

Veri Temizliđi ve Etiketleme:

Bu alıřma, grnt tabanlı yapay zek modellerinde (zellikle sınıflandırma ve segmentasyon gibi bilgisayarlı gr grevlerinde) kullanılacak veri setinin hazırlanması srecini kapsamaktadır. Bir modelin bařarı s yalnızca kullanılan algoritmaya deđil, byk lde eđitildiđi verinin kalitesine bađlıdır. Bu nedenle veri seti oluřturulurken yalnızca grsellerin toplanması yeterli deđildir; verilerin temizlenmesi, dzenlenmesi, dođru biimde etiketlenmesi ve varsa maskelerin kalitesinin kontrol edilmesi gerekir. Bu blmde veri temizliđi ve etiketleme kavramları aıklanmıř, ardından izlenecek sre adım adım zetlenmiřtir.



Veri temizliđi (data cleaning), veri seti iindeki hatalı veya modelin đrenmesini olumsuz etkileyebilecek rneklerin tespit edilip ayıklanması iřlemidir. Grsel veri setlerinde en sık karřılařılan problemler znrlk farklılıkları, pozlama sorunları ve tekrar eden grntlerdir.

znrlk, grselin piksel boyutunu ifade eder ve modelin detayları đrenebilmesi iin yeterli dzeyde olmalıdır. ok dřk znrlkl grsellerde nesne detayları kaybolduđu iin modelin dođru đrenmesi zorlařır. te yandan ařırı yksek znrlkler iřlem maliyetini arttırabilir ve standartlařtırmayı zorlařtırabilir. Bu nedenle veri setinde znrlklerin belirli bir aralıkta tutulması ve gerekirse yeniden boyutlandırma uygulanması nerilir.

Pozlama (exposure) ise grselin ıřık dengesini ifade eder. Ařırı karanlık veya ařırı parlak grntlerde kıyafet dokusu, renk ayrıntıları ve sınırlar kaybolabilir. Bu durum etiketleme srecini zorlařtırdıđu gibi modelin đrenmesini de olumsuz etkiler. Ayrıca bulanıklık, hareket kaynaklı bozulmalar ve dřk kalite (noise, sıkıřtırma artefaktları) gibi sorunlar da veri setinin genel kalitesini dřrr. Veri temizliđi ařamasında bu tr grseller tespit edilerek veri setinden ıkarılır veya ayrı bir kategori altında iřaretlenir.

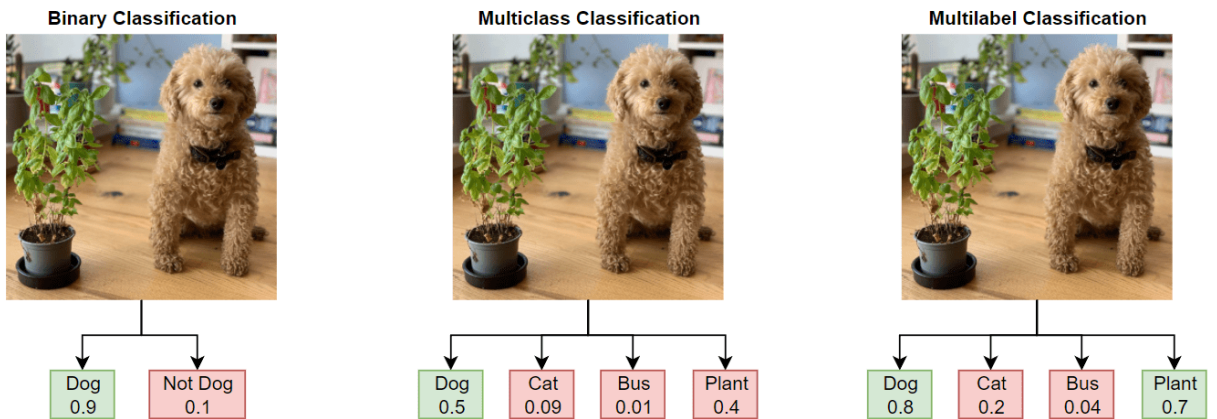
Grsel veri setlerinde dikkat edilmesi gereken bir diđer konu **tekrar veriler (duplicate/near-duplicate)** problemidir. Aynı grselin veya ok benzer grsellerin veri setinde fazla sayıda bulunması modelin genelleme kabiliyetini azaltabilir ve ezberlemeye yol aabilir. Bu nedenle veri setinde birebir tekrar eden grntler ve aynı fotođrafın farklı kırpma/yeniden boyutlandırma gibi varyasyonları kontrol edilerek azaltılmalıdır. Bylece veri seti daha dengeli ve temsil gc yksek hale gelir.

Veri seti temizlendikten sonra etiketleme (labeling) aşamasına geçilir. **Etiketleme**, görsellerin model tarafından öğrenilebilmesi için anlamlı sınıf bilgileriyle işaretlenmesi sürecidir.

- Örneğin kıyafet veri setlerinde görseller “hoodie”, “coat”, “dress” gibi sınıflarla etiketlenebilir. Bu etiketler, modelin eğitim sırasında hangi görselin hangi sınıfa ait olduğunu öğrenmesini sağlar. Etiketleme yapılırken en önemli nokta, etiketlerin tutarlı ve standart olmasıdır. Aynı sınıf için farklı yazımların kullanılması (örneğin “hoodie”, “hoodi”, “kapüşonlu”) veri setinde karışıklığa neden olabilir. Bu nedenle etiketleme öncesinde sınıf isimleri belirlenmeli ve bir etiket sözlüğü oluşturulmalıdır.

Etiketleme sürecinde **multi-label** yaklaşımı da kullanılabilir. Multi-label etiketleme, bir görselin tek bir etikete değil birden fazla etikete sahip olabilmesidir.

- Kıyafet örneğinde bir görsel yalnızca “hoodie” olarak değil, aynı zamanda “black”, “oversize”, “winter” gibi ek özelliklerle de etiketlenebilir. Bu yaklaşım, veri setinin daha zengin hale gelmesini sağlar ve modelin daha ayrıntılı öğrenmesine olanak tanır. Özellikle kıyafet öneri sistemleri, özellik çıkarımı ve detaylı sınıflandırma gibi uygulamalarda multi-label etiketleme önemli bir avantaj sağlar.



Bazı projelerde görsellere başlıklandırma (captioning) uygulanır. **Görsel başlıklandırma**, her görselin kısa bir açıklama cümlesiyle ifade edilmesidir.

- Örneğin “A woman wearing a black hoodie in a studio background” gibi bir cümle, görselin içerdiği temel bilgileri metin olarak temsil eder. Bu yöntem özellikle metin-görsel ilişkisi kuran generatif modellerde ve veri setinin dokümantasyonunda fayda sağlar.

Segmentasyon kullanılan veri setlerinde ayrıca **maske (segmentation mask)** kavramı bulunur. **Segmentasyon maskesi**, görseldeki hedef nesnenin piksel düzeyinde işaretlenmiş halidir. Bu maske sayesinde model yalnızca nesnenin ne olduğunu değil, görüntüde tam olarak hangi bölgede bulunduğunu da öğrenir. Segmentasyon maskelerinin kalitesi model performansını doğrudan etkiler. Kalitesiz maskeler genellikle sınır taşmaları, nesnenin bir kısmının eksik işaretlenmesi, maske içinde kopukluklar veya gereksiz küçük parçalar (gürültü) şeklinde kendini gösterir. Bu nedenle maske kalitesi, görsel üzerinde maske bindirme (overlay) yöntemiyle kontrol edilerek değerlendirilmelidir.



Etiketleme ve maske üretiminde kullanılan araçlar, süreci hızlandırmak ve standartlaştırmak açısından önemlidir. **Label Studio** ve **CVAT** bu amaçla kullanılan yaygın araçlardır.

Label Studio genellikle sınıflandırma, çoklu etiketleme ve görsel başlıklandırma gibi görevlerde pratik bir kullanım sunar.

CVAT ise özellikle bounding box ve segmentasyon gibi detaylı anotasyonlarda daha güçlü araçlara sahiptir. Bu tür araçlar sayesinde etiketleme süreçleri daha düzenli yürütülür, hata oranı azaltılır ve veri seti üretimi daha verimli hale gelir.

Sonuç olarak veri temizliği ve etiketleme aşamaları, model eğitiminden önce gerçekleştirilmesi gereken kritik hazırlık süreçleridir. Etiketleme sürecinde tutarlılık sağlanması, multi-label yapıların doğru kurgulanması ve caption kullanımının değerlendirilmesi veri setinin öğrenilebilirliğini yükseltir. Segmentasyon maskelerinin kalite kontrolü ise piksel düzeyinde öğrenme yapan modellerde doğrudan performans artışı sağlar.