

Este manual oferece orientações práticas para ajustar os parâmetros de modelos populares de machine learning, com base em objetivos comuns de otimização e prevenção de problemas como overfitting e underfitting.

Random Forest e Decision Tree

- **max_depth**: Define a complexidade máxima da árvore. Árvores muito profundas podem aprender padrões detalhados que não se generalizam bem. Uma prática recomendada é experimentar diferentes profundidades ou definir **max_depth** com base em validação cruzada.
- **min_samples_split**: Este parâmetro ajuda a prevenir o aprendizado de relações muito específicas, definindo o número mínimo de amostras que um nó deve ter antes de ser dividido. Aumentar este valor pode suavizar o modelo.
- **min_samples_leaf**: Semelhante ao acima, impede a criação de folhas com poucas amostras, aumentando a generalização.

CatBoost Regressor

- **iterations**: Mais iterações permitem mais ajuste ao conjunto de dados, mas cuidado com overfitting. Use junto com técnicas de parada antecipada.
- **learning_rate**: Um valor menor pode melhorar a generalização, permitindo que o modelo aprenda de forma mais gradual.
- **depth**: Profundidades maiores permitem ao modelo capturar interações mais complexas, mas aumentam o risco de overfitting. Profundidades menores são recomendadas para dados com menor complexidade.
- **l2_leaf_reg**: Regularização L2 pode ser aumentada para controlar o overfitting, penalizando pesos maiores.
- **border_count**: Controla a granularidade com a qual as características são avaliadas. Valores maiores podem aumentar a precisão, mas também o tempo de computação.

LightGBM Regressor

- **n_estimators** e **learning_rate**: O equilíbrio entre esses dois parâmetros é crucial. Um número maior de estimadores com uma taxa de aprendizado baixa pode melhorar a precisão, mas aumentar o risco de overfitting.
- **num_leaves**: Deve ser ajustado com cuidado, pois muitas folhas podem causar overfitting. Geralmente, um número menor para dados menos complexos é preferível.
- **max_depth**: Limitar a profundidade ajuda a prevenir o overfitting e reduz o uso de memória.

LARS (Least Angle Regression)

- **n_nonzero_coefs**: Este parâmetro define um limite no número de variáveis no modelo final, promovendo a seleção de variáveis.
- **fit_intercept**: Decidir se um termo intercept deve ser incluído no modelo pode afetar a interpretação dos coeficientes, especialmente se os dados não estão centralizados.
- **eps**: Controla o critério de parada. Valores menores podem resultar em maior precisão, mas maior tempo de computação.

Passive Aggressive Regressor

- **C**: Controla o trade-off entre seguir o modelo atual e corrigir o erro total. Valores maiores dão mais liberdade para o modelo se ajustar a exceções.
- **max_iter** e **tol**: Esses parâmetros controlam o número de iterações e a tolerância para a convergência, respectivamente. Experimentar com diferentes configurações pode encontrar um balanço entre desempenho e tempo de treinamento.

Lasso e Ridge (L1 e L2 Regularization)

- **alpha**: Ajusta a força da regularização. Valores maiores aumentam a regularização, podendo reduzir o overfitting, mas também simplificar excessivamente o modelo.

Huber Regressor

- **epsilon**: Controla a sensibilidade a outliers ao definir o intervalo no qual os erros são considerados quadráticos. Fora desse intervalo, os erros são tratados linearmente, reduzindo o impacto de outliers.
- **alpha**: Ajusta a regularização do modelo. É útil quando se tem muitas características ou quando se deseja prevenir o overfitting.

SVR (Support Vector Regression)

- **C**: Ajusta a penalidade de erros fora da margem. Um valor maior enfatiza a correta classificação de todos os pontos de treinamento, enquanto um valor menor pode aumentar a generalização.
- **epsilon**: Define a largura da margem de erro; ajustes aqui podem ajudar a evitar o ajuste excessivo aos dados de ruído.
- **kernel**: Escolher o kernel correto (linear, polinomial, RBF) é crucial para o desempenho do modelo, pois define como os dados são projetados em um espaço dimensional superior para facilitar a regressão.

Gaussian Process Regressor

- **kernel**: A escolha do kernel afeta a suavidade da função aprendida. Kernels comuns incluem RBF, que é suave, e Matérn, que pode capturar variações mais abruptas.

- **alpha**: Adiciona um termo de ruído à diagonal do kernel; aumentar este valor pode ajudar a prevenir overfitting, especialmente em casos de ruído nos dados de treinamento.

Orthogonal Matching Pursuit (OMP)

- **n_nonzero_coefs**: Limita o número de variáveis a serem incluídas no modelo, promovendo a esparsidade e potencialmente melhorando a interpretação do modelo.
- **precompute**: Determina se a matriz de Gram deve ser pré-calculada, o que pode acelerar os cálculos em conjuntos de dados grandes, mas aumentar o uso de memória.

Bayesian Ridge

- **tol**: A tolerância para critérios de parada afeta quando o algoritmo para as iterações, baseado na convergência da solução.
- **alpha_1, alpha_2, lambda_1, lambda_2**: Esses hiperparâmetros definem as distribuições a priori dos coeficientes e são cruciais para ajustar a quantidade de encolhimento aplicada aos coeficientes.

ARD Regression

- Similar ao Bayesian Ridge, mas com a capacidade de ajustar automaticamente a relevância de cada variável. Isso permite uma seleção automática e eficiente de características, útil em conjuntos de dados com muitas variáveis irrelevantes.

Ridge (L2 Regularization)

- **alpha**: Controla a força da regularização L2, que penaliza os coeficientes grandes para reduzir o overfitting e melhorar a generalização do modelo.

ElasticNet

- **alpha**: Este parâmetro controla a intensidade total da regularização, combinando L1 e L2 para obter as vantagens de ambos.
- **l1_ratio**: Ajusta a mistura de regularização L1 (lasso) e L2 (ridge), oferecendo controle sobre como as características são penalizadas.

PLS Regression

- **n_components**: Determina o número de componentes para a redução dimensional, impactando diretamente a complexidade do modelo e a capacidade de capturar a variância dos dados.
- **scale**: Decidir se os preditores devem ser escalados pode afetar significativamente o desempenho do modelo, especialmente em dados com unidades de medida variadas.

Dummy Regressor

- **strategy**: A estratégia para fazer previsões pode ser baseada na média, mediana, quantil, ou um valor constante, dependendo do que faz sentido no contexto do problema.

KNN Regressor

- **n_neighbors**: Ajusta o número de vizinhos a serem considerados. Um número menor pode levar a um modelo mais sensível ao ruído, enquanto um número maior pode suavizar demais o resultado.
- **weights**: O peso dado aos vizinhos pode ser uniforme ou variar inversamente com a distância, afetando como as características locais são modeladas.
- **algorithm**: A escolha do algoritmo para computar os vizinhos mais próximos pode impactar a eficiência, dependendo da estrutura dos dados.

Extra Trees Regressor

- Semelhante ao Random Forest, mas com divisões escolhidas aleatoriamente, aumentando a diversidade entre as árvores e potencialmente aumentando a robustez contra overfitting.

LassoLars

- **alpha**: Controla a regularização L1, promovendo um modelo mais esparsos ao penalizar a magnitude dos coeficientes.
- **eps**: Determina a precisão do modelo, influenciando o critério de parada do algoritmo.

MLP Regressor

- **hidden_layer_sizes**: Configura o tamanho e o número de camadas ocultas, que definem a capacidade do modelo de aprender padrões complexos.
- **activation**: A função de ativação determina a transformação não linear aplicada às entradas, impactando a complexidade dos padrões que o modelo pode aprender.
- **solver**: O solucionador para otimização do backpropagation pode afetar a velocidade e qualidade da convergência do treinamento.
- **alpha**: Regularização aplicada aos pesos das camadas, ajudando a evitar o overfitting.

AdaBoost Regressor

- **n_estimators**: Define o número de modelos sequenciais; mais estimadores podem melhorar o desempenho, mas com risco de overfitting.
- **learning_rate**: Ajusta a contribuição de cada modelo subsequente, onde uma taxa menor pode levar a um aprendizado mais robusto.

Gradient Boosting Regressor

- **n_estimators, learning_rate**: Gerenciam o número de estágios de boosting e a taxa com que o boosting aprende, essenciais para balancear entre desempenho e overfitting.
- **subsample**: A fração de amostras usadas para treinar cada árvore. Valores menores podem levar a uma maior diversidade entre as árvores e ajudar a prevenir o overfitting.
- **max_depth**: Controla a profundidade das árvores individuais, com implicações similares às descritas para Random Forest.

XGBoost Regressor

- **n_estimators, learning_rate, max_depth, e subsample**: Semelhante ao Gradient Boosting, mas com uma implementação que suporta mais otimizações e eficiência em grandes conjuntos de dados.

Gaussian Process Regressor (alternativa)

- **kernel**: A escolha do kernel afeta diretamente a flexibilidade do modelo em se adaptar à complexidade dos dados.
- **alpha**: Define o nível de ruído que se espera nos alvos, influenciando como o modelo interpreta variações nos dados de entrada.
- **n_restarts_optimizer**: O número de reinicializações do otimizador pode ajudar a encontrar o ajuste de kernel ótimo, evitando mínimos locais na otimização.

Essas orientações visam fornecer um ponto de partida para a parametrização eficiente de modelos de machine learning, com exemplos práticos de como ajustar esses parâmetros pode impactar o desempenho do modelo em diferentes cenários. Experimentação e validação cruzada continuam sendo essenciais para otimizar os parâmetros para conjuntos de dados específicos.