

Machine Learning Regression Notes (Complete Revision)

These notes cover everything you asked: linear regression, formulas, R^2 , adjusted R^2 , polynomial regression, pipelines, cross-validation, residuals, plotting, and step-by-step explanations.

Use this as your **future revision notebook**.



1. TRAIN–TEST SPLIT

Code:

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, random_state=42)
```

Meaning:

- **X** → Input features
- **y** → Target/output values
- **X_train** → Used to train model
- **y_train** → Correct answers for training
- **X_test** → New unseen data for testing
- **y_test** → Real answers to check model accuracy

Why needed:

To test model on **new unseen** data to avoid cheating.



2. SIMPLE LINEAR REGRESSION

Prediction formula (very important):

$$\hat{y} = \beta_0 + \beta_1 x$$

- β_0 = intercept
- β_1 = slope
- x = input
- \hat{y} = predicted output

How LinearRegression() finds β_0 and β_1 ?

It uses **OLS (Ordinary Least Squares)**:

$$\beta = (X^T X)^{-1} X^T y$$

This formula finds the **best-fit line** with minimum squared error.



3. PREDICTION

```
y_pred = regression.predict(X_test)
```

Uses:

$$\hat{y} = \beta_0 + \beta_1 x$$

Model checks performance using **y_test** vs **y_pred**.



4. RESIDUALS (ERRORS)

`residuals = y_test - y_pred`

- **Positive** → model predicted low
- **Negative** → model predicted high
- **Zero** → perfect prediction

Residual = actual – predicted



5. R-SQUARED (R^2)

Formula:

$$R^2 = 1 - (SS_{\text{res}} / SS_{\text{tot}})$$

Where:

$$SS_{\text{res}} = \sum (y_i - \hat{y}_i)^2$$

$$SS_{\text{tot}} = \sum (y_i - \bar{y})^2$$

Meaning:

- **1** → perfect
 - **0** → model useless (same as mean)
 - **< 0** → model worse than mean
-



6. ADJUSTED R-SQUARED

Formula (copy-friendly):

$$\text{Adjusted } R^2 = 1 - (1 - R^2) * (n - 1) / (n - p - 1)$$

Where:

R^2 = normal R-squared

n = number of samples used for testing

p = number of features

Python equivalent:

```
1 - (1 - score) * (len(y_test) - 1) / (len(y_test) - X_test.shape[1] - 1)
```



7. CROSS-VALIDATION

```
validation_score = cross_val_score(regression, X_train, y_train,  
scoring='neg_mean_squared_error', cv=3)
```

Why needed?

- One train-test split may be lucky or unlucky
- CV tests the model **3 times (cv=3)**
- Average of results gives **true performance**

Why negative MSE?

Because sklearn wants **higher = better**.

MSE is **lower = better** → so sklearn flips the sign.



8. STANDARDIZATION (Scaling)

Formula:

$$z = (x - \text{mean}) / \text{std}$$

Important Rule:

```
scaler.fit(X_train)    # learn mean,std from training  
scaler.transform(X_test) # apply same scaling to test
```

Never fit on test data → avoids data leakage.



9. POLYNOMIAL REGRESSION

```
PolynomialFeatures(degree=n, include_bias=True)
```

Transforms:

$$X \rightarrow [1, X, X^2, X^3, \dots]$$

So linear regression becomes:

$$\hat{y} = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_n x^n$$

include_bias=True → adds 1 column so model can learn β_0 .



10. WHY PIPELINE?

```
Pipeline([  
    ("poly_features", PolynomialFeatures()),  
    ("lin_reg", LinearRegression())  
)
```

Pipeline links steps like LEGO:

1. Create polynomial features
2. Train linear regression

Benefits:

- Clean code
 - No data leakage
 - Works with cross-validation
-



11. SMOOTH CURVE GENERATION

```
X_new = np.linspace(-3, 3, 200).reshape(200, 1)
```

Meaning:

- Make **200 smooth numbers** from -3 to 3
 - Used to draw smooth polynomial curve.
-



12. COMPLETE POLYNOMIAL REGRESSION FUNCTION

Explanation of how a function fits model, predicts curve, plots training/testing points.

- Fit pipeline
 - Predict on X_new
 - Draw curve
 - Draw blue training dots
 - Draw green testing dots
-

13. HOW SYNTHETIC DATA (X, y) WAS CREATED

```
X = 6 * np.random.rand(100,1) - 3
```

```
y = 0.5 * X**2 + 1.5 * X + 2 + np.random.randn(100,1)
```

Meaning:

- X values between -3 and 3
- y created using equation:

$$y = 0.5x^2 + 1.5x + 2 + \text{noise}$$

- Added noise for realistic data.
-

14. PLOTTING KEY POINTS

```
plt.scatter(X_train, y_train) # blue dots  
plt.scatter(X_test, y_test) # green dots  
plt.plot(X_new, y_pred_new) # red curve
```

Shows learning and predictions clearly.

FINAL SUMMARY FOR REVISION

- **Linear Regression** → $\hat{y} = \beta_0 + \beta_1 x$
- **Polynomial Regression** → adds X^2, X^3, \dots

- **OLS formula** → $\beta = (X^T X)^{-1} X^T y$
 - **R² measures fit**, Adjusted R² penalizes extra features
 - **Residuals** → actual – predicted
 - **Cross-validation** → more reliable model testing
 - **Standardization** → fit on train only, transform on test
 - **Pipeline** → keeps process clean and avoids mistakes
 - **linspace** → used to draw smooth curves
-



15. ASSUMPTIONS OF LINEAR REGRESSION

Linear Regression works properly only when these assumptions hold:

- **Linearity**: Relationship between X and y must be linear.
 - **Independence**: Data points should be independent.
 - **Homoscedasticity**: Residuals should have constant variance.
 - **Normality**: Residuals should be normally distributed.
 - **No multicollinearity**: Features should not be strongly correlated.
-



16. LIMITATIONS OF LINEAR REGRESSION

- Cannot model curves (only straight lines)
- Sensitive to outliers
- Requires assumptions to hold

- Performs poorly on nonlinear data
-



17. WHEN TO USE POLYNOMIAL REGRESSION

Use polynomial regression when:

- Data shows a **curve**
- Simple linear regression underfits
- You want to capture **nonlinear patterns**

Avoid very high degrees (overfitting).



18. UNDERFITTING vs OVERFITTING (IMPORTANT)

- **Underfitting:** Model too simple (degree=1)
- **Overfitting:** Model too complex (degree very high)
- **Right fit:** Model captures true pattern without noise

Visualization:

- Low degree → straight line → underfits
 - Medium degree → curve → good fit
 - High degree → zig-zag curve → overfits
-



19. GRADIENT DESCENT (Short Note)

LinearRegression in sklearn uses OLS formula, BUT another way to find β values is **Gradient Descent**.

Gradient Descent updates β values using:

$$\theta_i = \theta_i - \alpha * (\partial J / \partial \theta_i)$$

Used in:

- Deep Learning
 - Very large datasets
 - SGDRegressor, etc.
-

20. SUPER-SIMPLE CHILD-FRIENDLY EXPLANATIONS (10-year-old level)

★ What is a Feature (X)?

A feature is like a **clue**.

Example: "How many hours you studied" is a clue to predict your marks.

★ What is a Target (y)?

The value we want to guess.

Example: Marks in the test.

★ What is a Model?

A smart robot that **learns from examples** and then makes predictions.

★ What is Linear Regression?

A method that draws **the best straight line** through your data.
The line helps the model predict new values.

Example:
More study hours → more marks.

★ What is Polynomial Regression?

Sometimes data is **curved**, not straight.
Polynomial regression draws a **curvy line**.

Example:
Jump height increases quickly at first, then slows down.

★ What are Residuals?

Residuals = **Mistakes**.

Residual = Actual value – Predicted value

If residual = 0 → Perfect prediction.

★ What is R²?

R² tells **how good the model is**, from 0 to 1.

- 1 → Perfect model
 - 0 → Bad model
 - <0 → Worse than guessing
-

★ What is Adjusted R²?

It is a **fair version of R²** when we have many features.
It punishes unnecessary features.

★ What is Cross-Validation?

Instead of checking the model one time, we check it **3–10 times** on different small parts of data.

Like giving your model many mini-tests.

★ What is Standardization?

Make all numbers "fair" so the model doesn't get confused by large values.

$$z = (\text{value} - \text{mean}) / \text{std}$$

★ What is a Pipeline?

A pipeline is like a **machine with steps**.

Example:

1. Make polynomial features
2. Train model

Instead of doing each step separately, the pipeline does it all automatically.

★ What is Underfitting?

Model is too simple → cannot understand the pattern.

★ What is Overfitting?

Model is too smart → memorizes even noise → bad at new data.

★ What is the Perfect Fit?

Model learns the real pattern, not too simple, not too complex.

★ What is Noise?

Random little bumps or mistakes in data.

Example: Measurement errors.

★ Simple Memory Tricks

- **Linear Regression:** best straight line.
 - **Polynomial Regression:** best curve.
 - **Residual:** prediction mistake.
 - **R²:** how good the model is.
 - **Pipeline:** automatic process.
 - **Standardization:** make numbers fair.
 - **Train–Test Split:** study vs exam.
 - **Cross-Validation:** many exams.
-

21. COST FUNCTIONS (LOSS FUNCTIONS) — SIMPLE EXPLANATION

Cost functions tell the model how "bad" its predictions are. The model tries to make the cost as small as possible.

★ Why do we use a Cost Function?

Because the model must **measure its mistakes**.

- Without a cost function, the model has **no idea** if it's doing good or bad.
- The model learns by **reducing this cost**.
- Cost function = teacher checking homework mistakes.

We use cost functions to:

- Compare models
 - Improve accuracy
 - Tune parameters
 - Guide Gradient Descent (the model moves toward lower cost)
-

Cost functions tell the model how "bad" its predictions are. The model tries to make the cost as small as possible.

★ Mean Squared Error (MSE)

Formula:

$$\text{MSE} = (1/m) * \sum (y_i - \hat{y}_i)^2$$

Child-version: "Square the mistakes, add them, and take the average." Squaring makes big mistakes hurt more.

When to use: Standard choice for regression and when you want large errors punished.

★ Mean Absolute Error (MAE)

Formula:

$$\text{MAE} = (1/m) * \sum |y_i - \hat{y}_i|$$

Child-version: "Take how much you missed each time, add them, then average." Treats all mistakes equally.

When to use: Good when outliers exist and you want robustness.

★ Root Mean Squared Error (RMSE)

Formula:

$$\text{RMSE} = \text{sqrt}(\text{MSE})$$

Child-version: "Square root of MSE — brings the error back to the same units as Y."

★ Huber Loss (Best of both worlds)

Huber loss uses MAE for large errors and MSE for small errors. It is robust yet smooth.

★ Why cost functions matter (child friendly)

Imagine the model is a student and cost is the number of pages of homework to fix mistakes. The model learns to reduce its homework (cost) by improving predictions.

■ 22. CONVERGENCE IN MACHINE LEARNING (Very Simple Explanation)

Convergence means **the model has learned enough** and **stops improving**.

★ What is Convergence?

Imagine climbing down a hill to reach the lowest point (minimum error). When you reach the lowest point and cannot go lower → **you converged**.

In training:

- Cost keeps decreasing

- Then becomes almost constant
- Model stops changing weights

This is convergence.

★ Gradient Descent Convergence (for regression)

Gradient Descent updates parameters using:

$$\theta_i = \theta_i - \alpha * (\partial J / \partial \theta_i)$$

It repeats this many times until it converges.

How do we know it converged?

1. Cost function stops decreasing
 2. Change in parameters becomes tiny
 3. Algorithm reaches max iterations
-

★ Learning Rate (α) and Convergence

Learning rate decides how big steps we take.

↑ High Learning Rate

- Big jumps
- May skip the best point
- Might never converge (too fast)

↓ Low Learning Rate

- Small steps
- Very slow
- Takes many iterations

✓ Perfect Learning Rate

- Smooth decreasing curve
 - Converges quickly
-

★ Types of Convergence Behavior

- **Fast convergence:** cost drops quickly
 - **Slow convergence:** cost decreases slowly
 - **Divergence:** cost increases (bad learning rate)
 - **Oscillation:** cost jumps up and down (learning rate too high)
-

★ Stopping Criteria (When to stop training)

Training stops when:

- Cost change < tiny threshold (e.g., 1e-6)
 - Maximum iterations reached
 - Gradient becomes close to zero
 - Validation error starts increasing (early stopping)
-

★ Child-Friendly Example

Imagine sliding down a slide.

- At the top → high cost
 - Sliding down → cost decreasing
 - When you reach bottom → convergence
 - If you jump too fast (high learning rate), you fall off.
 - If you go too slow (low learning rate), you take forever.
-

★ Visual: Cost vs Iterations

Cost keeps dropping like:

100 → 50 → 20 → 10 → 5 → 4.5 → 4.4 → 4.39 → 4.389 → (flat)

When numbers stop improving → **converged**.

If you want, I can add a small diagram showing cost decreasing over iterations or show how convergence looks in real code. 😊

Machine Learning Regression Notes (Complete Revision)

These notes cover everything you asked: linear regression, formulas, R^2 , adjusted R^2 , polynomial regression, pipelines, cross-validation, residuals, plotting, and step-by-step explanations.

Use this as your **future revision notebook**.



1. TRAIN–TEST SPLIT

Code:

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, random_state=42)
```

Meaning:

- **X** → Input features
- **y** → Target/output values
- **X_train** → Used to train model
- **y_train** → Correct answers for training
- **X_test** → New unseen data for testing
- **y_test** → Real answers to check model accuracy

Why needed:

To test model on **new unseen** data to avoid cheating.



2. SIMPLE LINEAR REGRESSION

Prediction formula (very important):

$$\hat{y} = \beta_0 + \beta_1 x$$

- β_0 = intercept
- β_1 = slope
- x = input
- \hat{y} = predicted output

How LinearRegression() finds β_0 and β_1 ?

It uses **OLS (Ordinary Least Squares)**:

$$\beta = (X^T X)^{-1} X^T y$$

This formula finds the **best-fit line** with minimum squared error.



3. PREDICTION

`y_pred = regression.predict(X_test)`

Uses:

$$\hat{y} = \beta_0 + \beta_1 x$$

Model checks performance using **y_test** vs **y_pred**.



4. RESIDUALS (ERRORS)

`residuals = y_test - y_pred`

- **Positive** → model predicted low
- **Negative** → model predicted high
- **Zero** → perfect prediction

Residual = actual – predicted



5. R-SQUARED (R^2)

Formula:

$$R^2 = 1 - (SSres / SSTot)$$

Where:

$$SSres = \sum (y_i - \hat{y}_i)^2$$

$$SSTot = \sum (y_i - \bar{y})^2$$

Meaning:

- **1** → perfect
 - **0** → model useless (same as mean)
 - **< 0** → model worse than mean
-



6. ADJUSTED R-SQUARED

Formula (copy-friendly):

$$\text{Adjusted } R^2 = 1 - (1 - R^2) * (n - 1) / (n - p - 1)$$

Where:

R^2 = normal R-squared

n = number of samples used for testing

p = number of features

Python equivalent:

$$1 - (1 - score) * (len(y_test) - 1) / (len(y_test) - X_test.shape[1] - 1)$$



7. CROSS-VALIDATION

```
validation_score = cross_val_score(regression, X_train, y_train,  
                                   scoring='neg_mean_squared_error', cv=3)
```

Why needed?

- One train-test split may be lucky or unlucky
- CV tests the model **3 times (cv=3)**
- Average of results gives **true performance**

Why negative MSE?

Because sklearn wants **higher = better**.

MSE is **lower = better** → so sklearn flips the sign.



8. STANDARDIZATION (Scaling)

Formula:

$$z = (x - \text{mean}) / \text{std}$$

Important Rule:

```
scaler.fit(X_train)    # learn mean,std from training  
scaler.transform(X_test) # apply same scaling to test
```

Never fit on test data → avoids data leakage.



9. POLYNOMIAL REGRESSION

```
PolynomialFeatures(degree=n, include_bias=True)
```

Transforms:

$$X \rightarrow [1, X, X^2, X^3, \dots]$$

So linear regression becomes:

$$\hat{y} = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_n x^n$$

include_bias=True → adds 1 column so model can learn β_0 .



10. WHY PIPELINE?

```
Pipeline([
    ("poly_features", PolynomialFeatures()),
    ("lin_reg", LinearRegression())
])
```

Pipeline links steps like LEGO:

1. Create polynomial features
2. Train linear regression

Benefits:

- Clean code
 - No data leakage
 - Works with cross-validation
-



11. SMOOTH CURVE GENERATION

```
X_new = np.linspace(-3, 3, 200).reshape(200, 1)
```

Meaning:

- Make **200 smooth numbers** from -3 to 3
 - Used to draw smooth polynomial curve.
-

12. COMPLETE POLYNOMIAL REGRESSION FUNCTION

Explanation of how a function fits model, predicts curve, plots training/testing points.

- Fit pipeline
 - Predict on X_new
 - Draw curve
 - Draw blue training dots
 - Draw green testing dots
-

13. HOW SYNTHETIC DATA (X, y) WAS CREATED

$X = 6 * np.random.rand(100,1) - 3$

$y = 0.5 * X^{**2} + 1.5 * X + 2 + np.random.randn(100,1)$

Meaning:

- X values between -3 and 3
- y created using equation:

$$y = 0.5x^2 + 1.5x + 2 + \text{noise}$$

- Added noise for realistic data.
-



14. PLOTTING KEY POINTS

```
plt.scatter(X_train, y_train) # blue dots  
plt.scatter(X_test, y_test) # green dots  
plt.plot(X_new, y_pred_new) # red curve
```

Shows learning and predictions clearly.



FINAL SUMMARY FOR REVISION

- **Linear Regression** → $\hat{y} = \beta_0 + \beta_1 x$
 - **Polynomial Regression** → adds X^2, X^3, \dots
 - **OLS formula** → $\beta = (X^T X)^{-1} X^T y$
 - **R² measures fit**, Adjusted R² penalizes extra features
 - **Residuals** → actual – predicted
 - **Cross-validation** → more reliable model testing
 - **Standardization** → fit on train only, transform on test
 - **Pipeline** → keeps process clean and avoids mistakes
 - **linspace** → used to draw smooth curves
-

15. ASSUMPTIONS OF LINEAR REGRESSION

Linear Regression works properly only when these assumptions hold:

- **Linearity:** Relationship between X and y must be linear.
 - **Independence:** Data points should be independent.
 - **Homoscedasticity:** Residuals should have constant variance.
 - **Normality:** Residuals should be normally distributed.
 - **No multicollinearity:** Features should not be strongly correlated.
-

16. LIMITATIONS OF LINEAR REGRESSION

- Cannot model curves (only straight lines)
 - Sensitive to outliers
 - Requires assumptions to hold
 - Performs poorly on nonlinear data
-

17. WHEN TO USE POLYNOMIAL REGRESSION

Use polynomial regression when:

- Data shows a **curve**
- Simple linear regression underfits
- You want to capture **nonlinear patterns**

Avoid very high degrees (overfitting).

18. UNDERFITTING vs OVERFITTING (IMPORTANT)

- **Underfitting:** Model too simple (degree=1)
- **Overfitting:** Model too complex (degree very high)
- **Right fit:** Model captures true pattern without noise

Visualization:

- Low degree → straight line → underfits
 - Medium degree → curve → good fit
 - High degree → zig-zag curve → overfits
-

19. GRADIENT DESCENT (Short Note)

LinearRegression in sklearn uses OLS formula, BUT another way to find β values is **Gradient Descent**.

Gradient Descent updates β values using:

$$\theta_i = \theta_i - \alpha * (\partial J / \partial \theta_i)$$

Used in:

- Deep Learning
 - Very large datasets
 - SGDRegressor, etc.
-

20. SUPER-SIMPLE CHILD-FRIENDLY EXPLANATIONS (10-year-old level)

★ What is a Feature (X)?

A feature is like a **clue**.

Example: "How many hours you studied" is a clue to predict your marks.

★ What is a Target (y)?

The value we want to guess.

Example: Marks in the test.

★ What is a Model?

A smart robot that **learns from examples** and then makes predictions.

★ What is Linear Regression?

A method that draws **the best straight line** through your data.

The line helps the model predict new values.

Example:

More study hours → more marks.

★ What is Polynomial Regression?

Sometimes data is **curved**, not straight.

Polynomial regression draws a **curvy line**.

Example:

Jump height increases quickly at first, then slows down.

★ What are Residuals?

Residuals = **Mistakes**.

Residual = Actual value – Predicted value

If residual = 0 → Perfect prediction.

★ What is R²?

R² tells **how good the model is**, from 0 to 1.

- 1 → Perfect model
 - 0 → Bad model
 - <0 → Worse than guessing
-

★ What is Adjusted R²?

It is a **fair version of R²** when we have many features.

It punishes unnecessary features.

★ What is Cross-Validation?

Instead of checking the model one time, we check it **3–10 times** on different small parts of data.

Like giving your model many mini-tests.

★ What is Standardization?

Make all numbers "fair" so the model doesn't get confused by large values.

$$z = (\text{value} - \text{mean}) / \text{std}$$

★ What is a Pipeline?

A pipeline is like a **machine with steps**.

Example:

1. Make polynomial features
2. Train model

Instead of doing each step separately, the pipeline does it all automatically.

★ What is Underfitting?

Model is too simple → cannot understand the pattern.

★ What is Overfitting?

Model is too smart → memorizes even noise → bad at new data.

★ What is the Perfect Fit?

Model learns the real pattern, not too simple, not too complex.

★ What is Noise?

Random little bumps or mistakes in data.

Example: Measurement errors.

★ Simple Memory Tricks

- **Linear Regression:** best straight line.
- **Polynomial Regression:** best curve.
- **Residual:** prediction mistake.

- **R²**: how good the model is.
 - **Pipeline**: automatic process.
 - **Standardization**: make numbers fair.
 - **Train–Test Split**: study vs exam.
 - **Cross-Validation**: many exams.
-

21. COST FUNCTIONS (LOSS FUNCTIONS) — SIMPLE EXPLANATION

Cost functions tell the model how "bad" its predictions are. The model tries to make the cost as small as possible.

★ Why do we use a Cost Function?

Because the model must **measure its mistakes**.

- Without a cost function, the model has **no idea** if it's doing good or bad.
- The model learns by **reducing this cost**.
- Cost function = teacher checking homework mistakes.

We use cost functions to:

- Compare models
 - Improve accuracy
 - Tune parameters
 - Guide Gradient Descent (the model moves toward lower cost)
-

Cost functions tell the model how "bad" its predictions are. The model tries to make the cost as small as possible.

★ Mean Squared Error (MSE)

Formula:

$$\text{MSE} = (1/m) * \sum (y_i - \hat{y}_i)^2$$

Child-version: "Square the mistakes, add them, and take the average." Squaring makes big mistakes hurt more.

When to use: Standard choice for regression and when you want large errors punished.

★ Mean Absolute Error (MAE)

Formula:

$$\text{MAE} = (1/m) * \sum |y_i - \hat{y}_i|$$

Child-version: "Take how much you missed each time, add them, then average." Treats all mistakes equally.

When to use: Good when outliers exist and you want robustness.

★ Root Mean Squared Error (RMSE)

Formula:

$$\text{RMSE} = \sqrt{\text{MSE}}$$

Child-version: "Square root of MSE — brings the error back to the same units as Y."

★ Huber Loss (Best of both worlds)

Huber loss uses MAE for large errors and MSE for small errors. It is robust yet smooth.

★ Why cost functions matter (child friendly)

Imagine the model is a student and cost is the number of pages of homework to fix mistakes. The model learns to reduce its homework (cost) by improving predictions.

■ 22. CONVERGENCE IN MACHINE LEARNING (Very Simple Explanation)

Convergence means the model has learned enough and stops improving.

★ What is Convergence?

Imagine climbing down a hill to reach the lowest point (minimum error). When you reach the lowest point and cannot go lower → **you converged**.

In training:

- Cost keeps decreasing
- Then becomes almost constant
- Model stops changing weights

This is convergence.

★ Gradient Descent Convergence (for regression)

Gradient Descent updates parameters using:

$$\theta^{(k)} = \theta^{(k)} - \alpha * (\partial J / \partial \theta^{(k)})$$

It repeats this many times until it converges.

How do we know it converged?

-
1. Cost function stops decreasing
 2. Change in parameters becomes tiny
 3. Algorithm reaches max iterations
-

★ Learning Rate (α) and Convergence

Learning rate decides how big steps we take.

↑ High Learning Rate

- Big jumps
- May skip the best point
- Might never converge (too fast)

↓ Low Learning Rate

- Small steps
- Very slow
- Takes many iterations

✓ Perfect Learning Rate

- Smooth decreasing curve
 - Converges quickly
-

★ Types of Convergence Behavior

- **Fast convergence:** cost drops quickly
 - **Slow convergence:** cost decreases slowly
 - **Divergence:** cost increases (bad learning rate)
 - **Oscillation:** cost jumps up and down (learning rate too high)
-

★ Stopping Criteria (When to stop training)

Training stops when:

- Cost change < tiny threshold (e.g., 1e-6)
 - Maximum iterations reached
 - Gradient becomes close to zero
 - Validation error starts increasing (early stopping)
-

★ Child-Friendly Example

Imagine sliding down a slide.

- At the top → high cost
 - Sliding down → cost decreasing
 - When you reach bottom → convergence
 - If you jump too fast (high learning rate), you fall off.
 - If you go too slow (low learning rate), you take forever.
-

★ Visual: Cost vs Iterations

Cost keeps dropping like:

100 → 50 → 20 → 10 → 5 → 4.5 → 4.4 → 4.39 → 4.389 → (flat)

When numbers stop improving → **converged**.

If you want, I can add a small diagram showing cost decreasing over iterations or show how convergence looks in real code. 😊