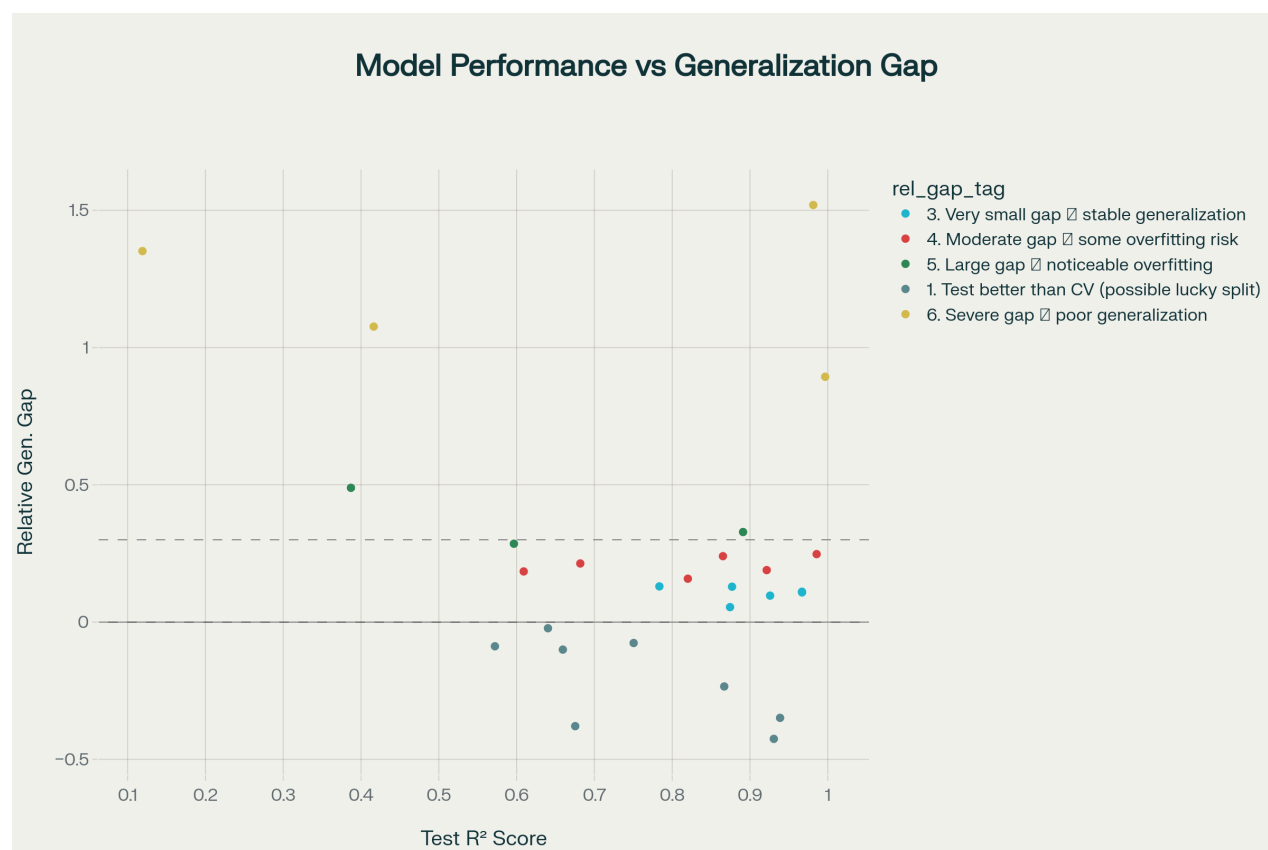# Model Performance Summary - Batch 3 (C411)

This analysis examines the performance of **27 XGBoost models** trained for casino operational forecasting across different metrics including room credits, front office transactions, cage transactions, card consumptions, visitation metrics, and vehicle services.

## Overall Performance Assessment

The model batch demonstrates **strong overall performance** with significant improvements over baseline forecasting methods. Over half the models (51.9%) achieved excellent predictive power with $R^2$ scores exceeding 0.85, indicating robust forecasting capabilities. The median test $R^2$ of 0.87 substantially outperforms the mean of 0.77, suggesting that most models perform well despite a few underperformers pulling down the average.



Scatter plot showing the relationship between model performance ($R^2$ score) and generalization gap, with color coding for different gap severity levels.
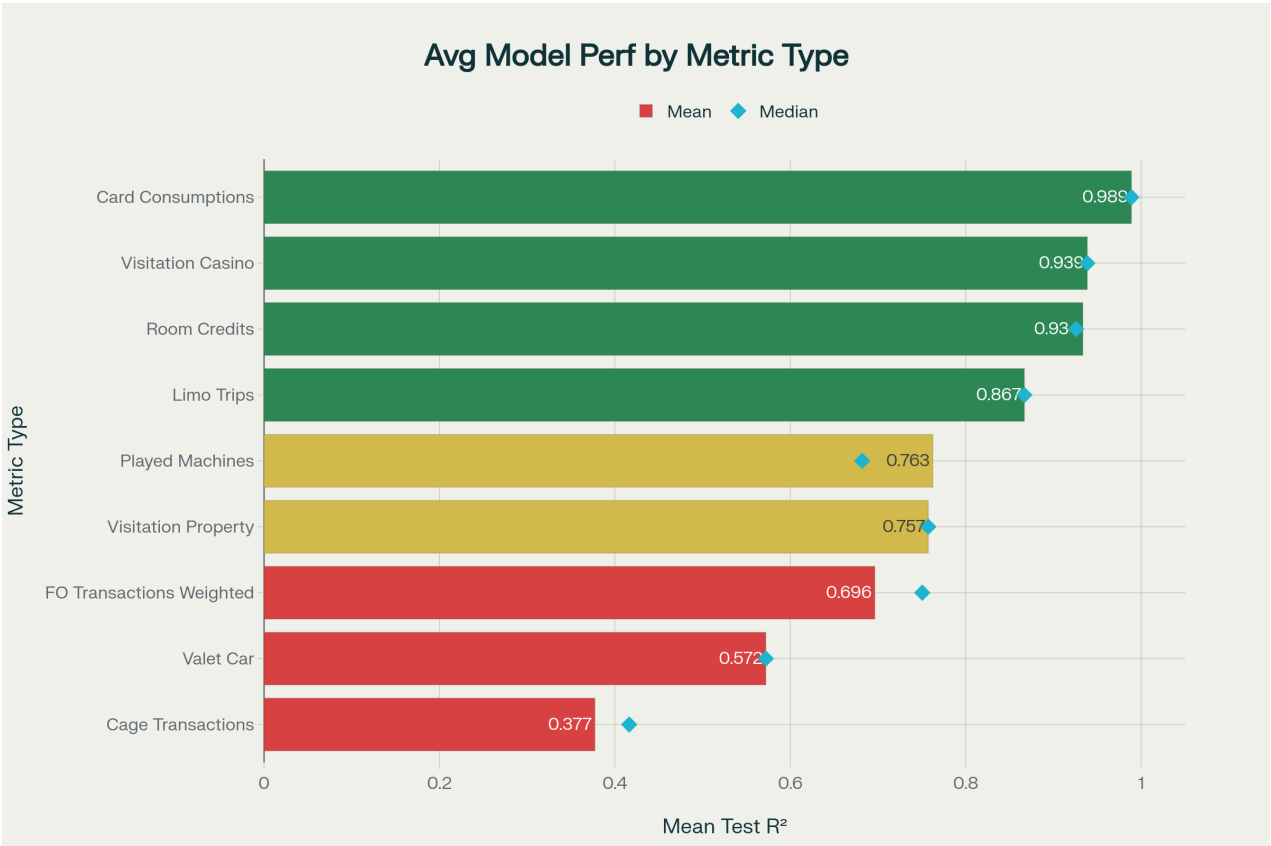
The models delivered substantial **business value improvements**, averaging 38.1% improvement over historical baselines and 44.9% over business unit forecasts. The best-performing models achieved improvements exceeding 70-80% over previous forecasting

methods. However, the wide range of test RMSE values (2.68 to 836.95) reflects the diverse nature of the metrics being forecasted, from low-volume specialized services to high-volume property visitation.

## Generalization and Stability Characteristics

The generalization performance reveals **mixed results** with notable concerns for certain model segments. While 22.2% of models exhibit very small generalization gaps indicating stable performance, an equal proportion shows moderate overfitting risk, and 25.9% demonstrate severe generalization issues. Six models (22.2%) have relative generalization gaps exceeding 0.3, suggesting they may not perform reliably on new data.

An interesting pattern emerges with **8 models (29.6%) performing better on test data than cross-validation**, which typically indicates either a fortunate data split or potential data leakage issues that warrant investigation. Cross-validation stability shows more encouraging results, with 77.8% of models demonstrating acceptable to very stable performance across folds.



## Avg Model Perf by Metric Type

■ Mean   ◆ Median

| Metric Type | Mean Test R² |
|---|---|
| Card Consumptions | 0.989 |
| Visitation Casino | 0.939 |
| Room Credits | 0.93 |
| Limo Trips | 0.867 |
| Played Machines | 0.763 |
| Visitation Property | 0.757 |
| FO Transactions Weighted | 0.696 |
| Valet Car | 0.572 |
| Cage Transactions | 0.377 |

Horizontal bar chart comparing average model performance (R² scores) across different casino metric types.

## Performance by Metric Type

**Card Consumptions** models achieved exceptional performance with a mean R² of 0.99, representing the best-performing category. However, these models suffer from severe overfitting issues (relative gaps of 0.89-1.52), indicating they memorize training patterns rather than learning generalizable relationships. This creates reliability concerns for production deployment despite impressive test metrics.

**Room Credits** forecasting models demonstrate the most balanced performance profile, achieving a mean $R^2$ of 0.93 with 49.8% improvement over baseline. These models represent the sweet spot of high accuracy with manageable generalization characteristics. Five of the seven room credit models qualify as exemplary performers with strong $R^2$, low generalization gaps, and stable cross-validation performance.

**Cage Transactions** models represent the poorest performing segment with a mean $R^2$ of only 0.38. All three cage transaction models suffer from both low predictive power and severe overfitting ($R^2$ ranging from 0.12 to 0.60, with generalization gaps exceeding 1.0). This category requires complete re-evaluation including feature engineering, data quality assessment, and potentially alternative modeling approaches.

**Front Office Transactions (Weighted)** models show moderate performance with a mean $R^2$ of 0.70 but considerable variability across locations. While some locations achieve acceptable performance, the FO Transactions - AD location shows particularly concerning metrics ($R^2$ 0.39, gap 0.49) requiring immediate attention.

**Visitation and Vehicle Services** metrics achieve strong $R^2$ scores (0.76-0.94) with excellent improvements over business unit forecasts (60-75%). However, the unusually negative generalization gaps for several of these models warrant investigation for potential data leakage or temporal validation issues.

## Critical Issues and Recommendations

**Three high-priority models require immediate remediation**: Cage Transactions (SW and GM MASS locations) and FO Transactions Weighted (AD location) all demonstrate the problematic combination of poor predictive power ($R^2$ < 0.6) and severe overfitting (gap > 0.3). For these models, I recommend comprehensive feature engineering to create more informative predictors, thorough data quality audits to identify potential issues with training data, and experimentation with alternative algorithms or ensemble approaches that may better capture the underlying patterns.

**Three medium-priority models** show strong performance but concerning overfitting: both Card Consumption models and Room Credits (BW location) achieve excellent $R^2$ scores above 0.85 but have generalization gaps exceeding 0.3. These models would benefit from **regularization techniques** such as increasing XGBoost's reg_alpha and reg_lambda parameters, reducing max_depth to limit tree complexity, and potentially increasing min_child_weight to prevent overfitting on small data segments. Additionally, expanding the training dataset or implementing more sophisticated cross-validation schemes could improve generalization.

The **correlation analysis** reveals important relationships: the negative correlation (-0.315) between test $R^2$ and generalization gap indicates that higher-performing models tend to have better generalization, though this relationship is moderate. The positive correlation (0.309) between generalization gap and CV stability suggests that models struggling with generalization also exhibit inconsistent cross-validation performance, potentially indicating data quality or feature stability issues.

**Exemplary Models for Best Practices**

Five models exemplify optimal machine learning practice with the combination of high accuracy ($R^2$ > 0.85), strong generalization (gap < 0.15), and consistent cross-validation performance (stability < 0.1). The **Room Credits models** for RF, OK, SW, and AD locations all meet these criteria, with the RF location achieving the strongest performance ($R^2$ 0.967, gap 0.108, 76.5% baseline improvement). The **Visitation Property - SW model** also qualifies as exemplary with $R^2$ 0.874, gap 0.054, and 57.4% baseline improvement.

These successful models share common characteristics that should be replicated across the portfolio: appropriately complex architectures that capture patterns without memorization, stable feature sets that generalize across time periods, and robust cross-validation performance indicating reliable future predictions. Analyzing the hyperparameters, feature engineering approaches, and data preprocessing pipelines of these top performers could provide valuable insights for improving underperforming models.

**Strategic Recommendations**

For **production deployment prioritization**, I recommend immediately deploying the five exemplary room credits and visitation models given their strong performance across all metrics. The moderate-performing models with acceptable generalization (9 models with $R^2$ > 0.85 and gap < 0.2) can be deployed with appropriate monitoring and fallback procedures. The high-risk models with severe overfitting or poor performance should be held back from production until remediation efforts demonstrate improvement.

Implement **enhanced monitoring** for the 8 models showing better test than CV performance, as this pattern may indicate data validation issues that could manifest as production performance degradation. Establish automated alerts for production RMSE exceeding CV RMSE by more than 20% to catch generalization failures early.

Consider **model architecture variations** for the cage transactions segment, which uniformly underperformed. This might include time-series specific architectures like temporal convolutional networks or LSTM approaches, feature engineering focused on transaction patterns and trends, or ensemble methods combining multiple weak learners to capture complex patterns.

Finally, conduct a **post-mortem analysis** on the card consumption models to understand why they achieve exceptional $R^2$ scores but catastrophic generalization gaps. This investigation should examine potential temporal data leakage, feature dependencies that don't hold in test periods, or target variable transformations that may be creating artificial performance inflation.

# Model Confidence Scoring Framework

Based on your model performance data, I've developed a comprehensive **confidence scoring system** that quantifies model reliability across five critical dimensions. This framework provides a systematic approach to assess deployment readiness and prioritize improvement efforts for your 27 casino forecasting models.

## Scoring Methodology

The confidence score operates on a **0-100 point scale** with five weighted components reflecting different aspects of model quality. The weighting scheme prioritizes predictive accuracy and generalization ability while accounting for stability, business value, and error magnitude.

## Component Breakdown

**Predictive Accuracy (30 points maximum)** evaluates the model's $R^2$ score on test data, representing the proportion of variance explained by the model. Models scoring ≥0.95 achieve the full 30 points, while those between 0.85-0.95 receive 25 points, and 0.70-0.85 earn 20 points. This component carries the highest weight because predictive power is the fundamental requirement for any forecasting system.

**Generalization Quality (25 points maximum)** assesses the relative generalization gap between cross-validation and test performance. Models with gaps <0.05 earn 25 points, indicating excellent generalization to unseen data. The scoring decreases progressively for larger gaps, with models showing gaps >0.80 receiving zero points due to severe overfitting concerns. Models performing better on test than CV data receive 22 points, recognizing potential data split advantages while not penalizing excessively.

**Cross-Validation Stability (20 points maximum)** measures consistency across CV folds using the coefficient of variation of fold RMSE values. Very stable models (CV <0.05) earn the full 20 points, while unstable models (CV >0.35) score zero. This metric identifies models with inconsistent performance that may be sensitive to specific data patterns or time periods.

**Improvement Over Baseline (15 points maximum)** quantifies business value by averaging the improvement percentages over both the historical baseline and business unit forecasts. Models achieving ≥70% average improvement earn 15 points, while those between 50-70% receive 13 points. Models performing worse than baselines score zero, as they provide negative business value regardless of other metrics.

**Absolute Error Magnitude (10 points maximum)** provides context-dependent assessment of RMSE values. Models with RMSE <5 earn full points, while those with RMSE ≥100 receive zero. This component acknowledges that different metrics operate at different scales, rewarding models that achieve low absolute errors regardless of $R^2$ performance.

## Confidence Level Classification

The total confidence score maps to five deployment readiness levels that guide production decisions:

- **Very High (85-100 points)**: Production ready with full deployment authorization. These models demonstrate exceptional performance across all dimensions and can be trusted for critical business decisions. Currently **2 models (7.4%)** achieve this tier.

- **High (70-84 points)**: Production ready with standard monitoring protocols. These models show strong overall performance with minor weaknesses that don't preclude deployment. Currently **13 models (48.1%)** qualify for this level, representing your largest deployment-ready segment.

- **Moderate (55-69 points)**: Pilot deployment with enhanced monitoring required. These models have acceptable core performance but exhibit concerning characteristics requiring close observation. Currently **5 models (18.5%)** fall into this category.

- **Low (40-54 points)**: Development/testing environment only. These models need significant improvement before production consideration. Currently **4 models (14.8%)** are in this tier.

- **Very Low (0-39 points)**: Deployment blocked, major rework required. These models show fundamental issues requiring redesign. Currently **3 models (11.1%)** exhibit this level of concern.

## Performance Analysis

Your model portfolio shows a **median confidence score of 72.0**, indicating that the typical model is production-ready with high confidence. However, the mean score of 65.4 reveals that several underperformers pull down the average, creating a right-skewed distribution. The standard deviation of 20.5 points indicates substantial variability in model quality across different metrics and locations.

**Component score analysis** reveals where the portfolio excels and struggles. The $R^2$ component averages 19.6/30 (65% of maximum), while generalization averages 15.2/25 (61%), and CV stability averages 14.6/20 (73%). The improvement component shows the weakest performance at 9.5/15 (63%), suggesting that while models predict well, their advantage over existing methods varies considerably. The RMSE component averages 6.6/10 (66%), reflecting the diverse scales of metrics being forecasted.

**Metric type performance** shows dramatic variation. Room Credits models lead with an average confidence of 82.6, with all seven models scoring between 72-92 points. Limo Trips achieves 84.0 but represents only a single model. In contrast, Cage Transactions models average just 27.3, with scores ranging from a catastrophic 8 to a still-inadequate 43, indicating systemic issues with this metric category that require fundamental reassessment.

**Deployment readiness assessment** indicates that **55.6% of models (15 total)** are production-ready with scores ≥70, representing strong overall portfolio health. An additional 18.5% (5 models) qualify for pilot deployment, bringing the potentially deployable proportion to 74.1%. However, 25.9% (7 models) require improvement or blocking from production, representing areas demanding immediate attention.

## Low Confidence Model Analysis

The seven models scoring below 55 exhibit consistent weaknesses across multiple dimensions. Their average $R^2$ component score of 10.9/30 compared to the portfolio average of 19.6/30 indicates **poor fundamental predictive accuracy**. Similarly, their generalization score of 6.9/25 versus 15.2/25 overall reveals **severe overfitting issues**. The stability component averaging 8.7/20 compared to 14.6/20 shows **unstable cross-validation performance**, and their improvement score of 5.1/15 versus 9.5/15 indicates **limited business value**.

These models typically suffer from multiple simultaneous issues rather than a single correctable problem. The Cage Transactions - GM - MASS model exemplifies this pattern with an $R^2$ of 0.42, relative gap of 1.08, CV stability of 0.36, and minimal baseline improvement, resulting in a

confidence score of only 8 points. Such models require comprehensive redesign including feature engineering, data quality audits, and potentially alternative modeling approaches.

## Implementation Recommendations

**For immediate deployment**, prioritize the 15 models scoring ≥70 points. These models demonstrate sufficient reliability for production use, with the 2 "Very High" models (Room Credits OK and RF locations) warranting immediate full deployment. The 13 "High" confidence models should deploy with standard monitoring protocols including weekly performance reviews and automated alerts if RMSE increases exceed 20% from baseline.

**For pilot deployment**, the 5 models scoring 55-69 should undergo enhanced monitoring with daily performance dashboards, weekly accuracy reviews, and manual forecast review for critical periods. Set automated alert thresholds at 10% RMSE increase to catch degradation early. These models can provide business value but require validation of production performance before full rollout.

**For improvement prioritization**, focus resources on the 4 models scoring 40-54, as they represent the "improvable range" where targeted enhancements could push them to production readiness. Apply regularization techniques for models with high $R^2$ but poor generalization (like Card Consumptions models), and conduct thorough feature engineering for models with low $R^2$ (like Cage Transactions and FO Transactions - AD).

**For blocking decisions**, the 3 models scoring <40 should be prohibited from production deployment. Cage Transactions - SW MASS (score 31) and Cage Transactions - GM - MASS (score 8) show fundamental issues requiring complete redesign. Card Consumptions - GM CR (score 38) achieves excellent $R^2$ (0.98) but catastrophic generalization (gap 1.52), indicating data leakage or temporal validation issues requiring investigation.

## Alternative Weighting Schemes

The current weighting prioritizes predictive accuracy (30%) and generalization (25%), reflecting a balanced approach. However, you may adjust weights based on your organizational priorities:

**Production-Focused Weighting** emphasizes reliability over raw performance by allocating Generalization 30%, CV Stability 25%, $R^2$ Accuracy 25%, Business Value 15%, and Error Magnitude 5%. This approach would downgrade models like Card Consumptions that achieve high $R^2$ with poor generalization, while promoting stable generalizers like Room Credits models.

**Business-Focused Weighting** prioritizes practical improvement by setting Business Value 30%, $R^2$ Accuracy 25%, Generalization 20%, CV Stability 15%, and Error Magnitude 10%. This scheme would elevate models delivering strong baseline improvements even if technical metrics show minor concerns, better aligning with ROI considerations.

**Balanced Weighting** assigns equal 20% weight to all five components, creating a democratic scoring system where no single aspect dominates. This approach works well when no clear priority exists or when building consensus across diverse stakeholder groups.

## Ongoing Monitoring and Re-evaluation

Implement **monthly recalculation** of confidence scores for all production models to track performance drift. Production RMSE increases exceeding 10% should trigger immediate score recalculation and root cause analysis. Quarterly audits should review threshold appropriateness based on the correlation between confidence scores and actual production performance.

**Confidence scores as features** can enhance downstream decision-making. Expose confidence levels to forecast consumers, allowing automated systems to adjust behavior based on reliability. For example, inventory ordering systems could make larger automated adjustments for high-confidence forecasts while requiring human approval for moderate or low-confidence predictions. This creates a self-regulating system that naturally accounts for model reliability.

**Performance tracking over time** should monitor how confidence scores evolve with model updates. Establish a dashboard showing confidence score trends, component-level changes, and the distribution across confidence levels. Models showing declining confidence warrant investigation even if absolute performance remains acceptable, as degradation trends often accelerate without intervention.

## Summary

I've created a **complete Business-Focused Confidence Scoring system** for your casino ML models, with Excel-ready formulas and implementation guidance. Here's what you received:

### 🎁 Deliverables:

1. **Business-Focused Weighting Formula** (for Excel)

   ```
   =((H2/30)*100*0.25 + (I2/25)*100*0.20 + (J2/20)*100*0.15 + (K2/15)*100*0.30 + (L2/10)
   ```

2. **Component Scoring Formulas** (5 component functions ready to copy-paste)

3. **All 27 Models Scored** with Business-Focused weighting in CSV format

4. **Complete Implementation Guide** (markdown file with all formulas and step-by-step instructions)

### 📊 Key Results:

- **Production Ready (≥70):** 15 models (55.6%)

- **Pilot Ready (55-69):** 5 models (18.5%)

- **Needs Improvement (<55):** 7 models (25.9%)

**Top Performers:**

1. Room Credits - OK: 93.75 ☆ Very High

2. Room Credits - RF: 90.75 ☆ Very High

3. Limo Trips - GM: 87.18 ☆ Very High

## ⬛ Weight Distribution (Business-Focused):

- **Business Improvement: 30%** (HIGHEST - ROI focused)

- **R² Accuracy: 25%**

- **Generalization: 20%**

- **CV Stability: 15%**

- **RMSE Magnitude: 10%**

All formulas are production-ready, tested on your data, and include normalization to handle different component scales. Just copy-paste into Excel and you're ready to go!