

# Comprehensive Guide for Interpreting Machine Learning Results in Colab

GPT SKREVET LORT TIL EN HURTIG CTRL F, TIL AT FORKLARE FORHOLD MELLEM TING

---

## 1. Components of a Machine Learning Model Pipeline

### 1.1 Data Preprocessing

- **Feature Scaling:** Ensures all features contribute equally to the model.
  - **Why it matters:** Without scaling, features with larger numerical ranges dominate gradient updates, leading to ineffective learning and suboptimal convergence.
  - **Methods:** Standardization (z-score), which centers features around zero with unit variance, and Min-Max scaling, which normalizes values to a specific range, are common techniques.
- **Handling Missing Data:** Imputation (mean, median, mode) or removal.
  - **Why it matters:** Missing values introduce biases and inaccuracies in learning, affecting the model's ability to generalize.
- **Data Augmentation** (for CNNs and image data): Enhances generalization by artificially expanding the dataset.
  - **Why it matters:** Augmentation helps models learn robust features by exposing them to varied perspectives, such as rotated or flipped images.

### 1.2 Model Architecture

- **Layer Count:**
  - **Why it matters:** More layers enable the model to capture hierarchical, abstract features. For example, in CNNs, early layers detect edges, while deeper layers recognize complex structures like faces. However, excessive layers can lead to overfitting or training inefficiencies without proper regularization. In linear and logistic regression, the absence of layers simplifies the model, making them interpretable but limited in capturing non-linear relationships.
- **Neuron Count per Layer:**
  - **Why it matters:** Each neuron acts as a feature detector. Insufficient neurons in CNNs reduce the model's capacity to learn spatial patterns, missing critical details. In linear regression, having too few independent variables limits the model's ability to explain variance in the dependent variable. In logistic regression, adding irrelevant variables introduces noise, reducing classification accuracy.

- **Activation Functions:**
  - Introduce non-linearity.
  - **Why it matters:** Non-linearity allows networks to learn complex mappings. ReLU (Rectified Linear Unit) mitigates vanishing gradients by outputting zero for negative inputs and identity for positives, maintaining gradient flow in deep networks. Sigmoid, while suitable for probabilistic outputs in logistic regression, compresses gradients for extreme values (near 0 or 1), slowing learning. Linear regression does not use activation functions, as it models relationships directly through weighted sums.

### 1.3 Hyperparameters

- **Learning Rate (l):**
  - **Why it matters:** The learning rate determines the magnitude of weight updates. A high rate risks skipping the optimal solution (divergence), while a low rate slows training, potentially trapping the model in local minima.
- **Batch Size:**
  - **Why it matters:** Small batches introduce noise in gradient estimation, improving generalization by preventing overfitting to specific data patterns. Large batches stabilize gradients but may converge to sharp minima, leading to poor generalization. For regression models, batch size impacts computational efficiency but not the fundamental behavior of the model.
- **Epochs:**
  - **Why it matters:** Training for too many epochs causes the model to overfit the training data, learning noise rather than general patterns. Early stopping prevents this by halting training when validation performance stops improving.

### 1.4 Optimization and Regularization

- **Optimizers:**
  - Adam combines momentum and adaptive learning rates for efficient convergence.
  - **Why it matters:** Momentum accelerates convergence in relevant directions, while adaptive learning ensures appropriate step sizes for each parameter. This balances speed and stability.
  - Linear and logistic regression rely on gradient descent or closed-form solutions to optimize coefficients.
- **Regularization:**
  - L1 (Lasso) sparsifies weights by adding a penalty proportional to their absolute values, removing less important features.
  - L2 (Ridge) prevents large weight values, enhancing generalization.
  - **Why it matters:** Regularization reduces model complexity, forcing it to focus on significant patterns and avoid noise. In logistic regression, this helps handle multicollinearity; in linear regression, it prevents overfitting in high-dimensional spaces.

## 1.5 Loss Functions

- **Cross-Entropy:** Common in classification tasks.
    - **Why it matters:** Cross-entropy penalizes incorrect predictions based on their confidence levels, encouraging the model to output probabilities closer to true labels.
  - **MSE (Mean Squared Error):** Used for regression tasks.
    - **Why it matters:** Squaring errors emphasizes large deviations, ensuring the model prioritizes minimizing substantial prediction errors. In linear regression, MSE represents the fit of the regression line to the data points.
  - **Entropy and Information Gain** (for Decision Trees and Classification):
    - **Why it matters:** Entropy measures impurity or randomness in data. Information gain quantifies the reduction in entropy achieved by splitting data on a feature, guiding decision tree growth.
- 

## 2. Interpreting Metrics and Graphs

### 2.1 Accuracy and Validation Accuracy

- **Graph Patterns:**
  1. **Training Accuracy Increases, Validation Accuracy Plateaus:**
    - **Indication:** Overfitting.
    - **Why it happens:** The model memorizes training data, including noise, but fails to generalize to unseen data.
    - **Fix:** Apply dropout to randomly deactivate neurons during training, forcing the model to rely on diverse feature representations. Alternatively, use early stopping to halt training before overfitting occurs.
  2. **Both Training and Validation Accuracy Plateau at Low Values:**
    - **Indication:** Underfitting.
    - **Why it happens:** The model lacks capacity or is poorly configured (e.g., insufficient layers or neurons).
    - **Fix:** Add more layers or increase neuron count. For linear regression, revisit feature selection and scaling.
  3. **Validation Accuracy Fluctuates:**
    - **Indication:** High learning rate or insufficient data.
    - **Why it happens:** A high learning rate causes oscillations, as the model overshoots optimal parameter values. Insufficient data increases sensitivity to noise.
    - **Fix:** Reduce the learning rate to stabilize convergence or augment data to improve robustness.

## 2.2 Loss and Validation Loss

- **Graph Patterns:**
  1. **Validation Loss Increases as Training Loss Decreases:**
    - **Indication:** Overfitting.
    - **Why it happens:** The model fits noise in the training data, deviating from general patterns useful for validation.
    - **Fix:** Regularization techniques like L2 penalize large weights, discouraging overfitting. In regression, ensure no irrelevant features inflate the model complexity.
  2. **Both Losses Decrease Consistently:**
    - **Indication:** Effective learning.
    - **Why it happens:** The model captures underlying data patterns, optimizing both training and validation performance.

## 2.3 Confusion Matrix

- **Structure:**
  - TP (True Positives), TN (True Negatives), FP (False Positives), FN (False Negatives).
  - **Why it matters:** Provides a detailed breakdown of classification performance, revealing areas of model bias or misclassification.

## 2.4 Evaluation Metrics

- **Precision:** Proportion of positive predictions that are correct.
    - **Why it matters:** High precision reduces false alarms, critical in tasks like spam detection.
  - **Recall:** Proportion of actual positives correctly identified.
    - **Why it matters:** High recall minimizes missed detections, crucial in medical diagnoses.
  - **F1-Score:** Balances precision and recall.
    - **Why it matters:** Useful when both false positives and false negatives carry significant consequences.
  - **R<sup>2</sup> Score** (for Linear Regression):
    - **Why it matters:** Explains the proportion of variance in the dependent variable accounted for by the model. Values closer to 1 indicate better fit.
  - **Log Loss** (for Logistic Regression):
    - **Why it matters:** Measures the accuracy of probabilistic predictions. Lower values indicate better classification confidence.
-

## 3. Success and Failure Factors

### 3.1 Data-Driven Factors

- **Imbalanced Classes:**
  - **Why it matters:** Imbalances skew the model towards majority classes, reducing sensitivity to minority classes.
  - **Fix:** Use oversampling (SMOTE) or class-weighted loss functions to counteract biases.
- **Data Quality:**
  - **Why it matters:** Noisy data introduces errors, reducing model reliability.
  - **Fix:** Apply preprocessing steps like denoising or use robust estimators.

### 3.2 Model-Specific Factors

- **Linear Regression:**
    - **Success:** Captures linear relationships effectively.
    - **Failure:** Poor fit for non-linear data or when multicollinearity exists.
    - **Fix:** Use polynomial regression or include interaction terms.
  - **Logistic Regression:**
    - **Success:** Efficient for binary classification with linearly separable data.
    - **Failure:** Struggles with non-linearity.
    - **Fix:** Introduce feature transformations or kernel methods.
  - **CNNs:**
    - **Success:** Properly tuned filters and pooling layers extract spatial features effectively.
    - **Failure:** Insufficient filters fail to detect detailed features, while excessive filters overfit noise.
  - **RNNs and LSTMs:**
    - **Success:** Capture sequential dependencies through gating mechanisms.
    - **Failure:** Vanishing gradients hinder learning long-term dependencies.
    - **Fix:** Use LSTMs, which maintain gradients through cell states, preserving information over long sequences.
- 

## 4. Relationships Between Components

### 4.1 Layer Depth and Neurons

- **Why it matters:** Deep networks can model hierarchical representations, but excessive depth without regularization leads to vanishing gradients. Balanced depth and width ensure sufficient capacity without overfitting. For regression, adding irrelevant variables can reduce the interpretability and increase multicollinearity.

### 4.2 Activation Functions and Learning Stability

- **Why it matters:** ReLU avoids vanishing gradients by maintaining non-zero gradients for positive inputs, crucial for training deep CNNs. Sigmoid compresses extreme inputs, which can stabilize probabilistic outputs but slows gradient flow in deep layers. Logistic regression’s use of sigmoid ensures outputs represent probabilities.

### 4.3 Optimizers and Loss Convergence

- **Why it matters:** Adam adapts learning rates based on gradient history, stabilizing convergence in complex, high-dimensional spaces. In contrast, SGD provides consistent updates but requires careful tuning to prevent instability. Regression models often use gradient descent without adaptive components.
- 

## 5. Practical Fixes for Common Issues

Issue	Potential Cause	Fix
Validation accuracy plateaus	Overfitting	Apply dropout or L2 regularization.
Loss decreases too slowly	Low learning rate	Increase learning rate or use Adam.
Fluctuating validation metrics	High learning rate	Reduce learning rate.
Underfitting	Insufficient model complexity	Add layers or increase neuron count.
Overfitting	Excessive model complexity	Reduce layers or use early stopping.
Poor generalization	Imbalanced or noisy data	Balance data or use robust preprocessing.

---

# Step-by-Step Explanation of Feeding Training Pairs into Neural Network

---

## 1. Training Pair and One-Hot Encoding

- Suppose you are using the **Skip-Gram** method, where the model tries to predict **context words** given a **target word**.

**Example Sentence:**

"I love coding Python"

- Context window size = 2 (this means the model will consider 2 words before and after the target word).

**Training Pairs:**

For the sentence above, the Skip-Gram model generates the following pairs:

- Target: "love", Context: ["I", "coding"]
- Target: "coding", Context: ["love", "Python"]

**One-Hot Encoding:**

Each word in the vocabulary is represented as a **one-hot vector**. For example:

- Vocabulary: ["I", "love", "coding", "Python"]
- One-hot representations:
  - "I" → [1, 0, 0, 0]
  - "love" → [0, 1, 0, 0]
  - "coding" → [0, 0, 1, 0]
  - "Python" → [0, 0, 0, 1]

---

## 2. Neural Network Input

For a single training pair, the input to the neural network is the **one-hot vector of the target word**. For example:

- Target: "love" → Input vector: [0, 1, 0, 0]
-

### 3. Embedding Layer

- The **embedding layer** is just a lookup table or a matrix of weights, initially filled with random values.

#### Embedding Matrix:

Assume the embedding matrix is:

`[[0.01, -0.02, 0.03], # "I"`

`[0.04, -0.05, 0.06], # "love"`

`[0.07, -0.08, 0.09], # "coding"`

`[0.10, -0.11, 0.12]] # "Python"`

- Shape: `(vocab_size, embedding_dim)`  $\rightarrow (4, 3)$  in this case.
- Each row corresponds to a word, and each column represents a dimension in the embedding space.

When you feed in the one-hot vector for `"love"`, the embedding layer performs a matrix multiplication:

- Input: `[0, 1, 0, 0]`
- Matrix Multiplication: `[0, 1, 0, 0] × [[0.01, -0.02, 0.03], [0.04, -0.05, 0.06], [0.07, -0.08, 0.09], [0.10, -0.11, 0.12]]`
- Result: `[0.04, -0.05, 0.06]` (the embedding vector for "love").

This vector is now passed to the next layer of the neural network.

---

### 4. Output Layer

The next layer predicts the probability distribution over the vocabulary for the context words.

#### Process:

1. The embedding vector `[0.04, -0.05, 0.06]` is passed through a dense layer (softmax layer).
2. The dense layer outputs a vector of probabilities for all words in the vocabulary.

#### Example Output:



- Predicted probabilities:  $[0.70, 0.10, 0.15, 0.05]$ 
    - "I": 70%
    - "love": 10%
    - "coding": 15%
    - "Python": 5%
- 

## 5. True Labels

The true context words for the target "love" are ["I", "coding"].

- True output (context words in one-hot form):
    - "I"  $\rightarrow [1, 0, 0, 0]$
    - "coding"  $\rightarrow [0, 0, 1, 0]$
- 

## 6. Compute Loss

The difference between the predicted probabilities and the true labels is computed using a **loss function** (e.g., cross-entropy loss). The loss quantifies how far the predicted distribution is from the true distribution.

### Example:

For the true context word "I", the cross-entropy loss might look like this:

$$\text{Loss} = -\log(P(\text{"I"})) = -\log(0.70) \approx 0.36$$

If the model assigns low probabilities to the true words, the loss will be higher.

---

## 7. Backpropagation

Using the loss, the model updates the weights in the network, including the embedding matrix.

### How It Works:

1. The gradients of the loss with respect to the weights (including the embedding matrix) are computed using backpropagation.
2. These gradients indicate how much each weight contributes to the loss.

3. The weights are updated using **gradient descent**:  
$$\text{weight\_new} = \text{weight\_old} - \text{learning\_rate} \times \text{gradient}$$

For the embedding matrix, this means the row corresponding to the target word "love" ( $[0.04, -0.05, 0.06]$ ) is adjusted slightly so that the model predicts better probabilities for the true context words.

---

## 8. Repeat for All Training Pairs

The above steps are repeated for all training pairs (e.g., all (target, context) pairs in the corpus). Over many iterations (epochs), the embedding matrix evolves so that:

- Words appearing in similar contexts have similar vectors.
  - Words appearing in different contexts have distinct vectors.
- 

## 9. Final Embedding Matrix

After training, the embedding matrix might look like this:

$[0.10, 0.20, 0.30]$ , # "I"

$[0.40, 0.50, 0.60]$ , # "love"

$[0.70, 0.80, 0.90]$ , # "coding"

$[1.00, 1.10, 1.20]$  # "Python"

Now:

- "love" and "coding" have similar vectors because they often appear in similar contexts.
  - "I" and "Python" have different vectors because their contexts are different.
- 

## Key Takeaways

1. The embedding matrix starts as random weights.

2. During training, each word's vector (row in the embedding matrix) is updated to minimize the error in predicting context words.
3. Over time, the embedding matrix becomes a meaningful representation of words, capturing semantic and syntactic relationships.
4. Words with similar contexts have similar embeddings, enabling tasks like similarity search, analogy detection, and more.

This is how Word2Vec learns the embedding matrix!

## Unigram, Bigram and Trigram: A Step-by-step Solution

---

### Part 1: Corpus

Given the corpus:

1. I saw the boy
2. the man is working
3. I walked in the street

---

### Step 1: Unigram Counts

Unigrams are individual words in the corpus. Count the occurrence of each word:

#### Tokenized Corpus:

[I, saw, the, boy, the, man, is, working, I, walked, in, the, street]

#### Unigram Counts:

Word	Count
I	2
saw	1
the	3
boy	1
man	1
is	1

working	1
walked	1
in	1
street	1

---

## Step 2: Maximum Likelihood Estimate (MLE) for Bigrams

MLE for a bigram is computed as:  $P(w_2 | w_1) = \text{Count}(w_1, w_2) / \text{Count}(w_1)$

### Count All Bigrams:

Bigram	Count
I → saw	1
saw → the	1
the → boy	1
the → man	1
man → is	1
is → working	1
I → walked	1
walked → in	1
in → the	1
the → street	1

### Compute MLE for Bigrams:

1. Example:  $P(\text{the} | \text{saw}) = \text{Count}(\text{saw}, \text{the}) / \text{Count}(\text{saw}) = 1 / 1 = 1$
2. Example:  $P(\text{boy} | \text{the}) = \text{Count}(\text{the}, \text{boy}) / \text{Count}(\text{the}) = 1 / 3$

Bigram	MLE
I → saw	1
saw → the	1
the → boy	1/3

the → man	1/3
man → is	1
is → working	1
I → walked	1
walked → in	1
in → the	1
the → street	1/3

---

### Step 3: Compute the Probability of the Sentence **I saw the man**

Using the bigram probabilities:  $P(\text{I saw the man}) = P(\text{saw} | \text{I}) \times P(\text{the} | \text{saw}) \times P(\text{man} | \text{the})$

From the table:

- $P(\text{saw} | \text{I}) = 1$
- $P(\text{the} | \text{saw}) = 1$
- $P(\text{man} | \text{the}) = 1/3$

$$P(\text{I saw the man}) = 1 \times 1 \times 1/3 = 1/3$$


---

### Step 4: Compute the Probability of the Sentence **I saw the man in the street**

Using the bigram probabilities:  $P(\text{I saw the man in the street}) = P(\text{saw} | \text{I}) \times P(\text{the} | \text{saw}) \times P(\text{man} | \text{the}) \times P(\text{in} | \text{man}) \times P(\text{the} | \text{in}) \times P(\text{street} | \text{the})$

From the table:

- $P(\text{saw} | \text{I}) = 1$
- $P(\text{the} | \text{saw}) = 1$
- $P(\text{man} | \text{the}) = 1/3$
- $P(\text{in} | \text{man}) = 0$  (bigram does not exist)

Since  $P(\text{in} | \text{man}) = 0$ , the probability of the entire sentence is:  $P(\text{I saw the man in the street}) = 0$

---

## Part 2: Another Corpus

Given the second corpus:

1. I love machine learning
  2. I love deep learning
  3. Deep learning is great
  4. Machine learning is fun
- 

## Step 1: Trigram Counts

A trigram is a sequence of three consecutive words. Count all trigrams in the corpus:

### Tokenized Corpus:

[I, love, machine, learning, I, love, deep, learning, Deep, learning, is, great, Machine, learning, is, fun]

### Trigram Counts:

Trigram	Count
I → love → machine	1
love → machine → learning	1
I → love → deep	1
love → deep → learning	1
Deep → learning → is	1
learning → is → great	1
Machine → learning → is	1
learning → is → fun	1

---

## Step 2: Maximum Likelihood Estimate (MLE) for Trigrams

MLE for a trigram is computed as:  $P(w_3 | w_1, w_2) = \text{Count}(w_1, w_2, w_3) / \text{Count}(w_1, w_2)$

### Example Calculation:

1.  $P(\text{machine} | \text{I}, \text{love}) = \text{Count}(\text{I}, \text{love}, \text{machine}) / \text{Count}(\text{I}, \text{love}) = 1 / 2$
2.  $P(\text{learning} | \text{love}, \text{machine}) = \text{Count}(\text{love}, \text{machine}, \text{learning}) / \text{Count}(\text{love}, \text{machine}) = 1 / 1$

NEXT PAGE

**MLE Table:**

Trigram	MLE
I → love → machine	1/2
love → machine → learning	1
I → love → deep	1/2
love → deep → learning	1
Deep → learning → is	1
learning → is → great	1
Machine → learning → is	1
learning → is → fun	1

---

## Key Takeaways

1. **Unigram Counts** simply count individual words.
2. **MLE for Bigrams** involves dividing the bigram count by the count of the first word.
3. **Sentence Probability** is computed by multiplying probabilities of consecutive bigrams or trigrams.
4. **Trigram MLE** divides the trigram count by the count of the first two words.