

“

JNOTES

CREATIVE NOTES

>

Logistic Regression is a discriminative approach for classification. In fact unlike generative model that try to model how data is generated ($P(x|y=c)$) logistic regression directly models the probability of class membership (probabilities) given the feature ($P(y=c|x=x)$).

So generative models have the necessity to learn how data was created because for computing the actual probability $P(c|x)$ they use the Bayes rule, so they compute (extreme) $P(x, c) = P(c|x)P(x)$ where $P(x|c)=0$ if data is generated with each class.

The discriminative models instead directly model the posterior distribution $c|x$.

Binary Classification

In generative approach with gaussian distribution having the same covariance matrix (tied covariance) we discovered that the optimal decision boundary is linear and can be written as:

$$\log \frac{P(C=1|x)}{P(C=0|x)} = \log \frac{P(x|C=1)}{P(x|C=0)} + \log \frac{\pi_1}{\pi_0} = w^T x + b = \text{log-odds}$$

log posterior probability log likelihood ratio log prior ratio
 $\frac{P(C=1|x)}{P(C=0|x)}$ $\frac{P(x|C=1)}{P(x|C=0)}$ $\frac{\pi_1}{\pi_0}$

So the log-odds is a linear function of our input feature.

Flipping the Approach with Logistic Regression

With logistic regression we assume that (constant odds ratio) our decision boundary should be linear.

We need to convert this linear score into a probability

if log-odds equal $w^T x + b$

$$\log \frac{P(C=1|x)}{P(C=0|x)} = w^T x + b \Rightarrow \frac{P(C=1|x)}{P(C=0|x)} = e^{w^T x + b} \quad \left. \right\} =$$

if we want to compute $P(C=1|x)$ since we know that $P(C=0|x) = 1 - P(C=1|x)$

$$= P(C=1|x) = e^{w^T x + b} [1 - P(C=1|x)] \\ = P(C=1|x)(1 + e^{w^T x + b}) = e^{w^T x + b} =$$

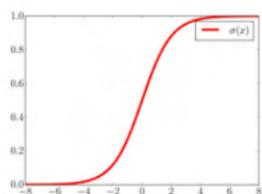
$$= P(C=1|x) = \frac{e^{w^T x + b}}{1 + e^{w^T x + b}} = \frac{1}{1 + e^{-(w^T x + b)}} = \sigma(w^T x + b)$$

$\sigma(w^T x + b)$ is called sigmoid function or logistic function

More is more appropriate to call $P(C=1|x)$ as $P(C=1|x, w, b)$ because we are assuming the value w, b so the probability depends on them.

Sigmoid Function

The sigmoid function $\sigma(x) = \frac{e^x}{1+e^{-x}}$ is the mathematical function that transforms a linear function into a probability, it has several important properties



- it maps any real number to [0,1] [probability]
- Symmetry: $\sigma(-x) = 1 - \sigma(x)$ the function is symmetric around the y axis
- $\frac{d\sigma(x)}{dx} = \sigma(x)(1 - \sigma(x))$

The expression $p(c_i | x; w, b) = \sigma(w^T x + b)$ provides a model to directly calculate the probability that a direct observation x belongs to class c_i ; assuming that the decision boundaries is a linear hyperplane in the feature space.

the decision rule are linear surface orthogonal to w (\perp) with b as the offset from the origin.

So if we knew (w, b) then we could compute the predictive distribution for the class labels $P(c_i | x; w, b)$.

We have seen an **indirect** way of computing (w, b) with a generative model

- The tied covariance Gaussian feature model, combined with application class prior, can be cast as a linear decision rule of the form $w^T x + b \leq 0$

However, in the following we are interested to find an alternative way to **estimate an effective classification rule** that **does not require** an explicit model of the feature vectors distribution

We thus ignore generative models for X and concentrate directly on the form of the class posterior probabilities

Again, we follow a frequentist approach, i.e. compute an estimate for w and b from a set of training samples

Using again likelihood

We start with a labeled dataset:

$$D = \{(x_1, c_1), \dots, (x_n, c_n)\}$$

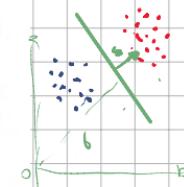
Where:

- Each x_i is a feature vector (sample)
- Each c_i is a class label

We make a critical assumption: the data points are independent and identically distributed (iid) given model parameters. This means:

$$\{ (x_i, c_i) \mid (x_j, c_j) \} \perp \!\!\! \perp \{ (x_i, c_i) \mid \forall i \}$$

So knowing the outcome of one data point doesn't tell us anything about the outcome of another data point, once we know the model parameters $\theta = (w, b)$



b can increase or decrease depending on the cost

So the likelihood function represent the probability of observing our entire dataset given the model parameters

$$L(\theta) = f_{x_1, \dots, x_n, c_1, \dots, c_n}(x_1, \dots, x_n, c_1, \dots, c_n | \theta)$$
$$= \prod_{i=1}^n f_{x_i, c_i}(x_i, c_i | \theta)$$

because of the iid assumption \approx

The main difference from the generative models is in the way we decompose the joint density factor

+

Generative Decomposition: $P(x, c) = P(x|c) P(c)$

Discriminative Decomposition: $P(x, c) = P(c|x) P(x)$

So in our likelihood function, each data point contributes:

$$f_{x_i, c_i | \theta}(x_i, c_i | \theta) = P(c=c_i | X=x_i, \theta) f(x_i)$$

An Analogy to Understand the Difference

Imagine you're a doctor trying to diagnose a disease:

Generative approach: You study healthy people and people with the disease separately. You learn the typical patterns for each group, understanding what's "normal" and what indicates disease. Then when a new patient arrives, you compare their symptoms to both patterns and see which they match better.

Discriminative approach: You directly study the boundary between healthy and sick patients, focusing only on which symptoms best distinguish between the two groups. You don't need to fully understand what "typical" looks like for each

Where:

- $P(c=c_i | X=x_i, \theta)$ is the conditional probability on the class given the features
- $f(x_i)$ is the marginal density

Note: We will ignore $f(x_i)$ it not depend on θ

For simplicity we consider the logarithm of our likelihood function:

$$\log L(\theta) = \log \prod_{i=1}^n f_{x_i, c_i | \theta}(x_i, c_i | \theta) = \sum_{i=1}^n \log P(c=c_i | X=x_i, \theta) + \sum_{i=1}^n \log f(x_i)$$

Since θ only influence the first sum and we're trying to maximize this function, we can focus on:

$$L(\theta) = \sum_{i=1}^n \log P(c=c_i | X=x_i, \theta)$$

$$\theta_{\text{ML}} = \arg \max_{\theta} \log L(\theta)$$

The parameter that best fit the distribution

We now need to express the probability of each class label in terms of model parameters (w, b)

[Considering binary problem]

- For positive class ($c=1$): $y_i = P(c_i=1 | x_i, w, b) = \sigma(w^T x_i + b)$
- For negative class ($c=0$): $1 - y_i = 1 - \sigma(w^T x_i + b) = P(c_i=0 | x_i, w, b) = \sigma(-w^T x_i - b)$
[σ = sigmoid function]

A key insight is that our classification binary problem follow the Bernoulli distribution pattern. The Bernoulli distribution describe a random variable that can take only two values (0 or 1) with some probability p .

$$C|x, w, b \sim \text{Ber}(\sigma(w^T x + b)) = \text{Ber}(y_i)$$

We are saying that given a specific feature vector x and our model parameters w and b

1. The class label C follows a Bernoulli distribution
2. The parameter of the Bernoulli distribution is $\sigma(w^T x + b)$

So

$$\begin{cases} P(C=1|x, w, b) = y \\ P(C=0|x, w, b) = 1-y \end{cases}$$

The compact formula

$$P(C=c|x, w, b) = y^c (1-y)^{1-c}$$

Now our log-likelihood function over all data points becomes:

$$l(w, b) = \log \prod_{i=1}^n P(C_i = c_i | x_i, w, b) = \log \prod_{i=1}^n y_i^{c_i} (1-y_i)^{1-c_i}$$

$$= \sum_{i=1}^n [c_i \log y_i + (1-c_i) \log (1-y_i)]$$

$$\text{Where } y_i = \sigma(w^T x_i + b)$$

We've stated that our goal in logistic regression is to find the parameters w and b that maximize the log-likelihood function:

$$l(w, b) = \sum_{i=1}^n [c_i \log y_i + (1-c_i) \log (1-y_i)]$$

Since optimization algorithm typically minimize rather than maximize we can convert the minimization problem by simply negating the function:

$$S(w, b) = -l(w, b) = \sum_{i=1}^n [-c_i \log y_i - (1-c_i) \log (1-y_i)]$$

Now the expression (looking at a single data point):

$$H(c_i, y_i) = -[c_i \log y_i + (1-c_i) \log (1-y_i)]$$

Per true noise x noise noise (from same or 2 classes noisy)

between the true label c_i and the predicted probability y_i .

What is cross-entropy?

The cross-entropy measure the "distance" between two probability distributions. Given two distributions P and Q over the same set of outcomes the cross-entropy is defined as:

$$H(P, Q) = -E_{P(x)} [\log Q(x)] = -\sum_{x \in S} P(x) \log Q(x)$$

Where:

P is the "true" distribution

Q is our approximation or prediction

The sum is over all possible outcomes x

$E_{P(x)}$ represent the expected value with respect to distribution P

In our case for each sample in our dataset we have:

- The actual outcome (0 or 1)
- Our model guess about how likely they were to pass.

Our goal is to make our guesses match reality as close as possible, this is possible by minimizing the cross-entropy.

- P is the "true" distribution: the empirical of class labels, from the point of view of an observer i (evaluator) who knows the true (the actual) labels

Given the data point (x_i, c_i) the true distribution puts all the probability mass on the observed class:

$$P(C_i=1 | X_i=x_i, \xi) = \begin{cases} 1 & \text{if } c_i=1 \\ 0 & \text{if } c_i=0 \end{cases}$$

$$P(C_i=0 | X_i=x_i, \xi) = \begin{cases} 0 & \text{if } c_i=1 \\ 1 & \text{if } c_i=0 \end{cases}$$

$$\Rightarrow P(C_i=1 | X_i=x_i, \xi) = c_i, P(C_i=0 | X_i=x_i, \xi) = 1 - c_i$$

↓
Bernoulli distribution with parameter c_i :

- Q is our model predicted distribution: this is the distribution given by our model parameters which we call R (Recognizer)

$$Q(C_i=1 | X_i=x_i, R(w, b)) = q_i = \sigma(w^T x_i + b)$$

$$Q(C_i=0 | X_i=x_i, R(w, b)) = 1 - q_i = 1 - \sigma(w^T x_i + b)$$

Logistic regression looks for the **minimizer of the average cross-entropy** between the distribution for the training set labels of an evaluator \mathcal{E} who knows the real label and the distribution for the training set labels as predicted by the model $\mathcal{R}(w, b)$ itself

The cross-entropy is a measure of goodness of the predictions, and the evaluation is performed over the training data itself

[Cross-entropy properties]

- ① it penalizes wrong predictions heavily: if the true label is t and our model gives a small probability to that class ($y = 0.1$) the cross-entropy is very large ($-\log(0.1) = 2.27$)
- ② it rewards correct predictions: with $y=0.5$ the cross-entropy is $(-\log(0.5)) = 0.69$

We are finding the parameters that minimize the coding inefficiency when using our model to predict the true label.

So when our predicted distribution Q exactly matches the true distribution P , ($P=Q$) the cross-entropy is minimum

When we minimize cross-entropy, we're trying to make our model's predictions match the empirical distributions of our training dataset as closely as possible, given the constraints of our model

In the case of logistic regression, these constraints are significant - we're limited to linear decision boundaries. The model can't perfectly match any arbitrary data distribution, but it finds the best linear approximation possible

Now, instead of using tables 0 and 1, we can use -1 and 1

The terms z_i still represent class labels, however for samples of class h_1 we have $z_i = 1$, whereas for samples of class h_0 we have $z_i = -1$:

$$z_i = \begin{cases} 1 & \text{if } c_i = 1 \\ -1 & \text{if } c_i = 0 \end{cases}$$

We can rewrite the cross-entropy following this conversion

$$z_i = 2c_i - 1$$

Now let's rewrite our objective function.

$$H(c_i, y_i) = -[c_i \log y_i + (1-c_i) \log(1-y_i)]$$

Where:

c_i is the class label (0 or 1) of the sample

$y_i = \sigma(w^T x_i + b)$ our predicted probability

$s_i = w^T x_i + b$ is our linear score (logit)

following the formula $z_i = c_i - b \Rightarrow c_i = 1, z_i = 1$

$c_i = 0, z_i = 0$

$$H(c_i, y_i) = \begin{cases} \log(y_i) + (1-y_i) \log(1-y_i) = -\log y_i = -\log \sigma(s_i) & c_i = 1, z_i = 1 \\ \log(1-y_i) + \log(1-\sigma(s_i)) = -\log \sigma(-s_i) & c_i = 0, z_i = -1 \end{cases} \quad \Delta$$

x property

Using $z_i = 1$ we can combine this in $H(c_i, y_i) = -\log \sigma(z_i s_i) \quad \} =$

Now because $\sigma(x) = \frac{1}{1+e^{-x}}$

$$\Rightarrow H(c_i, y_i) = \log(1+e^{-z_i s_i}) \stackrel{s_i = w^T x_i + b}{=} \log(1+e^{-z_i(w^T x_i + b)})$$

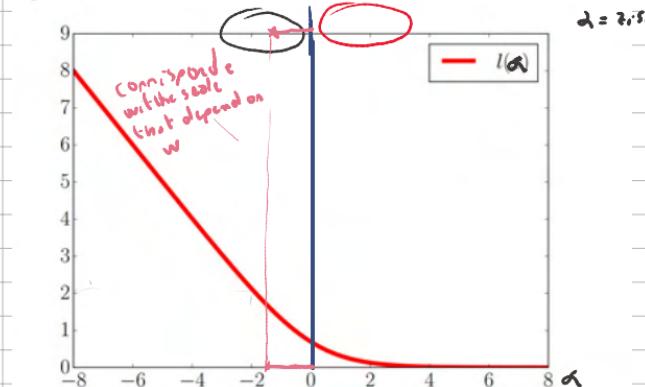
Logistic loss functions

The expression $\log(1+e^{-z_i(w^T x_i + b)})$ is called the logistic loss function often denoted as $l(a; s_i)$ or $l(a; (w^T x_i + b))$

Our entire objective function then becomes:

$$S(w, b) = \sum_{i=1}^n l(a_i; (w^T x_i + b)) = \sum_{i=1}^n \log(1+e^{-z_i(w^T x_i + b)})$$

Our goal is to find the minimizer of $S(w, b)$



Note:

In logistic regression, the linear score $s_i = w^T x_i + b$ has a powerful geometric interpretation:

① Decision boundary: The equation $w^T x + b = 0$ defines an hyperplane in the feature space.

② Orientation: The vector w is orthogonal (perpendicular) to this hyperplane

③ Signel Distance: The value of s_i is proportional to the signed distance from point x_i to the decision boundary.

When we say "Signed distance", we mean:

- Positive values ($s_i > 0$) indicate the point is on the positive side (class h_1) [incorrect]
- Negative values ($s_i < 0$) indicate the point is on the negative side (class h_0) [incorrect]
- The magnitude $|s_i|$ relates to how far the point is from the boundary

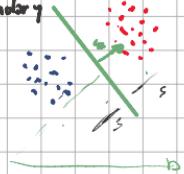
If you remember the score s_i has a direct probabilistic interpretation. It represent the log-ratio of the posterior probabilities:

$$\log \left(\frac{P(C_i=1|x_i=x_i)}{P(C_i=0|x_i=x_i)} \right) = w^T x_i + b = s_i$$

When $s_i=0$ the two class are equally likely

$s_i < 0$ Class h_0 more likely

$s_i > 0$ Class h_1 more likely



The magnitude $|s_i|$ indicate how confident the model is in the predictions

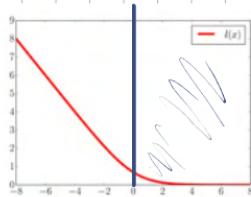
Going back to the logistic loss function $\ell(z; s_i)$, this function measures the cost of our predictions for each samples

① When predictions and actual Class Agree $z_i s_i > 0$

This happens in two cases:

- When $z_i = 1$ (actual class is h_1) and $s_i > 0$ (predicted class h_1)
- $z_i = -1$ (actual class is h_0) and $s_i < 0$ (predicted class h_0)

In this case the cost is relatively low. As $|s_i|$ increase (the point moves further from the boundary to the correct side) the cost decrease exponentially approaching 0 for point far from the boundary



$\ell(z)$ vs. value $\rightarrow 0$, so $S(w, b) \rightarrow 0$ cross-entropy $\rightarrow 0$. We are predicting well

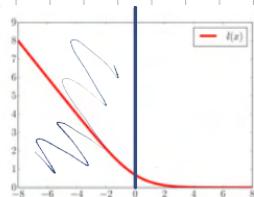
② When predictions and Actual Class Disagree: $z_i s_i < 0$

This happen in two cases:

- When $z_i = 1$ (actual class is h_1) and $s_i < 0$ (predicted class is h_0)

• When $z_i \geq 0$ (actual class is $y_i = 1$) and $s_i > 0$ (predicted class is $\hat{y}_i = 1$)

In this case we pay a much higher cost. As $|s_i|$ increase (the point moves further from the boundary to the wrong side) we pay a cost that increase (asymptotically) linearly with $|s_i|$:



$l(x)$ for value $P = \log(1+b)$ \rightarrow cross-entropy \rightarrow we are predicting bool

(Connection with empirical) risk

Logistic regression fits within a broad statistical learning framework (known as Agnostic, NTO statistics or ANOVA) called Empirical Risk Minimization (ERM). In this framework:

- ① We define a loss function $l(w, x_i, z_i)$ that measures the cost of our prediction for each sample
- ② The empirical risk is the average loss over our training data: $R(\theta) = \frac{1}{m} \sum_i l(w, x_i, z_i)$
- ③ We find the model parameters θ that minimize this empirical risk

For logistic regression, the loss function is the logistic loss: $l(w, x_i, z_i) = \log(1 + e^{-z_i(w^T x_i + b)})$

This framework connects logistic regression to other classification algorithms:

- SVM use the hinge loss: $l(w, x_i, z_i) = \max(0, 1 - z_i(w^T x_i + b))$
- Linear Regression can be viewed as using squared loss: $l(w, x_i, y_i) = (y_i - w^T x_i - b)^2$

Why "Empirical" Risk?

The term "empirical" emphasizes that we're calculating the risk based only on our observed training samples, not the true underlying data distribution. In statistical learning theory:

The true risk would be the expected loss over the entire data distribution

The empirical risk is an approximation based on our finite training set

Minimizing the empirical risk doesn't guarantee minimizing the true risk (which would give the best generalization), but with enough data and appropriate regularization, the two often align well.

The infinite Solution Problem

Linearly separable data means we can draw a straight line (or hyperplane in higher dimension) that perfectly separates the two classes. Mathematically, there exists some combinations of weights w and bias b such that:

- For all positive examples ($z_i = 1$): $w^T x_i + b \geq 0$
- For all negative examples ($z_i = -1$): $w^T x_i + b < 0$

This means all training points lie on the correct side of the decision boundary, making $z_i s_i \geq 0$ for all training points.

Here is the key insight: if we find such a separating hyperplane defined by w and b , we can create an infinite number of equivalent solutions by simply scaling w and b . For example, if (w, b) gives perfect separation then $(2w, 2b)$ also give perfect separation.

with even greater margins.

As we increase $\|w\|_2$ (the norm of w), the value of $s_i = w^T x_i + b$ grows larger for all points, pushing them further from the decision boundary. Looking at our logistic loss function $l(s_i, s_i) = \log(1 + e^{-s_i})$ we can see that:

- As s_i becomes increasingly positive, the loss approaches zero.
- The total loss function $J(w)$ approaches zero as $\|w\|_2 \rightarrow \infty$.



This is a problem!! our optimization has no minimum, only a greatest lower bound (unlike previous missing) that goes to zero when $w \rightarrow \infty$

So the algorithm would keep increasing $\|w\|_2$ indefinitely, leading to numerical issues.

The Regularization Solution

To solve this problem, we add a regularization term to our objective function that penalize large values of $\|w\|_2$. The most common term is L2 regularization (also called Ridge regularization):
The objective function that we want to minimize becomes:

$$R(w, b) = \frac{1}{2} \|w\|^2 + \sum_{i=1}^n \log(1 + e^{-s_i(w^T x_i + b)})$$

Or in the averaged form:

$$R(w, b) = \frac{1}{2} \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \log(1 + e^{-s_i(w^T x_i + b)})$$

Where:

- λ is the regularization parameter (hyperparameter)
- $\|w\|^2$ is the squared L2 norm of w (sum of squared weights)
- The second term is the logistic loss function

Regularization creates a trade-off between two competing objectives:

- Minimize the classification error (logistic loss term) ①
- Keep the model simple (regularization term) ②

As $\|w\|_2$ grows, the first term decreases (better classification), but the second term increases (more complex model). The optimum occurs at a finite value of $\|w\|_2$, ensuring a well-defined solution.

The regularization parameters control the balance between these two objectives:

- Large λ : Strongly favors simpler models (smaller $\|w\|_2$) potentially at the cost of classification accuracy
- Small λ : Focuses more on minimizing classification error, potentially leading to overfitting
- $\lambda = 0$: Revert to the origin, with the infinite problem.

Importantly, we cannot determine the optimal λ by minimizing the same objective function, that would give always $\lambda = 0$, defeating the purpose of the validation. Instead we typically use cross-validation to find the value of λ that gives the best performance on held-out validation set.

Understanding Contour Lines in the Logistic Regression Visualization

Those lines in the plots are contour lines, which are a powerful way to visualize the behavior of the logistic regression model. Let me explain exactly what they represent and why they're so revealing about regularization's effects.

What Are Contour Lines?

In these plots, the contour lines represent locations where the function $w^T x + b$ has the same value. Each line corresponds to a specific value of this function, which you can see labeled on the lines (+15, +10, +5, etc. in the left plot; +3, +2, +1, etc. in the right plot).

Think of these elevation lines on a topographical map. Just as elevation lines on a hiking map show points of equal height, these contour lines show points where the model's "score" is equal.

The Physical Landscape Analogy

Imagine the decision function $w^T x + b$ creates a 3D landscape above the 2D feature space:

The decision boundary (purple line) is where this landscape has height zero.

Areas above zero (positive values) belong to one class.

Areas below zero (negative values) belong to the other class.

The contour lines are like the elevation lines on this landscape:

Red contour lines show positive values (hills)

Blue contour lines show negative values (valleys)

The purple line is the zero contour (sea level)

Reading the Spacing Between Lines

The spacing between contour lines tells us about the steepness of this landscape:

Tightly packed lines (left plot, $\lambda = 0$) indicate a steep slope

Widely spaced lines (right plot, $\lambda = 1$) indicate a gentle slope

This spacing directly relates to w :

A large w creates a steep landscape with tightly packed contour lines

A small w creates a gentle landscape with widely spaced contour lines

Why the Values Are Different Between Plots

Notice how the values on the contour lines differ dramatically between the two plots:

Without regularization ($\lambda = 0$):

The contour lines include extreme values: +15, +10, +5, -5, -10

This indicates large coefficient values (large w)

The model is making very confident predictions far from the boundary

With regularization ($\lambda = 1$):

The contour values are more moderate: +3, +2, +1, -1, -2, -3

This indicates smaller coefficient values (smaller w)

The model makes more conservative predictions

The Decision-Making Process

These contour lines also reveal how the model makes decisions:

For a new data point, the model calculates $w^T x + b$

If this value is positive (red region), it predicts one class

If this value is negative (blue region), it predicts the other class

The magnitude tells us about the model's confidence

In the unregularized model, a point just slightly away from the boundary might already have a value of +5 or -5, indicating high confidence. In the regularized model, you'd need to go further from the boundary to reach the same confidence level.

Practical Significance

This visualization demonstrates why regularization improves generalization:

The unregularized model ($\lambda = 0$) is overconfident. It rapidly switches from strongly predicting one class to strongly predicting the other with only small movements in feature space. This makes it sensitive to small perturbations in the data. **Oversensitive problem**

The regularized model ($\lambda = 1$) is more cautious. It transitions more gradually between classes, requiring larger movements in feature space to achieve high confidence. This creates a more robust model that better handles variations in data.

The gray shading in the right plot likely represents a region of uncertainty where the model isn't extremely confident in either direction, which is a hallmark of a well-regularized model that knows when to be unsure.

Understanding these contour lines gives you deep insight into how regularization shapes the decision landscape and why this leads to models that generalize better to new data.

Some Consideration

The Regularization Mechanism

When we add L2 regularization to logistic regression, we're adding a penalty term λw^2 to our objective function. This penalizes large weights, encouraging the model to keep all weights small. But here's the crucial point: this penalty treats all weights equally, regardless of which feature they correspond to. The regularization doesn't know or care about the meaning or scale of your features—it simply pushes all weights toward zero with the same force.

The Problem With Different Feature Scales

Imagine we're predicting house prices using two features:

Feature 1: House size (in square feet), ranging from 1,000 to 5,000

Feature 2: Number of bedrooms, ranging from 1 to 6

Without regularization, the model might learn weights like:

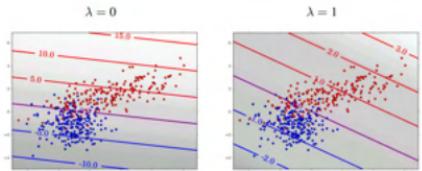
$w_1 = 0.01$ (for house size)

$w_2 = 10,000$ (for bedrooms)

These weights could give similar importance to both features:

A 1,000 sq ft increase contributes $0.01 \times 1,000 = 10$ to the prediction

An extra bedroom contributes $10,000 \times 1 = 10,000$ to the prediction



SO SO OVERCONFIDENT
ANGULAR SF IL CONSISTENT
A LOT SO CONFIDENT
WEAK IN MM CLASSIFICATION
ZONE SIN HUSTLE

So it's useful to preprocess data to improve the performance of logistic regression.

① When we center the data, we subtract the mean from each feature:

$$x'_i = x_i - \mu$$

This shifts all feature distributions to have a mean of zero. The mean can be calculated in different ways:

- Using the entire training set (most common)
- Using a weighted mean, such as the average of class means (useful when class are imbalanced)

(Centering helps remove arbitrary offsets in the data. For example, if one feature is temperature in Celsius (perhaps ranging from 20°C to 30°C), centering would transform it to range from roughly -5 to +5, making 0 a meaningful reference point.)

② Standardizing Variances

Standardizing goes beyond centering by also scaling each feature to have unit variance:

$$x'_{i(j)} = \frac{x_{i(j)}}{\sigma_{(j)}}$$

Where:

$x_{i(j)}$ is the j -th feature of the i -th sample, and $\sigma_{(j)}$ is the standard deviation of the j -th feature across the training set.

This ensures all features contribute equally to the distance calculations and regularization penalties.

For instance, if we have income (in thousands, perhaps ranging from 30-200) and age (ranging from 20-60), standardization would put both on comparable scales.

③ and ② → can be used together = zero normalization

③ Whitening the covariance matrix (so invariance)

Whitening standardizes variances and decorrelates features

$$x'_i = Ax_i \text{ where } A = Z^{-1} \quad \text{the } x'_i \text{ data will have the } I \text{ as covariance matrix}$$

Z is the covariance matrix of the training data, and Z^{-1} is the matrix square root inverse.

This transformation creates features that are:

- Uncorrelated with each other (orthogonal)
- All have unit variance

This technique can be particularly helpful when features have strong correlations, which is common in image data, audio signals, or financial time series. A variation uses the within-class covariance instead of the total covariance, which can better preserve class-discriminative information.

④ L2 Normalization

L2 Normalization scales each sample to have unit length:

$$x'_i = x_i / \|x_i\|$$

This makes all samples lie on the unit sphere. It's always performed after centering and whitening.



L₂ normalization is especially useful when the overall magnitude of each features doesn't matter, but rather prioritizes the relative proportions between features.

Will we choose one method or another using cross-validation.

logistic regression optimizes the probability of observing the training labels. As a result it naturally incorporates the empirical class distribution (prior) from the training data.

This mean if your training set has 90% negative samples and 10% positive samples the model implicitly learn this 9:1 ratio and adjust its predictions accordingly.

However this became problematic when the real world application has a different class distributions than the training data. For example my model can be deployed in a setting the true probability of positive class is 50% not 10%.

A solution to this can be adjusting the model score to behave like a log-likelihood ratio by subtracting the empirical prior log-odds (of the training dataset)

$$S_{LR} = w^T x + b - \log \frac{n_p}{n_n} = w^T x + b - \log \frac{P_{Temp}}{1 - P_{Temp}}$$

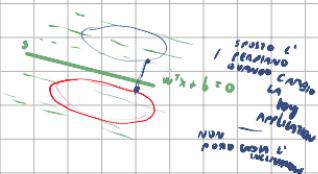
Where:

- n_p is the number of positive samples in training
- n_n " negative samples "
- $\log(n_p/n_n)$ is the log-odds of the empirical prior in the training set.

When we want to use this model for a different application with prior P_T we can make decision by comparing:

$$S_{LR} + \log \left(\frac{P_T}{1 - P_T} \right)$$

Anyways this is not perfect in a nutshell we are only moving the decision boundary.



FOR GAUSSIAN NOW È UN PROBLEMA NON SONO LINEARI ALTA LINEARITÀ



$$w^T x + b - \log \frac{P_{Temp}}{1 - P_{Temp}} \leq 0 \quad \text{we can't incline the line}$$

Prior-Weighted Logistic Regression

An even better approach is to directly incorporate the target application's class distribution during training.

This is called prior-weighted logistic regression

$$RLW = \frac{1}{2} \|w\|^2 + \frac{\pi_1}{n_1} \sum_{i:z_i=1} l(z_i, s_i) + \frac{1-\pi_1}{n_0} \sum_{i:z_i=0} l(z_i, s_i)$$

this effectively reweights the loss distribution from positive and negative examples to match the target application's class distribution. Positive samples are weighted by π_1 and negative examples by $(1-\pi_1)$

Even with this approach we can still compute an adjusted log-likelihood ratio score for use with other applications if needed

Note: π_1 can be estimated

We test the model on MNIST digit pairs (e.g. 0 vs 1, 0 vs 2, ...)

MNIST — Average Pairwise EER for Logistic Regression

DimRed	$\lambda = 0$	$\lambda = 0.0001$	$\lambda = 0.001$	$\lambda = 0.1$	Tied Gau
RAW [768]	1.7%	1.4%	1.2%	2.0%	—
PCA [50]	1.4%	1.4%	1.4%	2.1%	1.7%
PCA [100]	1.3%	1.2%	1.2%	2.0%	1.5%

LogReg obtains better performance than the Gaussian model

Regularization is important, especially when we do not reduce the dimensionality (we have more parameters to estimate, so over-fitting is more severe)

If we regularize too much the model performs poorly again

Multiclass Logistic Regression

[Take inspiration]

We begin with insights from the linear Gaussian classifier with uniform priors. For this model, the log posterior probability for any class j takes the form:

$$\log P(C=j|x) = w_j^T x + b_j + h(x) \quad \begin{cases} \text{for extended formula} \\ \text{see gaussian model} \end{cases}$$

Where

w_j is the weight vector for class j

b_j is the bias term for class j

$h(x)$ is a term that depends only on the input x , not on the class

$h(x)$ is the same for all classes, meaning it cancels out when comparing classes. This gives us the proportionality relationship.

$$P(C=j|x) \propto e^{w_j^T x + b_j}$$

Since we are working with a closed-set classification problem (the sample must belong to one of our classes) the probabilities must sum to 1:

$$\sum_{j=0}^n P(C=j|x) = 1$$

The normalized factor for $P(C=j|x)$ is thus: $\frac{(\text{answer})}{\sum_{j=0}^n e^{w_j^T x + b_j}}$

so the posterior probability corresponds to:

$$P(C=k|x) = \frac{e^{w_k^T x + b_k}}{\sum_j e^{w_j^T x + b_j}} \stackrel{(\text{answer})}{=} \frac{e^{w_k^T x + b_k}}{\sum_j e^{w_j^T x + b_j}}$$

Note: This ensures that the sum is $\frac{e^{w_0^T x + b_0}}{\sum_j e^{w_j^T x + b_j}} + \dots + \frac{e^{w_n^T x + b_n}}{\sum_j e^{w_j^T x + b_j}} = 1$.

The function that performs this normalization transforming raw scores into proper probabilities is called softmax.

$$f(s) = \left[\frac{e^{s_0}}{\sum_j e^{s_j}}, \dots, \frac{e^{s_n}}{\sum_j e^{s_j}} \right]$$

The softmax is a natural extension of the sigmoid function we used in binary logistic regression. In fact for k=2 softmax reduces to sigmoid function when we constrain $w_0=w_1, b_0=-b_1$ (we have two classes).

Given the model parameters

$$W = [w_1 \dots w_K], \quad b = \begin{bmatrix} b_1 \\ \vdots \\ b_K \end{bmatrix}$$

the logistic regression model allows computing the probability of each class

$$P(C = k|W, b, x) = \frac{e^{w_k^T x + b_k}}{\sum_j e^{w_j^T x + b_j}}$$

If we consider sample x_i , its class posterior distribution is thus a categorical distribution

VC: c_i van voorstelte
where $C_i|W, b, X_i = x_i \sim \text{Cat}(y_i) \quad x_i \rightarrow y_i = w_i^T x_i + b_i$

$$y_{ik} = \frac{e^{w_k^T x_i + b_k}}{\sum_j e^{w_j^T x_i + b_j}}$$

41/56

Note:

A categorical distribution is a natural extension of the Bernoulli distribution to multiple outcome. Instead of just two possible outcome, we now have K possibilities, and we need a specific probabilities for each one.

A categorical distribution is specified by a vector of K probabilities:

$$y = [y_0, y_1, \dots, y_K]$$

Where:

- each y_n represent the probability of outcomes n
- All probabilities are non-negative: $y_n \geq 0$ for all n
- The $\sum_n y_n = 1$

When we say C follow the categorical distribution with parameters y , written as $C \sim \text{Cat}(y)$ we mean $P(C=n) = y_n$ for $n=0, 1, \dots, k$

In logistic regression we compute this probability using the softmax function applied to our linear scores:

$$y_{in} = \frac{e^{w_i^T x_i + b_i}}{\sum_j e^{w_j^T x_i + b_j}}$$

This gives us a vector of probabilities $y_i = [y_{i0}, y_{i1}, \dots, y_{ik}]$ for each sample x_i . When we write:

$$C_i | W, b, X_i = x_i \sim \text{Cat}(y_i)$$

We're saying: "Given the input features x_i and our model parameters W and b , the class label C_i follows a categorical distribution with probabilities given by the vector y_i ".

As for the binary case we can express the log-probability of the training class labels as:

$$\ell(W, b) = \sum_{i=1}^n \log P(C_i=c_i | X_i=x_i, W, b)$$

To better represent this we can use the 1-of-N encoding.

1-of-N encoding is used to represent the true class.

For example, with a sample with true class c_i (which is a single number between 0 and k) we create a binary vector z_i where:

- $z_{ik}=1$ if $c_i=k$ (k is the true class.)
- $z_{ik}=0$ otherwise (for all other class.)

$$z_i = [0 \dots 0, 1, 0 \dots 0], \quad z_{ik} = \begin{cases} 1 & \text{if } c_i = k \\ 0 & \text{otherwise} \end{cases}$$

So the log-probability become:

$$\log P(C_i=c_i | X_i=x_i, W, b) = \log P(C_i=c_i | y_i) = \sum_{n=0}^k z_{in} \log y_{in}$$

(for a sample)

This sum picks out exactly the log-probability of the true class.

The log-likelihood of the entire training set becomes:

$$\ell(W, b) = \sum_{j=1}^J \sum_{n=0}^k z_{jn} \log y_{jn}$$

Where

- n is the number of the training samples
- K is the number of classes
- z_{ik} indicates whether sample i belongs to class k
- y_{ik} is the model predicted probability that sample i belong to class k

As for the binary case, the expression

$$H(z_i, y_i) = - \sum_{k=1}^K z_{ik} \log y_{ik}$$

SC STATION
CME IN CHIE
C 2
 $P(C=0/x) = 0$
 $P(C=1/x) = 1$

represents the (multiclass) cross-entropy between the observed and predicted label distributions for sample x_i

As for the binary case, we estimate W and b as to maximize the likelihood for the training labels

The ML solution is again the solution that minimizes the (average) cross-entropy:

$$\arg \max_{W,b} \ell(W, b) = \arg \max_{W,b} \left[- \sum_{i=1}^n H(z_i, y_i) \right] = \arg \min_{W,b} \sum_{i=1}^n H(z_i, y_i)$$

$H(z_i, y_i)$ measure the dissimilarity between:

- The true distributions π_i (all the mass to the true class)
- The predicted distributions y_i (softmax probabilities from our model)

When we model multiple classes, an interesting mathematical property emerges: our model has more parameters than it actually needs this is called over-parametrization

In multiclass logistic regression, we have a weight vector w_k , and bias term b_k , for each class k giving us K weight vectors and K bias terms. However not all of these parameters are independent

Why? Because only the difference of the class score matters, not their absolute values if we add the same constant vector to all the weight vectors w_k , the softmax probabilities remain exactly the same

$$P(C=k/x) = \frac{e^{w_k^T x + b_k}}{\sum_j e^{w_j^T x + b_j}}$$

if we can add c
we are using more
parameters

if we add c (constant) to every w_k :

$$P(C=k/x) = \frac{e^{(w_k+c)^T x + b_k}}{\sum_j e^{(w_j+c)^T x + b_j}} = \frac{e^{w_k^T x + b_k} e^{c^T x}}{e^{c^T x} \sum_j e^{w_j^T x + b_j}}$$

For a binary classification problem ($K=2$) multiclass logistic regressions should reduce to standard binary case
 if we subtract w_2 from both w_0 and w_1 , we recover the binary lr formula.

$$P(C=1|x) = \frac{e^{w_0^T x + b_0}}{e^{w_0^T x + b_0} + e^{w_1^T x + b_1}}$$

We redefine the parameters $w = w_0 - w_1$ $b = b_0 - b_1$

$$\begin{aligned} P(C=1|x) &= \frac{(w+w_0)^T x + b + b_2}{e^{(w+w_0)^T x + b + b_2} + e^{w_1^T x + b_1}} = \frac{\cancel{e^{w_0^T x + b_2}}}{\cancel{e^{w_0^T x + b_2}} \frac{e^{w_1^T x + b_1}}{e^{w_1^T x + b_1} + e^{w_0^T x + b_2}}} = \\ &= \frac{e^{w_1^T x + b_1}}{1 + e^{w_1^T x + b_1}} \end{aligned}$$

Finally, as for the binary class, we can cast the problem as a minimization of a loss function

We rewrite the objective in terms of class labels c as

$$\begin{aligned} J(W, b) &= - \sum_{i=1}^n \sum_{k=1}^K z_{ik} \log y_{ik} \\ &= - \sum_{i=1}^n \log \frac{e^{w_i^T x_i + b_{ci}}}{\sum_{c'=1}^K e^{w_{c'}^T x_i + b_{c'}}} \\ &= - \sum_{i=1}^n \left[\log \left(\sum_{c'=1}^K e^{w_{c'}^T x_i + b_{c'}} \right) - w_{ci}^T x_i - b_{ci} \right] \\ &= - \sum_{i=1}^n l(x_i, c_i, W, b) \end{aligned}$$

l is also called softmax loss

44/56

Where c_i is the true label of sample i , and c_i subscript in w_{ci} and b_{ci} , refers to the weights and the bias of the true class.

The interpretation of the softmax loss is:

- ① The first term $\log \left(\sum_{c'=1}^K e^{w_{c'}^T x_i + b_{c'}} \right)$ is the soft maximum of all the class scores, representing the normalized factor
- ② The second term $-(w_{ci}^T x_i + b_{ci})$ is the negative of the score for the true label
- ③ Together they encourage the model to increase the score of the correct class while decreasing the scores of incorrect classes

This loss function directly connects to the cross-entropy. It's the score entropy between a one-hot distribution (the true class) and the softmax probabilities (the model predictions)

As for the binary classifications we add regularization to avoid overfitting

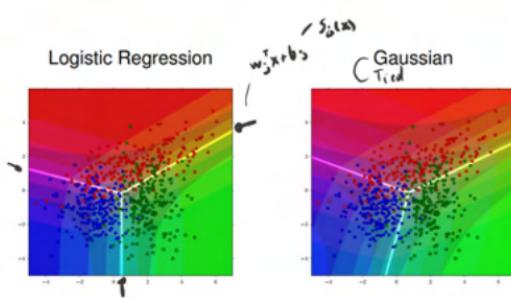
$$R(w, b) = \Omega(w) + \frac{1}{n} J(w, b)$$

A common regularization term is:

$$\Omega(w_1, \dots, w_n) = \frac{1}{2} \sum_i \|w_i\|^2$$

This L2 regularization penalizes large weights, encouraging simpler decision boundaries.

Again, we can replace the average cross-entropy with prior-weighted average cross entropy to account for priors that are different from the empirical training set prior



- ① Logistic regression create linear decision boundaries without making assumptions about feature distributions
- ② Gaussian models create quadratic boundaries (for equal covariance) base on distributional assumption

MNIST — Error rates for Logistic Regression

DimRed	$\lambda = 0$	$\lambda = 0.00001$	$\lambda = 0.001$	$\lambda = 0.1$	Tied Gau
RAW [768]	8.0%	7.4%	7.9%	12.9%	—
PCA [50]	8.8%	8.8%	8.9%	13.3%	12.6 %
PCA [100]	7.8%	7.8%	8.2%	12.9%	12.3%
PCA+LDA [9]	10.9%	10.9%	11.0%	12.4%	12.3 %

The multiclass logistic regression performs better than the Gaussian model — indeed, the Gaussian assumption is not very accurate for the features we are considering. LogReg assumes linear separation, but does not assume a specific distribution for the features

Again, regularization is important, especially when the feature space is large

On the MNIST dataset, linear logistic regression outperforms the tied covariance Gaussian Classifier but fall short compared to the non-tied covariance Gaussian Classifier (TWS)

Why does this happen? The answer lies in the decision boundaries

- Linear Logistic regression can only create linear separation boundaries (straight lines in 2D, flat planes in 3D)
- Tied covariance Gaussian models also create linear boundaries
- MyL create quadratic boundaries (curved in 2D)

For complex data, not linear boundary can be better.

We can transform our features to enable logistic regression to create non-linear decision boundaries, this is how.

For MyL a non linear Gaussian classifier:

$$\ln \frac{P(C=1|x)}{P(C=0|x)} = x^T A x + b^T x + c = s(x, A, b, c) \quad [\text{Quadratic form}]$$

- $x^T A x$ is a quadratic term (like x_1^2 or $x_1 x_2$)
- $b^T x$ is a linear term
- c is a constant

This creates curved decision boundaries that can separate more complex patterns.

We use a trick, while the expression is quadratic in the features x , it's actually linear in the parameter A and b .

Take a 2D example, the quadratic form $x^T A x$ expand to:

$$\begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = a_{11} x_1^2 + 2a_{12} x_1 x_2 + a_{21} x_2 x_1 + a_{22} x_2^2 + (a_{11} + a_{22})x_1 x_2 + a_{12} x_2^2$$

the coefficient a_0, a_1, \dots are linear [feature $\dots) = f(a_0) + f(a_1)x_1 + \dots]$

So the function is linear in the parameters but "quadratic" in the features"

We can so rewrite the quadratic term using a matrix inner product:

$$s(x, A, b, c) = \langle x x^T, A \rangle + b^T x + c$$

The inner product multiplies corresponding element of the matrices $x x^T$ and A , and sums them all

$$\langle A, B \rangle = \sum_{i,j} A_{ij} B_{ij}$$

We can further express $\langle A, x x^T \rangle$ as

$$\langle A, x x^T \rangle \gg \text{vec}(x x^T)^T \text{vec}(A) - [\text{dot product}]$$

prologue name

`vec(n)` stack the columns of A into a single column vector.

Finally we can combine all our term by defining:

$$\Phi(x) = \begin{bmatrix} \text{vec}(x x^T) \\ x \end{bmatrix}$$

$$w = \begin{bmatrix} \text{vec}(A) \\ 1 \end{bmatrix}$$

the log posterior ratio can be expressed as

$$s(x, w, c) = w^T \Phi(x) + c$$

Now instead of feeding our original feature x to logistic regression we can:

- ① Transform our features using the mapping $\Phi(x)$ that include all the quadratic product
- ② Train the standard logistic regression model on this expanded features.
- ③ The resulting model will have linear boundaries in the expanded space but quadratic boundaries in the original space

For example, if our original feature were $[x_0, x_1]$ our expanded features would be $[x_0^2, x_0 x_1, x_1^2, x_0, x_1]$ capturing all quadratic terms.

The LR model (both binary and multiclass) allows computing linear separation rules for the transformed features $\phi(x)$

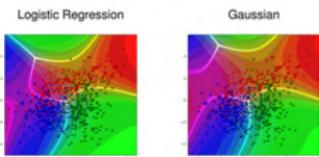
Since expressions $w^T \phi(x) + c$ correspond to quadratic forms in the original feature space, we are actually estimating **quadratic separation surfaces** in the original space

In general, we can consider a transformation $\phi(x)$ of our feature space such that our classes are (approximately) linearly separable in the expanded feature space

This approach is powerful but it comes at a computational costs

For original features of dimension M , a polynomial expansion of degree d result in an expanded feature space of dimension $O(d^M)$. In this case we are using quadratic form so $d=2$ ($\text{len}(x) \times \text{len}(x)$)

LA Dimensionality Reduction
Multidim



b1221
MNIST — Average pairwise EER for LR with quadratic feature expansion

DimRed	$\lambda = 0$	$\lambda = 1e^{-5}$	$\lambda = 1e^{-3}$	$\lambda = 1e^{-1}$	Gaussian
PCA [50]	1.0%	1.0%	0.9%	1.5%	0.8%

MNIST — Multiclass error rates for LR with quadratic feature expansion

DimRed	$\lambda = 0$	$\lambda = 1e^{-5}$	$\lambda = 1e^{-3}$	$\lambda = 1e^{-1}$	Gaussian
PCA [50]	2.3%	1.9%	1.7%	3.1%	3.6%

==