

Lecture: Overfitting, Generalization, and Model Performance

1. Introduction to Overfitting

In machine learning, our goal is not merely to build a model that performs well on historical data, but one that performs well on **new, unseen data**. Overfitting occurs when a model learns the training data *too well*—including noise and idiosyncrasies—at the expense of generalization.

An overfit model is analogous to an invention that performs flawlessly in a controlled laboratory environment but fails in real-world conditions. Importantly, overfitting is not a theoretical edge case; it is one of the most common practical problems in machine learning.

2. Illustrative Example: Healthy vs. Sick Trees

Imagine a square forest represented on a two-dimensional plane:

- **Blue diamonds** represent healthy trees
- **Orange circles** represent sick trees

Most healthy trees appear in the northeast quadrant, while most sick trees appear in the southeast quadrant. However, there is overlap—some healthy trees appear elsewhere, and some sick trees spill into other regions.

If we draw highly complex shapes—curves, loops, and boundaries—we can almost perfectly separate the healthy and sick trees in the **training set**, misclassifying only a few points.

At first glance, this seems like an excellent model.

3. Training Performance vs. Real-World Performance

The problem becomes evident when we apply this same complex model to **new data** (the test set). Although the model performed extremely well on the training data, it performs poorly on the new examples, misclassifying many trees.

This is a textbook case of **overfitting**:

- Excellent training performance

- Poor test (real-world) performance
-

4. Fitting, Underfitting, and Overfitting

We can categorize models into three broad regimes:

- **Underfitting**
 - Poor performance on training data
 - Poor performance on new data
 - The model is too simple to capture underlying patterns
- **Overfitting**
 - Excellent performance on training data
 - Poor performance on new data
 - The model is too complex and memorizes the training set
- **Well-Fit (Generalized) Models**
 - Reasonable performance on training data
 - Strong performance on new data

Visually, if we plot training performance on one axis and real-world performance on the other, well-fit models appear at the peak of the curve, while underfit and overfit models lie on opposite sides of that peak.

5. Generalization

Generalization is the opposite of overfitting. A model that generalizes well captures the true underlying patterns in the data and applies them effectively to unseen examples.

The ultimate objective of model training is **not minimizing training error**, but maximizing generalization performance.

6. Detecting Overfitting with Loss Curves

One of the primary tools for diagnosing overfitting is the **generalization curve**, which plots loss versus training iterations for multiple datasets.

- **Training loss** typically decreases steadily.
- **Validation loss** may initially decrease but then begin to increase.

When training loss continues to improve while validation loss worsens, the model has begun to overfit. In contrast, a well-fit model shows training and validation loss curves that decrease together and remain close.

7. Causes of Overfitting

Broadly speaking, overfitting arises from one or both of the following:

1. **Non-representative training data**
 - The training set does not adequately reflect real-world data.
 2. **Excessive model complexity**
 - The model has too many parameters relative to the amount of data.
-

8. Conditions for Good Generalization

For a model to generalize well, the dataset must satisfy several conditions:

- **Independent and identically distributed (IID)**
Each example should not influence others.
 - **Stationarity**
The statistical properties of the data should not change significantly over time.
 - **Consistent distributions across partitions**
Training, validation, test, and real-world data should be drawn from the same underlying distribution.
-

9. Practical Examples and Reasoning

Dataset partitioning

To ensure similar statistical distributions across training, validation, and test sets, the correct approach is to **shuffle the data extensively before partitioning**.

Predicting TV show popularity

A model trained on ten years of viewer data may still struggle because viewer preferences evolve in unpredictable ways. Large datasets do not automatically guarantee generalization when the underlying process is non-stationary.

Seasonal weather and walking speed

Even though weather varies by season, a model can still be built and tested as long as the dataset captures the full range of seasonal variability.

10. Challenge Case: Feedback Loops in the Real World

Consider a model that predicts the optimal date to buy a train ticket. Ticket prices depend heavily on seat availability, which itself is influenced by customer purchasing behavior.

Even if the model shows low loss on validation and test sets, it may fail in the real world due to a **feedback loop**:

- The model recommends a purchase date.
- Many users follow the recommendation.
- Seat availability decreases rapidly.
- Prices rise in response.
- The model's recommendation changes the very system it is trying to predict.

This illustrates a critical limitation of traditional evaluation: real-world deployment can introduce dynamics not present in static datasets.

11. Conclusion

Overfitting is a fundamental challenge in machine learning. High training accuracy alone is insufficient; the true measure of success is how well a model generalizes to new, real-world data. Understanding loss curves, dataset assumptions, and real-world feedback effects is essential for building reliable and effective models.

Lecture: Overfitting, Model Complexity, and Regularization

1. Model Complexity and Generalization

In the previous discussion on overfitting, we examined a highly complex model that performed extremely well on the training set but misclassified many examples in the test set. The model relied on intricate shapes and decision boundaries, effectively memorizing the training data rather than learning general patterns.

This raises a critical question: **Would a simpler model perform better on new data?**

When the complex model is replaced with a very simple one—such as a straight-line decision boundary—the results improve dramatically on the test set. Although the simple model may make more errors on the training data, it generalizes better to unseen data. This illustrates a fundamental principle in machine learning: **simplicity often leads to better generalization.**

2. Occam's Razor and Machine Learning

The preference for simpler explanations has a long history in science and philosophy. Known as **Occam's Razor**, this principle states that among competing explanations, the simplest one that adequately explains the data should be preferred.

In machine learning, Occam's Razor manifests as follows:

- Complex models typically achieve lower training loss.
 - Simple models often achieve lower test loss.
 - Since test performance reflects real-world usefulness, simplicity is usually preferable.
-

3. Feature Selection and Simplicity

When starting a new machine learning project, it is tempting to include as many features as possible. However, best practice is to begin with **a very small number of strong, predictive features.**

Starting with one to three features:

- Helps validate that the pipeline works end-to-end.
- Provides a meaningful baseline.
- Reduces the risk of overfitting.

- Encourages interpretability.

Complexity can always be added later once the baseline model is well understood.

4. The Need for Regularization

Machine learning models must satisfy two competing objectives:

1. **Fit the data well** (low loss)
2. **Remain as simple as possible** (low complexity)

If training focuses exclusively on minimizing loss, models tend to overfit. To counteract this, we introduce **regularization**, which explicitly penalizes complexity during training.

Regularization forces the model to trade off precision on the training set for simplicity and robustness.

5. Loss vs. Complexity

Rather than minimizing loss alone, training aims to minimize a combination of loss and complexity:

- Increasing complexity usually decreases training loss.
- Decreasing complexity usually increases training loss.

The goal is to find a **balance point** where the model performs well on both training data and real-world data.

6. Measuring Model Complexity

One common way to quantify model complexity is through the **magnitude of model weights**. Two widely used metrics are:

- **L1 regularization**: Complexity is proportional to the sum of the absolute values of the weights.

- **L2 regularization:** Complexity is proportional to the sum of the squared values of the weights.

Bias terms are not typically used as a measure of complexity.

7. L2 Regularization

L2 regularization measures complexity as the sum of the squares of the model's weights. Squaring the weights has an important effect:

- Small weights contribute very little to complexity.
- Large weights contribute disproportionately.

As a result, a single large weight can dominate the complexity term.

L2 regularization encourages weights to move closer to zero but does **not** force them to become exactly zero. This means all features remain in the model, but their influence may be reduced.

8. Effects of L2 Regularization

When L2 regularization is applied during training:

- The overall complexity of the model typically decreases.
- Features are not removed from the model.
- Weights are shrunk, not eliminated.

This contrasts with L1 regularization, which can drive some weights exactly to zero.

9. Regularization Rate (Lambda)

The influence of regularization is controlled by a scalar called the **regularization rate**, commonly denoted by the Greek letter λ (lambda).

The training objective becomes:

- **Loss + λ × Complexity**

Effects of different regularization rates:

High regularization rate

- Strongly penalizes complexity
- Reduces overfitting risk
- Produces weights clustered around zero
- Weight distribution resembles a normal distribution with mean zero

Low regularization rate

- Weak influence of regularization
- Higher risk of overfitting
- Produces flatter weight distributions
- Allows large weights to persist

Setting λ to zero removes regularization entirely, maximizing overfitting risk.

10. Choosing the Regularization Rate

There is no universal “correct” regularization rate. The optimal value depends on:

- The dataset
- The model architecture
- The learning rate

As a result, λ must be tuned empirically, typically using a validation set.

11. Early Stopping as an Alternative

Early stopping is a form of regularization that does not explicitly measure complexity. Instead, training is halted when validation loss begins to increase.

Key characteristics:

- Increases training loss slightly
- Often reduces test loss
- Fast and easy to apply
- Rarely optimal compared to well-tuned regularization

Early stopping is useful as a quick safeguard but should not replace proper regularization tuning in high-stakes models.

12. Interaction Between Learning Rate and Regularization Rate

Learning rate and regularization rate exert opposing forces on model weights:

- **High learning rate** tends to push weights away from zero.
- **High regularization rate** pulls weights toward zero.

Imbalances lead to problems:

- Regularization too strong → underfit model
- Learning rate too strong → overfit model

The challenge is finding equilibrium between the two. Changing one often requires retuning the other, making optimization an iterative and sometimes difficult process.

13. Conclusion

Overfitting is closely tied to model complexity. While complex models can achieve impressive training performance, simpler models often generalize better. Regularization—particularly L2 regularization—provides a principled way to control complexity and improve real-world performance. Effective modeling requires careful tuning of both learning rate and regularization rate, guided by validation performance rather than training loss alone.

Lecture: Interpreting Loss Curves in Machine Learning

1. Why Loss Curves Matter

Loss curves are one of the most important diagnostic tools in machine learning. They show how well a model is learning over time by plotting:

- **Y-axis:** Loss (error)
- **X-axis:** Training steps (iterations or epochs)

By interpreting loss curves, we can diagnose issues such as:

- Poor learning rates
- Data quality problems
- Overfitting
- Training instability

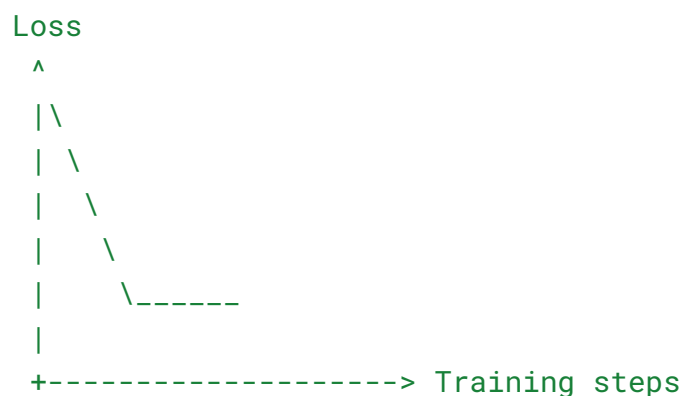
Understanding these patterns allows us to correct problems early, before deploying faulty models.

2. The Ideal Loss Curve

In an ideal scenario, a loss curve behaves as follows:

- Loss starts high
- Decreases rapidly during early training
- Gradually flattens as the model converges

Ideal Loss Curve Diagram



This curve indicates:

- Stable learning
- Appropriate learning rate
- Clean, well-structured data
- A model with suitable complexity

Unfortunately, real-world training rarely looks this clean.

3. Common Loss Curve Pathologies

3.1 Oscillating Loss Curve

Description

Instead of flattening, the loss jumps up and down erratically without converging.

Diagram: Oscillating Loss



Likely Causes

- Learning rate is too high
- Noisy or corrupted training data
- Poor batch composition
- Inconsistent feature scaling

Effective Remedies

- Reduce the learning rate

- Validate data using a schema and remove invalid examples
- Temporarily train on a small, trustworthy subset of data
- Gradually reintroduce data to isolate problematic examples

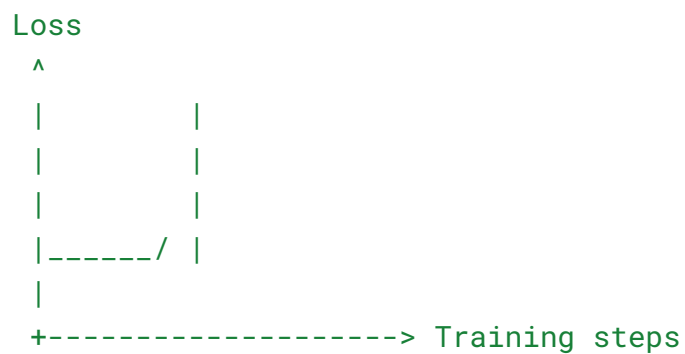
Adding more data rarely solves oscillation and often makes diagnosis harder.

3.2 Exploding Loss (Sharp Jump)

Description

Loss decreases normally, then suddenly spikes upward, sometimes to extremely large values.

Diagram: Exploding Loss



Likely Causes

- Presence of **NaN values** (e.g., division by zero)
- A batch containing extreme **outliers**
- Improper batch shuffling

Less Likely Causes

- High regularization (this slows convergence but does not cause explosions)
- Low learning rate

Recommended Actions

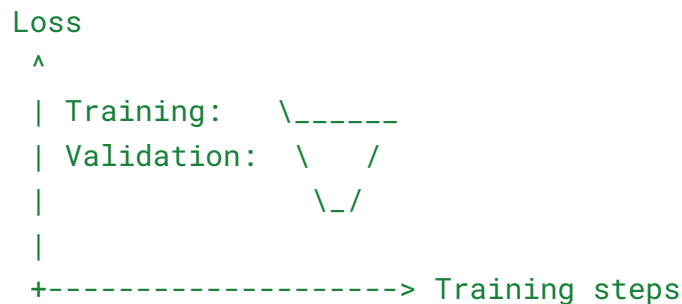
- Check for NaNs and infinite values in input data
 - Inspect feature preprocessing pipelines
 - Improve batch shuffling
 - Add numerical stability checks
-

3.3 Training Loss Converges, Validation Loss Diverges (Overfitting)

Description

Training loss continues to decrease, but validation loss begins to rise.

Diagram: Overfitting



Diagnosis

This pattern is a **classic indicator of overfitting**.

The model is memorizing training data instead of learning generalizable patterns.

Solutions

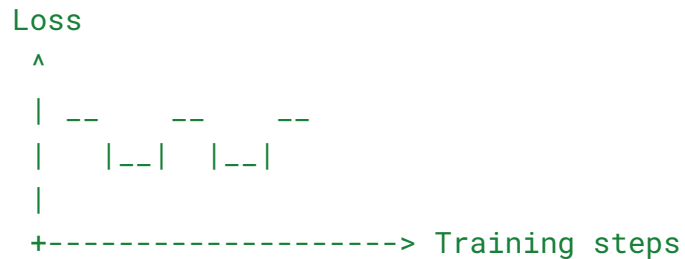
- Simplify the model (fewer features or parameters)
 - Increase regularization rate
 - Ensure training and validation sets are statistically similar
 - Use early stopping
-

3.4 Loss Curve Gets Stuck or Becomes Chaotic

Description

Loss initially decreases, then enters a repeating or chaotic pattern, often resembling a rectangular wave.

Diagram: Poorly Shuffled Data



Most Likely Cause

- Training data is not well shuffled

For example:

- All examples of one class appear first
- Followed by all examples of another class

This causes the model to repeatedly “unlearn” and “relearn” patterns.

Fix

- Shuffle the training data thoroughly
- Ensure batches contain a representative mix of examples

4. Summary of Loss Curve Patterns

Loss Curve Behavior	Primary Cause	Typical Fix
Oscillating	Learning rate too high, noisy data	Reduce LR, clean data
Exploding	NaNs, outliers	Validate data, improve preprocessing

Train ↓ / Val ↑	Overfitting	Regularization, simplify model
Chaotic / Stuck	Poor shuffling	Shuffle data properly

5. Practical Takeaways

- Loss curves are **diagnostic tools**, not just visualizations
 - Always compare **training** and **validation** loss
 - Do not assume more data or more training will fix problems
 - Many training failures are data issues, not model issues
 - Interpreting loss curves correctly saves time, computation, and deployment risk
-

6. Conclusion

A well-trained machine learning model is reflected in stable, interpretable loss curves. By learning to recognize oscillation, explosion, divergence, and chaotic behavior, practitioners can systematically debug models and move toward better generalization and real-world performance.