# Data Science Notes

Bhavjot Khurana

November 11, 2025

# Contents

# 1   Introduction

These notes follow the FreeCodeCamp data science crash course (`https://www.youtube.com/watch?v=XU5pw3QRYjQ`). The aim is to keep the technical pieces intact while explaining every term in plain language so you can read, remember, and speak about each topic without warming up elsewhere.

Use these pages to:

- Recall what a modeling tool does and why its assumptions matter.

- Keep core equations close at hand and understand the symbols they use.

- Practice interview-style explanations that sound confident yet approachable.

# 2   Linear Regression

## 2.1   Definition

Linear regression fits a straight line through paired data $(x_i, y_i)$ to predict a numeric response:

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad \varepsilon_i \sim \text{Normal}(0, \sigma^2).$$

$\beta_0$ is the intercept (predicted $y$ when $x = 0$), $\beta_1$ is the slope (change in $y$ per one-unit change in $x$), and $\varepsilon_i$ is a random error term assumed to be independent, mean-zero, and with constant variance. Those assumptions justify the straight-line summary and the inference that follows.

## 2.2   Key Formulas

Ordinary least squares (OLS) chooses the coefficients that minimize squared prediction errors:

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i.$$

Closed-form solutions connect back to sample moments:

$$\hat{\beta}_1 = \frac{\sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{n} (x_i - \bar{x})^2}, \qquad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}.$$

The numerator is the sample covariance between $x$ and $y$; the denominator is the variance of $x$. Their ratio captures how strongly $x$ pushes $y$ up or down.

## 2.3  Error Function

The residual for observation $i$ is the prediction error

$$e_i = y_i - \hat{y}_i.$$

Plot residuals versus fitted values to verify that the errors look like centered noise. Patterns, curvature, or fan shapes warn that the model is missing structure or violating homoscedasticity (constant variance).

## 2.4  Mean Squared Error (MSE)

Mean Squared Error averages the squared residuals:

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2.$$

Squaring keeps negative and positive mistakes from canceling, penalizes large misses more than small ones, and yields a smooth objective that calculus-based methods can optimize easily.

## 2.5  How the Line is Fitted

Fitting the line means solving

$$\min_{\beta_0, \beta_1} \sum_{i=1}^{n} (y_i - \beta_0 - \beta_1 x_i)^2.$$

Taking derivatives with respect to $\beta_0$ and $\beta_1$ produces the normal equations $(X^\top X)\hat{\beta} = X^\top y$. Software solves them with stable matrix decompositions (QR, SVD), while gradient descent variants handle massive datasets by streaming through mini-batches.

## 2.6  Hypothesis Testing and p-values

To check whether $x$ truly explains $y$, test the null hypothesis $H_0 : \beta_1 = 0$ using

$$t = \frac{\hat{\beta}_1}{\text{SE}(\hat{\beta}_1)}.$$

Under the null, $t$ follows a $t_{n-2}$ distribution. The p-value is the probability of seeing a $t$-statistic this extreme if the true slope were zero. A small p-value means the observed relationship is unlikely to come from noise alone.

## 2.7 Residual Standard Error (RSS and TSS)

Two sums of squares summarize variation:

$$\text{RSS} = \sum_{i=1}^{n}(y_i - \hat{y}_i)^2, \qquad \text{TSS} = \sum_{i=1}^{n}(y_i - \bar{y})^2.$$

Residual Standard Error (RSE) rescales RSS into the original units of $y$:

$$\text{RSE} = \sqrt{\frac{\text{RSS}}{n-2}}.$$

Think of RSE as the typical distance between the observed data and the fitted line; smaller values indicate a tighter fit.

## 2.8 Interpretations

- **Slope** $\hat{\beta}_1$: Expected change in $y$ for a one-unit increase in $x$; the sign shows direction, the magnitude shows strength.

- **Intercept** $\hat{\beta}_0$: Predicted $y$ when $x = 0$. Interpret only if $x = 0$ exists in the data or has practical meaning.

- $R^2 = 1 - \frac{\text{RSS}}{\text{TSS}}$: Fraction of the variance in $y$ that the model explains.

- **Prediction intervals**: Combine $\hat{y}_i$ with interval bands to communicate both the estimate and its uncertainty.

## 2.9 Example: Simple Linear Regression Output

Using the `Advertising` dataset (TV spend predicting sales), a `statsmodels` fit reports:

| Statistic | Estimate | Interpretation |
| --- | --- | --- |
| Intercept ($\hat{\beta}_0$) | 7.03 | Expected sales if TV spend were \$0k. |
| TV coefficient ($\hat{\beta}_1$) | 0.0475 | Each \$1k on TV adds about 0.048k units sold. |
| $R^2$ | 0.612 | TV spend accounts for 61% of sales variation. |
| Residual Std. Error | 3.26 | Typical prediction miss in sales units. |
| F-statistic | 312.1 | Strong evidence that the slope is not zero. |

Small standard errors create huge t-statistics, so the slope's p-value is effectively zero, and residual diagnostics (Omnibus, Jarque-Bera, Durbin-Watson) show no obvious violations.

## 2.10 Interview Questions and Answers

1. **Question:** When would you question the linearity assumption?

   **Solution:** Scatter plots or residual plots that curve, fan out, or show autocorrelation indicate the straight-line model is missing structure. Add features, transform variables, or switch to flexible methods.

2. **Question:** How do you construct a 95% confidence interval for $\beta_1$?

   **Solution:** Use $\hat{\beta}_1 \pm t_{0.975, n-2} \times \text{SE}(\hat{\beta}_1)$. The $t$-quantile reflects sampling variability with $n - 2$ degrees of freedom.

3. **Question:** Why does OLS minimize squared errors instead of absolute errors?

   **Solution:** Squared errors are differentiable and lead to closed-form solutions while emphasizing large deviations. Absolute errors give a more robust median-based fit but require linear programming.

4. **Question:** What does $R^2 = 0.70$ mean?

   **Solution:** The model explains 70% of the variance in $y$; the rest comes from noise or missing predictors. It is a descriptive statistic, not proof of causation.

5. **Question:** Residuals show a funnel pattern (variance increases with fitted values). What now?

   **Solution:** That is heteroscedasticity. Consider transforming $y$, modeling variance explicitly (weighted least squares), or using robust standard errors for inference.

# 3 Multiple Linear Regression

## 3.1 Definition

Multiple Linear Regression (MLR) extends the straight-line idea to $p$ predictors:

$$y_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} + \varepsilon_i$$

Each coefficient captures the expected change in $y$ for a one-unit bump in its predictor while everything else stays fixed (the "all else equal" clause). Writing the model in matrix form $y = X\beta + \varepsilon$ highlights that OLS still solves $(X^\top X)\hat{\beta} = X^\top y$ provided the design matrix has full rank.

## 3.2 F-statistic and Interpretation

The model-wide F-test compares explained to unexplained variance:

$$F = \frac{(\text{TSS} - \text{RSS})/p}{\text{RSS}/(n - p - 1)}$$

A large F-statistic (small p-value) means the predictor group explains far more variation than an intercept-only baseline. If the F-test is weak, the entire feature set fails to beat simply predicting the mean of $y$.

## 3.3 Adjusted $R^2$ and Collinearity

Adjusted $R^2 = 1 - (1 - R^2)\frac{n-1}{n-p-1}$ penalizes models for adding predictors that do not pull their weight. Monitor variance inflation factors (VIFs) to detect multicollinearity; values above 5–10 hint that coefficients may swing wildly because predictors overlap in information.

## 3.4 Example: Multiple Linear Regression Output

Fitting TV, radio, and newspaper spend simultaneously to predict sales yields:

| Statistic | Estimate | Interpretation |
|---|---|---|
| Intercept | 2.94 | Baseline sales when all spends are zero. |
| TV coefficient | 0.0458 | Positive, highly significant effect (p < $10^{-50}$). |
| Radio coefficient | 0.1885 | Strong positive effect (p < $10^{-40}$). |
| Newspaper coefficient | -0.0010 | Not significant ($p \approx 0.86$). |
| $R^2$ | 0.897 | 89.7% of sales variance explained jointly. |
| F-statistic | 570.3 (p < $10^{-90}$) | Model is overwhelmingly significant. |

The combined model lifts $R^2$ from 0.61 (TV only) to 0.90, but the newspaper coefficient's large p-value says it adds noise, not signal. Dropping that variable or gathering better newspaper features would simplify the model without hurting accuracy.

## 3.5 Interview Questions and Answers

1. **Question:** In the model $\hat{y} = 2.9 + 0.046\,\text{TV} + 0.189\,\text{Radio} - 0.001\,\text{Newspaper}$, how do you interpret the radio coefficient?

   **Solution:** Holding TV and newspaper spend fixed, an extra \$1k in radio spend raises expected sales by 0.189k units. Always mention the "holding others constant" piece to avoid implying causal isolation.

2. **Question:** Suppose $n = 60$, $p = 3$, and $R^2 = 0.82$. Compute the adjusted $R^2$.

   **Solution:** $R^2_{\text{adj}} = 1 - (1 - 0.82)\frac{59}{56} \approx 0.81$. The small decrease shows the added predictors still improve fit after accounting for model size.

3. **Question:** How can you detect multicollinearity, and what are two mitigation strategies?

   **Solution:** High VIFs or a large condition number for $X^\top X$ flag collinearity. Collapse redundant features, use domain-informed feature engineering, or apply ridge/elastic net regularization to stabilize coefficients.

4. **Question:** You fit a baseline model with TV and radio predictors ($\text{RSS}_0 = 120$, $p_0 = 2$) and an extended model that adds newspaper ($\text{RSS}_1 = 110$, $p_1 = 3$) on $n = 200$ observations. Conduct the partial F-test for the added variable.

   **Solution:** $F = \frac{(\text{RSS}_0 - \text{RSS}_1)/(p_1 - p_0)}{\text{RSS}_1/(n - p_1 - 1)} = \frac{10/1}{110/196} \approx 17.82$. Compare to an $F_{1,196}$ distribution; the large value yields a tiny p-value, supporting inclusion of the newspaper predictor.

5. **Question:** Why might you standardize predictors before fitting an MLR model, and how does that affect coefficient interpretation?

   **Solution:** Standardizing removes scale differences, making optimization numerically stable and coefficient magnitudes comparable. After standardization, each coefficient reflects the change in $y$ for a one-standard-deviation increase in its predictor, holding others fixed.

# 4 Logistic Regression

## 4.1 Problem Setup

Binary logistic regression models the log-odds of a "success" as a linear function of the features:

$$\log \frac{\Pr(y_i = 1 \mid \mathbf{x}_i)}{1 - \Pr(y_i = 1 \mid \mathbf{x}_i)} = \beta_0 + \mathbf{x}_i^\top \boldsymbol{\beta}$$

The left side is the log-odds: how many times more likely class 1 is than class 0. Passing the linear predictor through $\sigma(z) = 1/(1 + e^{-z})$ converts log-odds to probabilities $\hat{p}_i$. Predict class 1 whenever $\hat{p}_i$ exceeds a threshold $\tau$ (0.5 by default, but business costs often demand another value).

## 4.2   Estimation and Decision Rules

Parameters maximize the log-likelihood (equivalently, minimize cross-entropy loss)

$$\ell(\boldsymbol{\beta}) = \sum_{i=1}^{n} [y_i \log \hat{p}_i + (1 - y_i) \log(1 - \hat{p}_i)]$$

via iteratively reweighted least squares (Newton's method) or gradient-based solvers. Regularization (L1, L2, elastic net) shrinks coefficients to curb overfitting, especially when predictors are correlated or outnumber observations. Tune thresholds, class weights, or calibration curves so that predicted probabilities line up with real-world costs.

## 4.3   Interpreting Coefficients

Coefficient $\beta_j$ measures the change in log-odds caused by a one-unit increase in $x_{ij}$ while keeping other predictors constant. Exponentiating yields an odds ratio $e^{\beta_j}$: values above 1 increase the odds of the event, values below 1 decrease them. Wald tests or likelihood-ratio tests show whether each predictor adds meaningful signal.

## 4.4   Model Diagnostics and Metrics

Summarize predictions with the confusion matrix:

|                 | Predicted Positive | Predicted Negative |
|-----------------|--------------------|--------------------|
| Actual Positive | True Positive (TP) | False Negative (FN) |
| Actual Negative | False Positive (FP) | True Negative (TN) |

From these counts compute metrics:

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}, \quad \text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}, \quad \text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$\text{Specificity} = \frac{\text{TN}}{\text{TN} + \text{FP}}, \qquad F_1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

Accuracy summarizes overall hits; precision focuses on how many predicted positives were correct; recall (a.k.a. sensitivity) tracks how many actual positives we caught; specificity focuses on the negative class. $F_1$ balances precision and recall. ROC and precision-recall curves visualize how these quantities change with the threshold, and calibration plots verify that a predicted 0.8 probability really corresponds to 80% positives.

## 4.5  Handling Class Imbalance

With rare events, a model can achieve high accuracy by always predicting the majority class. Counter this by resampling (oversample the minority or undersample the majority), generating synthetic examples (SMOTE), using class-weighted losses so mistakes on the minority cost more, or moving the classification threshold toward the minority class to emphasize recall.

## 4.6  Example: Email Spam Detection

On a spam dataset, a logistic regression summary might look like: Coefficients spell

| Feature | Coefficient | Interpretation |
|---|---|---|
| Intercept | -1.92 | Base log-odds favor "ham" when no spam tokens appear. |
| ``free'' count | 0.78 | Each occurrence multiplies spam odds by $e^{0.78} \approx 2.18$. |
| ``money'' count | 0.55 | Positive association with spam likelihood. |
| email length | -0.21 | Longer emails slightly decrease spam odds. |
| ROC-AUC | 0.947 | Excellent separability between spam and ham. |

out how certain words nudge emails toward spam or ham. Pair coefficient reading with ROC/PR curves (for ranking quality) and calibration plots (for probability trustworthiness).

## 4.7  Interview Questions and Answers

1. **Question:** How does logistic regression differ from linear regression in terms of model output and optimization objective?

**Solution:** Logistic regression outputs probabilities via the sigmoid link and optimizes log-likelihood, while linear regression outputs unbounded numbers and minimizes squared error. The probabilistic view unlocks likelihood-based inference and probability-aware metrics.

2. **Question:** Interpret a coefficient $\beta_j = 0.4$ in a standardized logistic regression model.

   **Solution:** A one-standard-deviation increase in $x_j$ multiplies the odds of the positive class by $e^{0.4} \approx 1.49$ (roughly a 49% increase), assuming other predictors stay fixed.

3. **Question:** You observe 90% accuracy but only 40% recall on the positive class. What actions do you take?

   **Solution:** The model is missing positives. Lower the decision threshold, reweight classes, oversample the minority class, or tune directly on recall/$F_1$ until the business trade-off is acceptable.

4. **Question:** Why is feature scaling important for penalized logistic regression?

   **Solution:** L1/L2 penalties shrink coefficients based on magnitude. Without scaling, features measured in large units get penalized more heavily, skewing feature selection. Standardization lets the penalty treat predictors evenly.

5. **Question:** When would you apply probability calibration methods like Platt scaling?

   **Solution:** Use calibration when the score 0.8 corresponds to positives only 60% of the time. Well-calibrated probabilities support risk-based decisions (e.g., how aggressively to filter emails or approve loans).

# 5 Multiple Logistic Regression & Multiclass Classification

## 5.1 Multiple Logistic Regression

Multiple logistic regression keeps the binary target but adds several predictors at once:

$$\Pr(y_i = 1 \mid \mathbf{x}_i) = \sigma\left(\beta_0 + \sum_{j=1}^{p} \beta_j x_{ij}\right)$$

Interpretations still use odds ratios, but now every statement must include "holding the other predictors constant." Check for collinearity with VIFs and add interaction terms when two features only matter together (e.g., marketing spend and customer tenure).

## 5.2 Feature Engineering and Regularization

Indicator variables for categories, standardized numeric features, and interaction terms often unlock signal that simple models miss. Regularization matters: L1 (lasso) zeroes out weak predictors, L2 (ridge) stabilizes correlated features, and elastic net blends the two. Tune the penalty via cross-validation.

## 5.3 Multiclass Logistic Regression

For $K > 2$ classes, multinomial logistic regression uses the softmax function:

$$\Pr(y_i = k \mid \mathbf{x}_i) = \frac{\exp(\mathbf{x}_i^\top \boldsymbol{\beta}_k)}{\sum_{j=1}^{K} \exp(\mathbf{x}_i^\top \boldsymbol{\beta}_j)}, \qquad k = 1, \ldots, K$$

Training minimizes multiclass cross-entropy. Each coefficient vector $\boldsymbol{\beta}_k$ compares class $k$ to a reference class, so interpret coefficients as "how much more (or less) likely class $k$ is than the baseline when this feature increases."

## 5.4 One-vs-Rest and One-vs-One Strategies

Instead of softmax, you can break the problem into binaries. One-vs-rest (OvR) trains $K$ classifiers, each separating one class from the others. One-vs-one (OvO) trains a classifier for every class pair and uses voting at prediction time. OvR scales well but struggles with heavy imbalance; OvO can be more accurate but needs $\frac{K(K-1)}{2}$ models.

## 5.5 Evaluation for Multiclass Problems

Use a $K \times K$ confusion matrix to inspect which classes are confused. Calculate per-class precision/recall, then macro-average (treat each class equally) to highlight minority performance or micro-average (pool all TP/FP/FN before computing) to summarize overall accuracy when class frequencies differ.

## 5.6 Interview Questions and Answers

1. **Question:** Differentiate multiple logistic regression from simple logistic regression.

   **Solution:** Simple logistic regression uses one predictor; multiple logistic regression adds many, enabling "all else equal" statements but requiring checks for multicollinearity and interaction effects.

2. **Question:** Explain the one-vs-rest strategy and when it may struggle.

   **Solution:** OvR trains one classifier per class against the rest. It can struggle when classes overlap heavily or when one class is extremely rare because that classifier sees a lopsided dataset.

3. **Question:** How do macro-averaged and micro-averaged $F_1$ scores differ?

   **Solution:** Macro $F_1$ averages the per-class $F_1$ values, giving each class equal weight. Micro $F_1$ sums all TP/FP/FN first, then computes one $F_1$, so frequent classes dominate the score.

4. **Question:** Why should features be standardized before fitting k-NN or SVM classifiers for multiclass tasks?

   **Solution:** k-NN and SVMs rely on distance calculations; without scaling, features with large numeric ranges dominate the distance and distort decision boundaries.

5. **Question:** How would you evaluate a multiclass classifier when one class is rare yet critical?

   **Solution:** Inspect the $K \times K$ confusion matrix, focus on per-class recall for the rare class, and adjust thresholds or class weights until that class meets the required recall without exploding false positives.

# 6 Linear Discriminant Analysis

## 6.1 Motivation and Assumptions

Linear Discriminant Analysis (LDA) assumes that each class follows a multivariate Gaussian with its own mean vector $\boldsymbol{\mu}_k$ but all classes share the same covariance matrix $\Sigma$. In plain terms: each class forms an ellipsoid cloud with identical shape but different center. These assumptions lead to linear decision boundaries between classes.

## 6.2 Derivation of the Discriminant Function

Bayes' rule yields class posterior probabilities

$$\Pr(y = k \mid \mathbf{x}) \propto \pi_k \cdot \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k, \Sigma),$$

where $\pi_k$ is the prior probability of class $k$. Taking logs and discarding constants leads to a linear discriminant score:

$$\delta_k(\mathbf{x}) = \mathbf{x}^\top \Sigma^{-1} \boldsymbol{\mu}_k - \frac{1}{2} \boldsymbol{\mu}_k^\top \Sigma^{-1} \boldsymbol{\mu}_k + \log \pi_k.$$

Assign $\mathbf{x}$ to the class with the largest $\delta_k(\mathbf{x})$. The resulting boundary between any two classes $k$ and $j$ is a hyperplane.

## 6.3 Parameter Estimation

Estimate class priors via sample proportions $\hat{\pi}_k = n_k/n$. Sample means follow $\hat{\boldsymbol{\mu}}_k = \frac{1}{n_k} \sum_{i:y_i=k} \mathbf{x}_i$. The pooled covariance

$$\hat{\Sigma} = \frac{1}{n - K} \sum_{k=1}^{K} \sum_{i:y_i=k} (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_k)(\mathbf{x}_i - \hat{\boldsymbol{\mu}}_k)^\top$$

captures how points vary within classes. LDA is data efficient because it pools covariance information, but the shared-covariance assumption makes it sensitive to outliers and unequal spreads.

## 6.4 Dimensionality Reduction via Discriminant Axes

For $K$ classes, LDA finds up to $K-1$ discriminant directions that maximize between-class variance relative to within-class variance. Plotting data along the first few discriminant axes makes class separation easy to visualize and yields compact features for downstream models (e.g., run k-NN on the discriminant scores).

## 6.5 Comparison with Logistic Regression

Both LDA and logistic regression create linear boundaries. Logistic regression models $\Pr(y \mid \mathbf{x})$ directly and makes minimal distributional assumptions. LDA models $\Pr(\mathbf{x} \mid y)$ and can shine when the Gaussian assumption is reasonable and data are scarce because it borrows strength across the shared covariance. If covariances differ greatly or distributions are messy, logistic regression is typically safer.

## 6.6 Example: Iris Species Classification

Applying LDA to the Iris dataset (four measurements per flower) produces: Plotting the

| Statistic | Value | Interpretation |
| --- | --- | --- |
| Training accuracy | 0.973 | Strong in-sample separation along discriminant axes. |
| First discriminant variance ratio | 0.991 | Nearly all separability captured in a single axis. |
| Mean vector difference (Setosa vs. Versicolor) | Large in sepal length | Confirms visual separation. |
| Pooled covariance determinant | 0.14 | Shared covariance assumption reasonable. |

first discriminant axis separates Setosa immediately, while the second axis helps split Versicolor and Virginica—an easy visual talking point.

## 6.7 Interview Questions and Answers

1. **Question:** What assumptions does LDA make about the data generating process?

   **Solution:** Each class is Gaussian with its own mean and a shared covariance matrix, and class priors are fixed. If covariances differ, switch to QDA or regularized versions.

2. **Question:** How do you interpret the discriminant function $\delta_k(\mathbf{x})$?

   **Solution:** $\delta_k(\mathbf{x})$ is a score: large when $\mathbf{x}$ is close to class $k$'s mean and when that class is common (via $\log \pi_k$). Assign the label with the largest score.

3. **Question:** When might you prefer LDA over logistic regression?

   **Solution:** When classes look roughly Gaussian with similar spreads and you have limited data, LDA's shared covariance lowers variance and often outperforms logistic regression.

4. **Question:** How do you handle LDA when the number of predictors exceeds observations?

   **Solution:** The covariance estimate becomes singular. Use regularized LDA (shrink toward the identity), reduce dimensionality with PCA first, or adopt penalized discriminant analysis.

5. **Question:** How would you evaluate whether the equal covariance assumption holds?

**Solution:** Compare sample covariances, run Box's M test, or inspect spread along discriminant axes. Large differences say "try QDA or another heteroskedastic method."

# 7 Quadratic Discriminant Analysis

## 7.1 Motivation and Assumptions

Quadratic Discriminant Analysis (QDA) drops LDA's shared-covariance assumption. Each class $k$ gets its own covariance matrix $\Sigma_k$, so the Gaussian "ellipse" can stretch or rotate differently per class. Decision boundaries become curved surfaces, letting the model follow nonlinear separations.

## 7.2 Discriminant Function

Starting from Bayes' rule with class-specific covariances,

$$\Pr(y = k \mid \mathbf{x}) \propto \pi_k \cdot \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k, \Sigma_k).$$

Taking logs yields

$$\delta_k(\mathbf{x}) = -\frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^\top \Sigma_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) + \log \pi_k,$$

which is quadratic in $\mathbf{x}$. Predict by choosing the class with the largest $\delta_k(\mathbf{x})$.

## 7.3 Parameter Estimation and Complexity

Class priors and means match LDA: $\hat{\pi}_k = n_k/n$, $\hat{\boldsymbol{\mu}}_k = \frac{1}{n_k} \sum_{i:y_i=k} \mathbf{x}_i$. Covariance matrices use class-specific estimates

$$\hat{\Sigma}_k = \frac{1}{n_k - 1} \sum_{i:y_i=k} (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_k)(\mathbf{x}_i - \hat{\boldsymbol{\mu}}_k)^\top.$$

Estimating a full covariance per class introduces $\mathcal{O}(p^2 K)$ parameters, so QDA needs far more data than LDA to remain stable. Regularized QDA shrinks each $\hat{\Sigma}_k$ toward a diagonal or pooled version to tame variance.

## 7.4 When to Prefer QDA over LDA

Reach for QDA when classes have noticeably different spreads or the Bayes-optimal boundary is curved. Always validate on held-out data: the extra flexibility produces higher variance, so if covariances are actually similar, LDA may outperform.

## 7.5 Example: Credit Risk Segmentation

Applying QDA to a credit dataset with features like income, debt-to-income ratio, and credit score: After tuning the shrinkage parameter via cross-validation, the curved bound-

| Statistic | Value | Interpretation |
|---|---|---|
| Validation accuracy | 0.842 | QDA outperforms LDA (0.803) due to curved boundaries. |
| Average log determinant difference | 1.27 | Covariance volume differs meaningfully by class. |
| Regularization parameter | 0.05 | Mild shrinkage stabilizes covariance estimates. |
| Misclassified high-risk loans | Reduced by 18% | Better capture of heteroskedastic patterns. |

aries reduced costly high-risk false negatives compared with LDA.

## 7.6 Interview Questions and Answers

1. **Question:** How do the decision boundaries of QDA differ from LDA?

   **Solution:** QDA yields quadratic (curved) boundaries thanks to class-specific co-variances; LDA keeps them linear because the covariance is shared.

2. **Question:** What data conditions motivate QDA despite its higher variance?

   **Solution:** Use QDA when each class clearly has a different spread or covariance structure and you have enough samples per class to estimate those covariances reliably.

3. **Question:** How can you regularize QDA when sample size is limited relative to predictors?

   **Solution:** Shrink each $\hat{\Sigma}_k$ toward the identity or the pooled covariance (convex combination), assume diagonal covariances, or reduce dimensionality with PCA first.

4. **Question:** Compare QDA with logistic regression in terms of assumptions and flexibility.

   **Solution:** QDA models $\Pr(\mathbf{x} \mid y)$ with Gaussian densities and class-specific covariances, so it captures curved boundaries. Logistic regression models $\Pr(y \mid \mathbf{x})$ with linear log-odds; it is more robust but less flexible unless you engineer nonlinear features.

5. **Question:** How would you evaluate whether QDA is overfitting?

   **Solution:** Compare training vs. validation accuracy, inspect fold-to-fold variation in confusion matrices, and benchmark against LDA or regularized QDA. If validation metrics stagnate while training accuracy skyrockets, QDA is overfitting.

# 8 Resampling and Regularization

## 8.1 Cross-Validation Strategies

Resampling estimates how a model will behave on unseen data by repeatedly splitting the dataset. Common strategies:

- **Holdout**: One train/validation split; quick sanity check but high variance because results depend on that single split.

- $k$**-fold cross-validation**: Split into $k$ folds, train on $k-1$, validate on the remaining fold, and average. Provides a good bias-variance trade-off.

- **Stratified $k$-fold**: Keeps class ratios intact within each fold—crucial when classes are imbalanced.

- **Leave-one-out (LOOCV)**: $k = n$. Very low bias, very high variance and compute.

- **Time-series splits**: Only roll forward in time so the model never "sees the future."

Use cross-validation to compare models, tune hyperparameters, and estimate generalization error before final training on all data.

## 8.2 Bootstrap for Uncertainty

Bootstrap sampling draws datasets of size $n$ with replacement to approximate the sampling distribution of a statistic. For model assessment:

1. Draw $B$ bootstrap samples of size $n$.

2. Train the model on each sample and evaluate on the out-of-bag observations.

3. Aggregate metrics (mean, standard deviation) to quantify variability and construct confidence intervals.

Bootstrap shines when the math for standard errors is messy or unknown; the empirical distribution from resamples gives the uncertainty directly.

## 8.3   Regularization Overview

Regularization adds penalty terms to the loss function, trading bias for reduced variance to prevent overfitting:

$$\text{Loss}_{\text{reg}} = \text{Loss}_{\text{empirical}} + \lambda \cdot \Omega(\boldsymbol{\beta}),$$

where $\lambda$ controls penalty strength and $\Omega$ encodes the desired constraint.

- **Ridge (L2)**: $\Omega(\boldsymbol{\beta}) = \|\boldsymbol{\beta}\|_2^2$. Shrinks coefficients smoothly toward zero, stabilizing correlated predictors.

- **Lasso (L1)**: $\Omega(\boldsymbol{\beta}) = \|\boldsymbol{\beta}\|_1$. Drives some coefficients exactly to zero, acting as built-in feature selection.

- **Elastic Net**: Mixes L1 and L2 to get sparsity plus stability, great when predictors come in correlated groups.

- **Other models**: Apply similar penalties to logistic regression, LDA/QDA covariances, or neural nets to keep decision boundaries from getting too wiggly.

Tune $\lambda$ with cross-validation or pick using information criteria (AIC, BIC) when computation is tight.

## 8.4   Bias-Variance Trade-off

Resampling lets you watch training vs. validation error as the model becomes more flexible. Regularization increases bias slightly but slashes variance, often lowering validation error overall. Plot CV error against $\lambda$ to see where you move from underfitting to overfitting.

## 8.5 Model Selection Workflow

1. Define candidate models and features.

2. Standardize or normalize predictors where required (e.g., before applying L1/L2 penalties).

3. Use stratified $k$-fold cross-validation to tune hyperparameters such as $\lambda$, penalty mix, or tree depth.

4. Evaluate final model on a held-out test set to obtain an unbiased estimate of performance.

5. Retrain on the full dataset with chosen hyperparameters before deployment.

## 8.6 Case Study: Ridge Regression

Ridge regression extends linear regression by adding an L2 penalty:

$$\min_{\boldsymbol{\beta}} \sum_{i=1}^{n} (y_i - \beta_0 - \mathbf{x}_i^{\top}\boldsymbol{\beta})^2 + \lambda\|\boldsymbol{\beta}\|_2^2.$$

The normal equations become $(X^{\top}X + \lambda I)\boldsymbol{\beta} = X^{\top}y$, guaranteeing an invertible system even when predictors are collinear or $p > n$. Ridge shrinks coefficients smoothly while keeping all predictors in the model. Useful diagnostics:

- **Coefficient paths**: Show $\beta_j$ vs. $\lambda$ to visualize how each feature shrinks.

- **Effective degrees of freedom**: $\mathrm{df}(\lambda) = \mathrm{tr}\left(X(X^{\top}X + \lambda I)^{-1}X^{\top}\right)$; larger $\lambda$ means a simpler model.

- **Cross-validated error**: Pick the $\lambda$ with the lowest validation error or use the one-standard-error rule for a slightly simpler yet still strong model.

Standardize predictors before fitting to ensure the penalty treats features on comparable scales.

## 8.7 Interview Questions and Answers

1. **Question:** Why is stratified cross-validation important for classification?

   **Solution:** It preserves class proportions in each fold, yielding more stable estimates for metrics like recall/precision, especially when classes are imbalanced.

2. **Question:** Compare ridge and lasso regularization in terms of their effect on coefficients.

   **Solution:** Ridge shrinks coefficients toward zero but never eliminates them; lasso can zero them out entirely, creating sparse models. Elastic net blends the two behaviors.

3. **Question:** How would you choose the regularization strength $\lambda$?

   **Solution:** Evaluate a grid (or use algorithms like coordinate descent) within cross-validation, selecting the $\lambda$ that minimizes validation error or lies within one standard error of the minimum to favor simpler models.

4. **Question:** When is the bootstrap preferable to $k$-fold cross-validation?

   **Solution:** Bootstrap is useful for estimating parameter uncertainty or when the dataset is too small for reliable folds; it provides confidence intervals for metrics by resampling with replacement.

5. **Question:** How does regularization mitigate multicollinearity?

   **Solution:** Penalties like ridge constrain coefficient magnitude, reducing sensitivity to correlated predictors and stabilizing estimates, leading to lower variance in predictions.

# 9 Decision Trees

## 9.1 Structure and Intuition

Decision trees repeatedly split the feature space into rectangles (axis-aligned regions). Each internal node asks a question such as "is feature$_j \leq t$?" and routes the observation left or right; each leaf stores the prediction (a mean value or a class probability). Because splits can occur anywhere, trees capture nonlinear relationships and interactions automatically.

## 9.2 Recursive Binary Splitting (Top-Down Greedy)

The CART algorithm grows a tree top-down. At each node it considers every feature-threshold pair, picks the split with the biggest impurity reduction, and recurses on the children. This greedy strategy is fast even though it does not peek ahead. Stop when nodes are pure enough or too small, or when you reach a depth limit.

## 9.3 Regression Trees and RSS

Regression trees predict continuous outcomes by minimizing the residual sum of squares within each node. For a node with observations $\mathcal{D}$, the impurity is

$$\text{RSS}(\mathcal{D}) = \sum_{(x_i, y_i) \in \mathcal{D}} (y_i - \bar{y}_{\mathcal{D}})^2,$$

where $\bar{y}_{\mathcal{D}}$ is the mean response in the node. Candidate splits are scored by the drop in RSS between the parent and its children. Leaves store the average response of training points that fall into the region.

## 9.4 Classification Trees and Error Metrics

Classification trees rely on class proportions in each node. Common impurity measures include:

- **Classification error rate**: $1 - \max_k p_k$; easy to interpret but less sensitive to class changes.

- **Gini impurity**: $1 - \sum_k p_k^2$; smooth and fast to compute.

- **Entropy**: $-\sum_k p_k \log p_k$; strongly penalizes rare-class mixtures.

During pruning or model comparison, classification error rate provides a high-level accuracy view, while Gini/entropy guide split decisions.

## 9.5 Pruning and Regularization

Fully grown trees memorize the training data. Control complexity by pre-pruning (limit depth, minimum samples per split/leaf, max leaves) or by post-pruning (grow a large tree, then prune using cost-complexity pruning with parameter $\alpha$). Tune those settings via cross-validation. Ensembles like random forests or gradient boosting further tame variance by averaging many trees.

## 9.6 Handling Continuous and Categorical Features

Continuous features are split by thresholds; categorical features require either one-hot encoding or subset splits. When categories are high-cardinality, group infrequent levels or use target encoding carefully to avoid leakage.

## 9.7  Strengths and Limitations

- **Strengths**: Easy to visualize, handle numeric and categorical data, no scaling, capture nonlinear relationships.

- **Limitations**: High variance, piecewise-constant predictions, biased toward features with many split points, sensitive to small data changes.

Bagging, random forests, and boosting mitigate most limitations by averaging or sequencing many trees.

## 9.8  Ensemble Extensions: Bagging, Random Forests, Boosting

- **Bagging (Bootstrap Aggregation)**: Train many deep trees on bootstrap samples and average their predictions. Reduces variance and improves stability.

- **Random Forests**: Bagging plus random feature subsampling at each split, decorrelating trees and yielding better performance than plain bagging.

- **Boosting**: Sequentially fit shallow trees to residuals or gradients (e.g., AdaBoost, Gradient Boosted Trees). Each tree focuses on examples the previous ensemble mispredicted, reducing bias with careful learning-rate control.

Hyperparameters (number of trees, depth, learning rate) are tuned via cross-validation to balance bias and variance.

## 9.9  Example: Customer Churn Classification

Training a depth-limited tree on customer behavior metrics yields: Tree visualization

| Statistic | Value | | Interpretation |
|---|---|---|---|
| Max depth | 5 | | Controls model complexity. |
| Validation accuracy | 0.812 | | Competitive with logistic baseline (0.789). |
| Top feature importances | Tenure, calls, charges | Support Monthly | Key churn drivers. |
| Pruning parameter $\alpha$ | 0.01 | | Selected via cross-validation to reduce overfitting. |

highlights churn segments (e.g., short-tenure customers with high support call counts).

## 9.10 Interview Questions and Answers

1. **Question:** How does Gini impurity differ from entropy when splitting nodes?

   **Solution:** Both measure class mix; Gini is simpler $(1-\sum p_k^2)$, entropy is $-\sum p_k \log p_k$. They often select similar splits, but entropy penalizes rare classes slightly more.

2. **Question:** Why do decision trees tend to overfit, and how can you mitigate that?

   **Solution:** Trees can memorize training data by growing deep, low-sample leaves. Mitigate using pruning, depth/min-samples constraints, or ensemble methods (bagging/boosting) that reduce variance.

3. **Question:** What bias might arise when splitting on categorical variables with many levels?

   **Solution:** Features with many categories can appear more informative because they create numerous small pure splits, leading to biased feature selection. Mitigate by collapsing categories, using feature hashing, or applying corrected splitting criteria.

4. **Question:** When would you prefer a single decision tree over a random forest?

   **Solution:** When interpretability is paramount and the dataset is small or low-noise. Otherwise, random forests generally outperform single trees by averaging variance.

5. **Question:** How do you evaluate feature importance in tree models?

   **Solution:** Sum the impurity reduction contributed by each feature across splits, or use permutation importance by measuring performance drop when feature values are shuffled.

6. **Question:** Explain the top-down greedy nature of recursive binary splitting and one consequence of its greediness.

   **Solution:** The algorithm picks the best immediate split without considering future splits; this can miss globally optimal trees. Pruning and cross-validation help correct overfitting caused by greedy decisions.

7. **Question:** Contrast bagging, random forests, and boosting for tree-based models.

   **Solution:** Bagging averages many independent trees to lower variance; random forests add feature subsampling to further decorrelate trees; boosting fits trees sequentially to residuals to reduce bias. Choice depends on variance vs. bias trade-offs and computation.

# 10 Support Vector Machines

## 10.1 Maximal Margin Classifier

When classes are perfectly separable, Support Vector Machines (SVMs) look for the hyperplane with the largest possible margin (distance to the closest points). For labeled data $(\mathbf{x}_i, y_i)$ with $y_i \in \{-1, +1\}$, the hard-margin problem is

$$\max_{\mathbf{w}, b} \frac{2}{\|\mathbf{w}\|} \quad \text{subject to} \quad y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 \ \forall i.$$

The support vectors are the points lying exactly on the margin boundaries and fully determine the optimal separating hyperplane.

## 10.2 Hyperplanes and Geometric View

A hyperplane in $p$ dimensions is $\{\mathbf{x} : \mathbf{w}^\top \mathbf{x} + b = 0\}$. $\mathbf{w}$ controls orientation and $b/\|\mathbf{w}\|$ sets the offset from the origin. The signed distance from a point to the hyperplane is $\frac{\mathbf{w}^\top \mathbf{x} + b}{\|\mathbf{w}\|}$, so maximizing the margin equals minimizing $\|\mathbf{w}\|$ subject to the classification constraints.

## 10.3 Soft Margin Support Vector Classifier

Real-world data rarely are perfectly separable. The soft margin SVM introduces slack variables $\xi_i \geq 0$ to allow margin violations:

$$\min_{\mathbf{w}, b, \boldsymbol{\xi}} \frac{1}{2}\|\mathbf{w}\|^2 + C \sum_{i=1}^{n} \xi_i \quad \text{s.t.} \quad y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i.$$

The penalty parameter $C$ controls the trade-off between margin width and misclassification tolerance. In the dual formulation, the optimization depends only on inner products $\mathbf{x}_i^\top \mathbf{x}_j$, paving the way for kernel methods.

## 10.4 Kernel Trick and Nonlinear Decision Boundaries

Replacing inner products with kernel functions $K(\mathbf{x}_i, \mathbf{x}_j)$ implicitly maps points into a high-dimensional feature space without ever writing down the mapping.

- **Linear kernel**: recovers the soft margin classifier.

- **Polynomial kernel**: $K(\mathbf{x}, \mathbf{z}) = (\gamma \mathbf{x}^\top \mathbf{z} + r)^d$ captures feature interactions up to degree $d$.

- **Radial Basis Function (RBF)**: $K(\mathbf{x}, \mathbf{z}) = \exp(-\gamma \|\mathbf{x} - \mathbf{z}\|^2)$ yields flexible, localized decision boundaries.

Kernel choice and hyperparameters $(C, \gamma, d)$ are tuned via cross-validation.

## 10.5 Model Selection and Practical Considerations

- **Feature scaling**: Standardize inputs; otherwise the Euclidean margin gives more weight to large-scale features.

- **Class imbalance**: Use class-specific $C$ values or adjust thresholds so rare classes receive higher recall.

- **SVM regression (SVR)**: Uses an $\varepsilon$-insensitive loss so deviations within $\varepsilon$ are ignored; same margin logic applies.

- **Computational cost**: Kernel SVM training scales poorly with sample size; rely on linear SVM solvers or approximate kernels for very large datasets.

## 10.6 Example: RBF-SVM for Handwritten Digits

Using standardized pixel intensities to classify digits 0–9: Grid search over $(C, \gamma)$ with

| Setting | Value | Notes |
| --- | --- | --- |
| Kernel | RBF ($\gamma = 0.015$) | Captures nonlinear digit boundaries. |
| Penalty $C$ | 10 | Balances margin and errors. |
| Cross-validated accuracy | 0.983 | Beats logistic baseline (0.944). |
| Support vectors | 3,120 ($\approx$ 30% of data) | Points defining the decision surface. |

stratified folds identified the best trade-off between bias and variance.

## 10.7 Interview Questions and Answers

1. **Question:** Define the maximal margin classifier and explain when it is feasible.

   **Solution:** It is the hard-margin SVM that maximizes the geometric margin subject to perfect separation; feasible only when data are linearly separable.

2. **Question:** How does the parameter $C$ influence the soft margin SVM?

   **Solution:** Large $C$ penalizes misclassifications heavily, favoring narrower margins and lower bias but higher variance; small $C$ allows wider margins with more violations, increasing bias and robustness to noise.

3. **Question:** What role do support vectors play in prediction?

   **Solution:** Only support vectors have non-zero dual coefficients; predictions depend on their weighted contribution via $K(\mathbf{x}_i, \mathbf{x})$, making SVMs resilient to outliers far from the margin.

4. **Question:** Compare linear, polynomial, and RBF kernels in terms of flexibility.

   **Solution:** Linear kernels build linear boundaries; polynomial kernels add global interactions controlled by degree; RBF kernels generate highly flexible, localized boundaries. Choice depends on domain knowledge and data complexity.

5. **Question:** How would you diagnose overfitting in an SVM model?

   **Solution:** Monitor cross-validation gap versus training accuracy, inspect number of support vectors (very high proportion may signal overfitting), and visualize decision boundaries or learning curves across $C$ and kernel parameters.

# 11 Unsupervised Learning

## 11.1 Overview

Unsupervised learning searches for structure in unlabeled data. Two common goals: dimensionality reduction (compress features while keeping the big patterns) and clustering (group similar observations). Because there is no ground-truth label, we rely on internal metrics, visualization, or downstream tasks to judge quality.

## 11.2 Principal Component Analysis (PCA)

PCA finds orthogonal directions (principal components) that capture the largest variance. With centered data matrix $X$, solve

$$\left( \frac{1}{n-1} X^\top X \right) \mathbf{v}_k = \lambda_k \mathbf{v}_k.$$

Eigenvalue $\lambda_k$ tells you how much variance component $k$ explains, and the projection scores are $X\mathbf{v}_k$. Standardize features first so one column does not dominate simply because of its scale.

## 11.3 Clustering Fundamentals

Clustering groups points based on a similarity metric (Euclidean distance, cosine similarity, etc.). Keep in mind:

- Distance choice shapes the clusters (cosine for sparse text, Euclidean for continuous features, etc.).

- Feature scaling changes those distances, so standardize unless units already match.

- Use internal validity metrics (silhouette, Davies-Bouldin) or stability checks to see if clusters are meaningful.

## 11.4 k-means Clustering

k-means partitions observations into $K$ clusters by minimizing within-cluster sum of squares:

$$\min_{\{\mathcal{C}_k\}_{k=1}^K} \sum_{k=1}^K \sum_{\mathbf{x}_i \in \mathcal{C}_k} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|^2,$$

where centroid $\boldsymbol{\mu}_k$ is the mean of cluster $\mathcal{C}_k$. The standard algorithm alternates between assigning points to closest centroids and updating centroids. Initialization matters; k-means++ selects diverse seeds to improve convergence. The elbow method or silhouette analysis helps choose $K$.

## 11.5 Hierarchical Clustering

Hierarchical methods create nested clusterings visualized via dendrograms:

- **Agglomerative (bottom-up)**: Start with singletons; iteratively merge the closest clusters based on linkage (single, complete, average, Ward).

- **Divisive (top-down)**: Start with all points; recursively split clusters.

Ward linkage minimizes variance increase and behaves similarly to k-means. Cutting the dendrogram at a chosen height yields cluster assignments without pre-specifying $K$.

## 11.6 Example: Customer Segmentation Workflow

1. Standardize behavioral features (spend, frequency, tenure).

2. Apply PCA to retain components explaining 85% of variance (e.g., first three components).

3. Run k-means on PCA scores with $K = 4$; silhouette score 0.48 indicates reasonable separation.

4. Compare against hierarchical clustering (Ward linkage) to validate grouping stability.

5. Profile clusters (high-value loyalists, promotion responders, at-risk, new users) for marketing strategy.

## 11.7 Interview Questions and Answers

1. **Question:** Why is feature scaling important before PCA or k-means?

   **Solution:** Both rely on Euclidean geometry; unscaled features with large variances dominate principal components or cluster assignments, distorting structure.

2. **Question:** How do you choose the number of principal components to retain?

   **Solution:** Inspect the scree plot for an elbow, set a cumulative variance threshold (e.g., 90%), or validate with downstream performance (e.g., classification accuracy using reduced features).

3. **Question:** Compare k-means with hierarchical clustering.

   **Solution:** k-means scales well, assumes spherical clusters, and requires pre-specifying $K$. Hierarchical clustering handles varied shapes, produces dendrograms aiding interpretation, but is more computationally expensive.

4. **Question:** What challenges arise when clusters have different densities?

   **Solution:** k-means struggles; dense clusters dominate centroid placement. Alternatives include DBSCAN or Gaussian Mixture Models, which adapt to varying densities and covariance structures.

5. **Question:** How can you assess cluster quality without labels?

   **Solution:** Use internal metrics (silhouette, Calinski-Harabasz), stability analysis (re-run clustering on bootstrap samples), or evaluate business outcomes (e.g., uplift in targeted campaigns).

# 12   Conclusion

You now have a plain-language tour of the main models from the FreeCodeCamp course: linear and logistic regression (single and multiple), discriminant analysis, trees and ensembles, SVMs, resampling/regularization, and unsupervised standbys such as PCA and clustering. Interviewers usually probe assumptions, diagnostics, and "why this model?" reasoning, so rehearse those talking points as often as you review formulas.

- Re-run the companion notebooks in `data-science/Code/` to connect the prose with live code.

- Build one-page cheat sheets of metrics, assumptions, and failure modes for each family of models.

- Practice walking through an end-to-end project: feature engineering, validation strategy, model comparison, and business translation.

Regular quick reviews—paired with hands-on experimentation—will keep the ideas accessible and interview-ready.