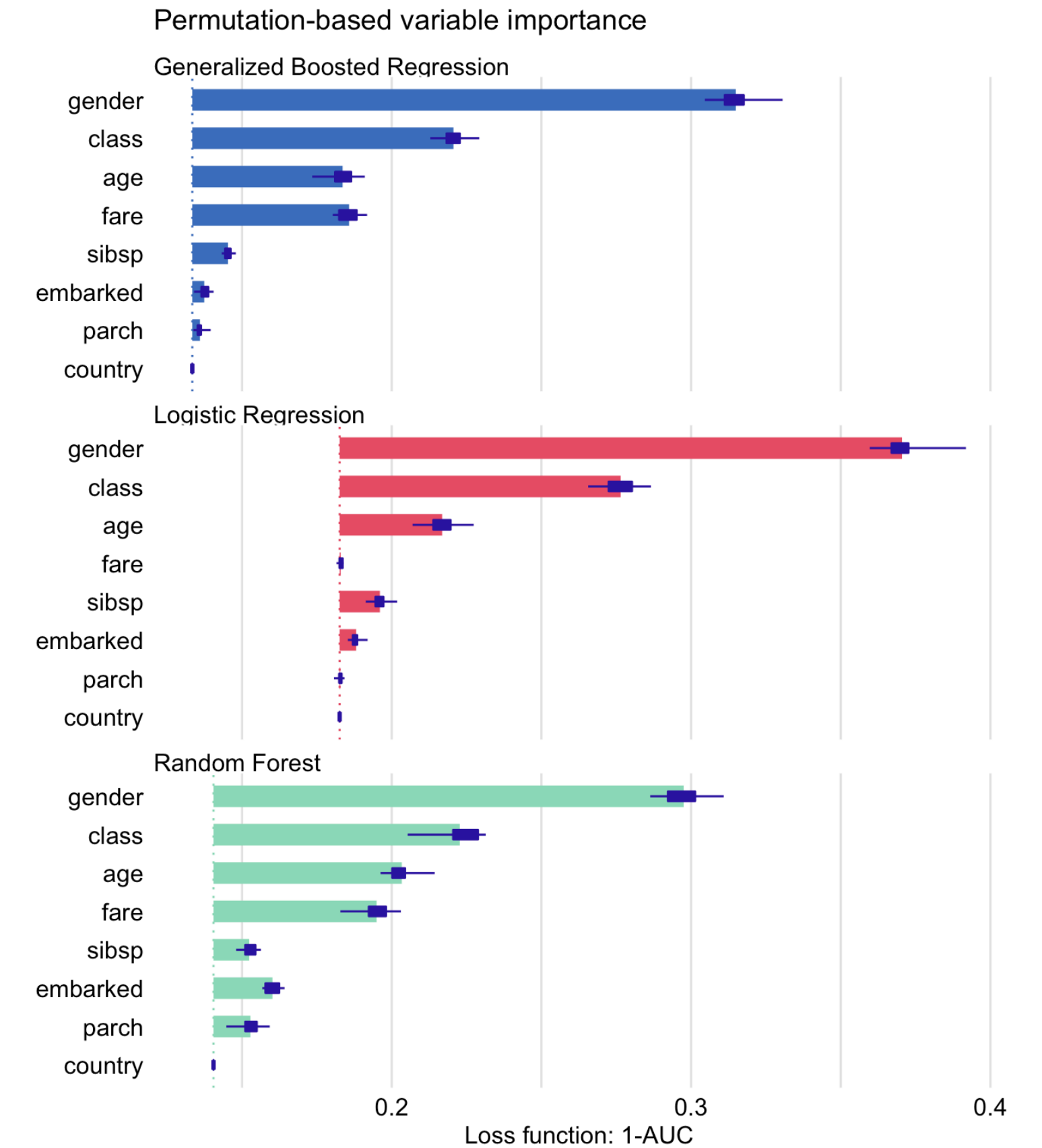
Bir önceki grafik ile karşılaştırıldığında göze çarpan tek fark, **sibsp** ve **parch** değişkenlerinin sıralamasındaki değişiklidir.



Kayıp fonksiyonu olarak 1-AUC kullanan Titanic verileri için **Random Forest** modelinde yer alan açıklayıcı değişkenler için 10 permütasyona dayalı değişken önem ölçümlerinin ortalamaları

Uygulama

- Üç modelde de **gender** en önemli açıklayıcı değişken olduğunu, ardından **class** ve **age** geldiği görülmektedir.
- **class** ile yüksek korelasyona sahip olan **fare** değişkeni, rastgele orman ve SVM modellerinde önemlidir ancak lojistik regresyon modelinde önemli değildir.
- parch aslında ne GBR ne de lojistik regresyon modelinde çok önemli değildir, ancak rastgele orman modelinde kayda değer bir öneme sahiptir.
- Hiçbir modelde **country** önemli değildir.
- Genel olarak rastgele orman modelinde tüm değişkenlerin (**country** hariç) bir miktar öneme sahip olduğunu, diğer iki modelde ise önemliliğin genel olarak **gender**, **class** ve **age** (GBR'da +**fare**) değişkenleri arasında dağıldığı görülmektedir.



Titanic verileri için **Random Forest**, **Generalized Boosted Regression** ve **Logistic Regression** modelleri için permütasyona dayalı değişken önem ölçümleri. Farklı modeller içinde edilen AUC değerindeki farklılıklar nedeniyle çubukların başlangıç noktalarının farklı olduğuna dikkat ediniz.

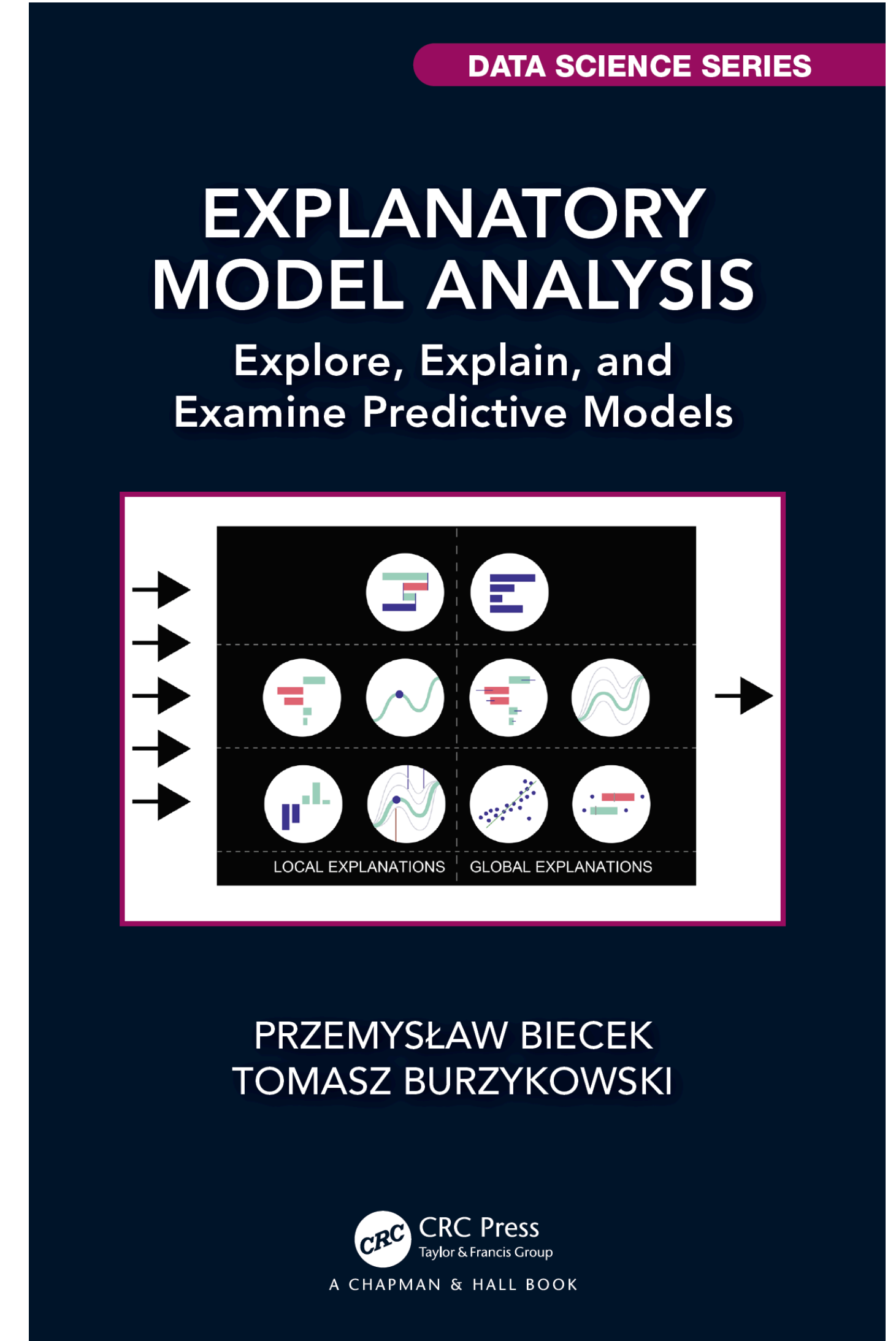
Artı ve eksileri

- Modelden bağımsız
- Görselleştirilebilir
- Anlaşılması kolay
- Değişkenler arası etkileşim hakkında ipucu verebilir.

- Permütasyonların rastgeleliğine bağlıdır. Farklı permütasyonlar ile farklı sonuçlar elde edilebilir.
- Kayıt fonksiyonunun seçimine bağlıdır.

Kaynaklar

Ders materyallerinin hazırlanmasında **Explanatory Model Analysis (Biecek and Burzykowski, 2021)** kitabından yararlanılmıştır. Kitabın ücretsiz online versiyonuna bağlantı üzerinden erişilebilir: <https://ema.drwhy.ai/>



Ders notlarına dersin **GitHub** reposu üzerinden ulaşabilirsiniz.

Ders ile ilgili sorularınız için **mustafacavus@eskisehir.edu.tr** adresi üzerinden benimle iletişime geçebilirsiniz.

Mustafa Cavus, Ph.D.

 Eskişehir Teknik Üniversitesi - İstatistik Bölümü

 mustafacavus@eskisehir.edu.tr

 linktr.ee/mustafacavus