

Logistic regression

CS 446 / ECE 449

2022-01-25 09:57:46 -0600 (4c2f15b)

Plan for today

- ▶ Linear **classifiers**.
- ▶ ERM for classification.
- ▶ Solving the ERM problem.

Linear regression (last lecture) vs logistic regression (this lecture)

1. Start from **training data** $((\mathbf{x}_i, y_i))_{i=1}^n$, with $\mathbf{x}_i \in \mathbb{R}^d$ and $y_i \in \mathbb{R}$.

Last lecture: $y_i \in \mathbb{R}$; **this lecture:** $y_i \in \{+1, -1\}$.

2. **Model** is a **linear predictor**: pick $\mathbf{w} \in \mathbb{R}^d$ with

$$\mathbf{x}_i \mapsto \mathbf{w}^\top \mathbf{x}_i =: \hat{y}_i \approx y_i.$$

3. Choose \mathbf{w} by minimizing **empirical risk** (average loss over training set):

$$\hat{\mathcal{R}}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}^\top \mathbf{x}_i, y_i).$$

Last lecture: squared loss ℓ_{sq} ; **this lecture:** logistic loss ℓ_{logistic} .

4. Basic method: **gradient descent**. Set $\mathbf{w}_0 = 0$, and thereafter

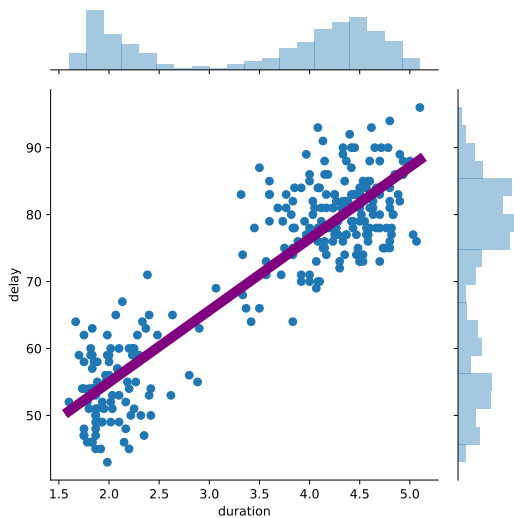
$$\mathbf{w}_{i+1} := \mathbf{w}_i - \eta \nabla \hat{\mathcal{R}}(\mathbf{w}_i),$$

where η is a **learning rate** (step size).

Same in both lectures, however **least squares also had SVD solution**.

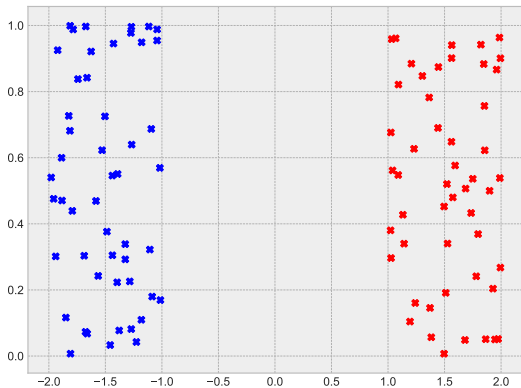
Linear regression

Last lecture, we studied [regression](#);
the output/label space was \mathbb{R} .



Linear classification

Today, the goal is a **classification**;
the output/label space is discrete.

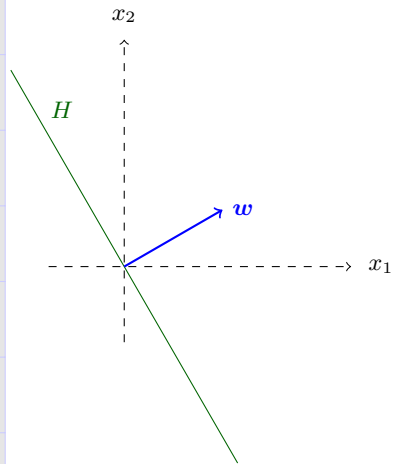


Binary classification means output space $\mathcal{Y} = \{-1, +1\}$.

A linear predictor $\mathbf{w} \in \mathbb{R}^d$ classifies according to $\text{sign}(\mathbf{w}^\top \mathbf{x}) \in \{-1, +1\}$.

Given $((\mathbf{x}_i, y_i))_{i=1}^n$ and a weight vector $\mathbf{w} \in \mathbb{R}^d$,
we want $\hat{y}_i := \text{sign}(\mathbf{w}^\top \mathbf{x}_i) \in \{-1, +1\}$ and y_i to agree.

Geometry of linear classifiers



Given $w \in \mathbb{R}^d$, predict with

$$\mathbf{x} \mapsto \text{sign}(\mathbf{w}^\top \mathbf{x}) \in \{\pm 1\}.$$

Let H be the hyperplane orthogonal to w :

$$H := \left\{ \mathbf{x} \in \mathbb{R}^d : \mathbf{x}^\top \mathbf{w} = 0 \right\}.$$

H splits \mathbb{R}^d into points we label positive,

$$\left\{ \mathbf{x} \in \mathbb{R}^d : \mathbf{w}^\top \mathbf{x} > 0 \right\},$$

and points we label negative,

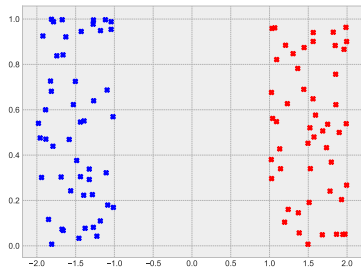
$$\left\{ \mathbf{x} \in \mathbb{R}^d : \mathbf{w}^\top \mathbf{x} < 0 \right\}.$$

The **decision boundary** is H .

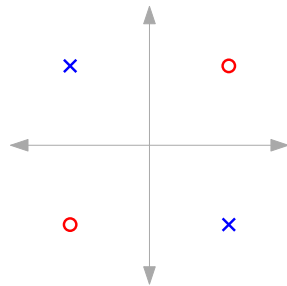
Linear separability

Is it always possible to find w with $\text{sign}(w^T x_i) = y_i$?

I.e., is it always possible to find a (homogeneous) hyperplane which **separates** the data?



Linearly separable.

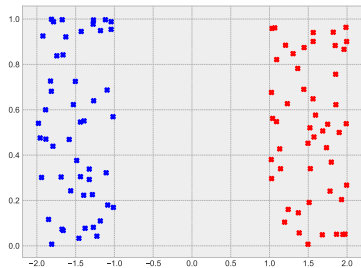


Not linearly separable.

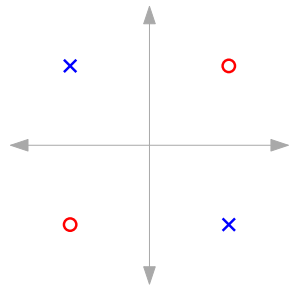
Linear separability

Is it always possible to find w with $\text{sign}(w^T x_i) = y_i$?

I.e., is it always possible to find a (homogeneous) hyperplane which **separates** the data?



Linearly separable.



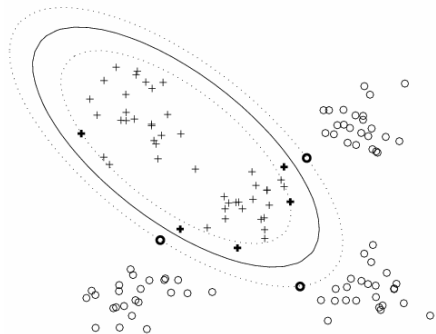
Not linearly separable.

Lecture 4: adding **features** can make things separable.

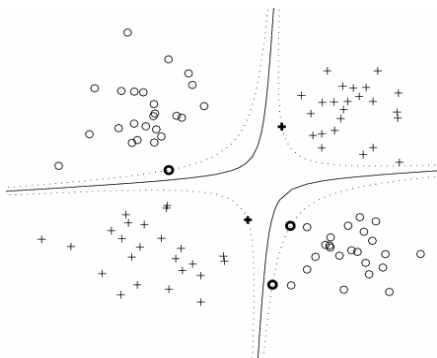
Decision boundary with quadratic feature expansion

With old faithful data, we appended 1 to x .

In general, we can map x to more exotic things, and separate more data.



elliptical decision boundary



hyperbolic decision boundary

We'll discuss this next lecture.

Finding linear classifiers with ERM

Finding linear classifiers with ERM

Why not try “pytorch meta-algorithm”, with empirical risk

$$\arg \min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \mathbb{1}[\text{sign}(\mathbf{w}^\top x_i) \neq y_i] ?$$

- ▶ Discrete/combinatorial search;
NP-hard in general; awkward for continuous optimization algorithms.

Relaxing the ERM problem

- ▶ **Step 1:** remove $\text{sign}(\cdot)$. Let's remove one source of discreteness:

$$\frac{1}{n} \sum_{i=1}^n \mathbb{1}[\text{sign}(\mathbf{w}^\top \mathbf{x}_i) \neq y_i] \quad \longrightarrow \quad \frac{1}{n} \sum_{i=1}^n \mathbb{1}[y_i(\mathbf{w}^\top \mathbf{x}_i) \leq 0] .$$

Are these equivalent?

- ▶ **Step 2:** remove $\mathbb{1}[\cdot]$. Rewrite the preceding as

$$\widehat{\mathcal{R}}_{z_0}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell_{z_0}(y_i \mathbf{w}^\top \mathbf{x}_i) \quad \text{where } \ell_{z_0}(z) = \mathbb{1}[z \leq 0].$$

Here, $y_i(\mathbf{w}^\top \mathbf{x}_i)$ is the (unnormalized) **margin** of \mathbf{w} on example i .

Next let's replace ℓ_{z_0} with something continuous!

(**Side note:** squared loss had two arguments, today we'll use one. We'll discuss this point next lecture.)

Logistic loss

We want to choose a nice loss ℓ in

$$\widehat{\mathcal{R}}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell(y_i \mathbf{w}^\top \mathbf{x}_i).$$

Key desired properties:

- ▶ ℓ is **continuous** (for sake of optimization);
- ▶ ℓ **prefers correct classifications**:
if $y_i \mathbf{w}^\top \mathbf{x}_i > 0$ (correct), then $\ell(y_i \mathbf{w}^\top \mathbf{x}_i) < \ell(-y_i \mathbf{w}^\top \mathbf{x}_i)$.

Logistic loss

We want to choose a nice loss ℓ in

$$\widehat{\mathcal{R}}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell(y_i \mathbf{w}^\top \mathbf{x}_i).$$

Key desired properties:

- ▶ ℓ is **continuous** (for sake of optimization);
- ▶ ℓ **prefers correct classifications**:
if $y_i \mathbf{w}^\top \mathbf{x}_i > 0$ (correct), then $\ell(y_i \mathbf{w}^\top \mathbf{x}_i) < \ell(-y_i \mathbf{w}^\top \mathbf{x}_i)$.

Examples.

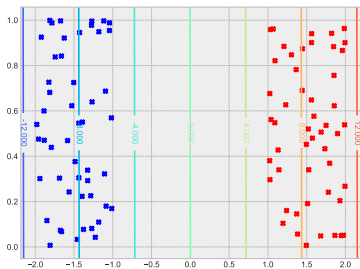
- ▶ **Squared loss**, written in margin form: $\ell_{\text{ls}}(z) := \frac{1}{2}(1 - z)^2$; note

$$2\ell_{\text{ls}}(y\hat{y}) = (1 - y\hat{y})^2 = y^2(1 - y\hat{y})^2 = (y - \hat{y})^2.$$

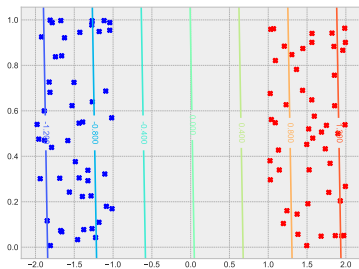
Squared loss doesn't just push $y_i \hat{y}_i$ positive; it wants $y_i \hat{y}_i = 1$!

- ▶ **Logistic loss**: $\ell_{\text{log}}(z) = \ln(1 + \exp(-z))$.
This one doesn't care so long as $z = y\hat{y} > 0$.

Squared and logistic losses on linearly separable data I

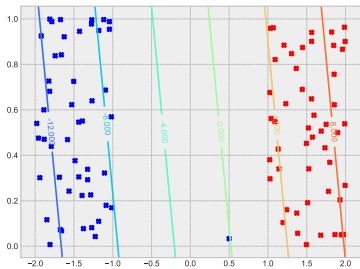


Logistic loss.

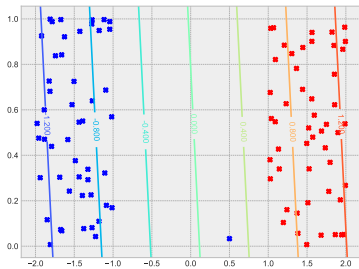


Squared loss.

Squared and logistic losses on linearly separable data II

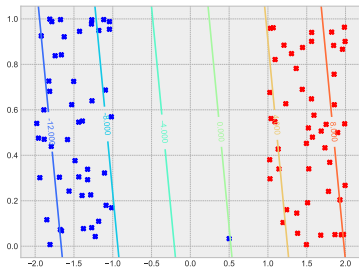


Logistic loss.

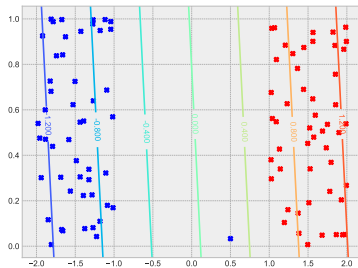


Squared loss.

Squared and logistic losses on linearly separable data II



Logistic loss.



Squared loss.

(Math note: it's easy to prove this.)

Least squares and logistic ERM

Least squares and logistic ERM

Least squares:

- ▶ Take gradient of $2n\widehat{\mathcal{R}}_{\text{ls}}(\mathbf{w}) = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$, set to 0; obtain **normal equations** $\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y}$. Suffice to take **OLS solution** $\hat{\mathbf{w}}_{\text{ols}} = \mathbf{X}^+ \mathbf{y}$.
- ▶ Alternatively, gradient descent.

Logistic loss:

- ▶ Unclear how to solve

$$\nabla_{\mathbf{w}} \widehat{\mathcal{R}}_{\text{log}}(\mathbf{w}) = \nabla_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \ln(1 + \exp(-y_i \mathbf{w}^T \mathbf{x}_i)) = 0.$$

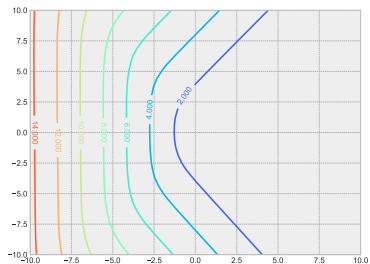
- ▶ Gradient descent still fine!

Given an empirical risk $\hat{\mathcal{R}} : \mathbb{R}^p \rightarrow \mathbb{R}$, **gradient descent** is the iteration

$$\mathbf{w}_{i+1} := \mathbf{w}_i - \eta_i \nabla_{\mathbf{w}} \hat{\mathcal{R}}(\mathbf{w}_i),$$

where \mathbf{w}_0 is given, and η_i is a **learning rate (step size)**.

Gradient descent goes down the contours of $\hat{\mathcal{R}}$:

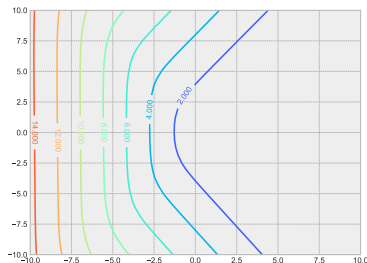


Given an empirical risk $\hat{\mathcal{R}} : \mathbb{R}^p \rightarrow \mathbb{R}$, **gradient descent** is the iteration

$$\mathbf{w}_{i+1} := \mathbf{w}_i - \eta_i \nabla_{\mathbf{w}} \hat{\mathcal{R}}(\mathbf{w}_i),$$

where \mathbf{w}_0 is given, and η_i is a **learning rate (step size)**.

Gradient descent goes down the contours of $\hat{\mathcal{R}}$:



Remarks.

- ▶ In the convexity lecture, we'll show this works for linear and logistic regression.
- ▶ For deep networks, \mathbf{w}_0 is random, and η_i is highly tuned/varying.

For logistic loss and linear predictors, typically $w_0 = 0$, and

For logistic loss and linear predictors, typically $\mathbf{w}_0 = \mathbf{0}$, and

$$\mathbf{w}_{i+1} := \mathbf{w}_i - \eta_i \nabla_{\mathbf{w}} \widehat{\mathcal{R}}_{\log}(\mathbf{w}_i) = \mathbf{w}_i - \frac{\eta_i}{n} \sum_{j=1}^n \ell'_{\log}(y_j \mathbf{x}_j^\top \mathbf{w}_i) y_j \mathbf{x}_j,$$

where $\ell'_{\log}(z) = \frac{-1}{1+\exp(z)}$.

For logistic loss and linear predictors, typically $\mathbf{w}_0 = 0$, and

$$\mathbf{w}_{i+1} := \mathbf{w}_i - \eta_i \nabla_{\mathbf{w}} \widehat{\mathcal{R}}_{\log}(\mathbf{w}_i) = \mathbf{w}_i - \frac{\eta_i}{n} \sum_{j=1}^n \ell'_{\log}(y_j \mathbf{x}_j^{\top} \mathbf{w}_i) y_j \mathbf{x}_j,$$

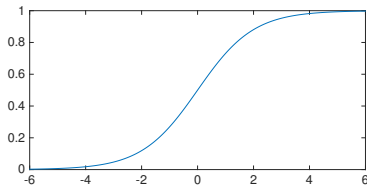
where $\ell'_{\log}(z) = \frac{-1}{1+\exp(z)}$.

But who cares, pytorch does it for us:

```
def GD(X, y, loss, step = 0.1, n_iters = 10000):  
    w = torch.zeros(X.shape[1], requires_grad = True)  
    for i in range(n_iters):  
        l = loss(X, y, w).mean()  
        l.backward()  
  
        with torch.no_grad():  
            w -= step * w.grad  
            w.grad.zero_()  
  
    return w
```

“Logistic” “regression”?

The (negative) derivative $-\ell'_{\log}(z) = \frac{1}{1+e^z}$ is the logistic function.



It can **suggest** a **probability / confidence** of the output label.
We'll revisit this next lecture.

Warning: many treat these explicitly as confidences, but they are not.

Summary for today

- ▶ Linear **classifiers**.
- ▶ ERM for classification.
- ▶ Solving the ERM problem.

(Appendix.)

Logistic risk and separation

If there exists a perfect linear separator, empirical logistic risk minimization should find it.

Theorem.

Logistic risk and separation

If there exists a perfect linear separator, empirical logistic risk minimization should find it.

Theorem. If there exists $\bar{\mathbf{w}}$ with $y_i \bar{\mathbf{w}}^\top \mathbf{x}_i > 0$ for all i , then every \mathbf{w} with $\widehat{\mathcal{R}}_{\log}(\mathbf{w}) < \ln(2)/2n + \inf_{\mathbf{v}} \widehat{\mathcal{R}}_{\log}(\mathbf{v})$, also satisfies $y_i \mathbf{w}^\top \mathbf{x}_i > 0$.

Logistic risk and separation

If there exists a perfect linear separator, empirical logistic risk minimization should find it.

Theorem. If there exists $\bar{\mathbf{w}}$ with $y_i \bar{\mathbf{w}}^\top \mathbf{x}_i > 0$ for all i , then every \mathbf{w} with $\hat{\mathcal{R}}_{\log}(\mathbf{w}) < \frac{\ln(2)}{2n} + \inf_{\mathbf{v}} \hat{\mathcal{R}}_{\log}(\mathbf{v})$, also satisfies $y_i \mathbf{w}^\top \mathbf{x}_i > 0$.

Proof.

Step 1: low risk implies few mistakes. For any \mathbf{w} with $y_j \mathbf{w}^\top \mathbf{x}_j \leq 0$ for some j ,

$$\hat{\mathcal{R}}_{\log}(\mathbf{w}) \geq \frac{1}{n} \ln(1 + \exp(-y_j \mathbf{w}^\top \mathbf{x}_j)) \geq \frac{\ln(2)}{n}.$$

By contrapositive, any \mathbf{w} with $\hat{\mathcal{R}}_{\log}(\mathbf{w}) < \frac{\ln(2)}{n}$ makes no mistakes.

Step 2: $\inf_{\mathbf{v}} \hat{\mathcal{R}}_{\log}(\mathbf{v}) = 0$. Note:

$$0 \leq \inf_{\mathbf{v}} \hat{\mathcal{R}}_{\log}(\mathbf{v}) \leq \inf_{r>0} \frac{1}{n} \sum_{i=1}^n \ln(1 + \exp(-r y_i \bar{\mathbf{w}}^\top \mathbf{x}_i)) = 0.$$

□

Remark. We didn't prove that gradient descent finds such a predictor, but that in turn is an easy consequence of the above and one bound from the upcoming convexity lecture.

The Perceptron algorithm (streamed data!)

Start from $\mathbf{w}_0 = 0$, and thereafter rotate towards $y_i \mathbf{x}_i$ if wrong:

$$\mathbf{w}_{i+1} := \mathbf{w}_i + y_{i+1} \mathbf{x}_{i+1} \mathbb{1}[y_{i+1} \mathbf{w}_i^\top \mathbf{x}_{i+1} \leq 0].$$

Remark: this is a specific subgradient of $\partial_{\mathbf{w}} \max\{0, -y_{i+1} \mathbf{x}_{i+1}^\top \mathbf{w}_i\}$.

Theorem (perceptron convergence, Novikoff '62). Suppose there exists \mathbf{u} with $\|\mathbf{u}\|_2 = 1$ and $\mathbf{u}^\top \mathbf{x}_i y_i \geq \gamma > 0$, and $\|\mathbf{x}_i\| \leq 1$. Then perceptron makes at most $1/\gamma^2$ mistakes:

$$\sum_{i < t} \mathbb{1}[\mathbf{w}_i^\top \mathbf{x}_{i+1} y_{i+1} \leq 0].$$

Proof. Define mistake set $\mathcal{M} := \{i < t : \mathbb{1}[\mathbf{w}_i^\top \mathbf{x}_{i+1} y_{i+1} \leq 0]\}$. Note

$$\|\mathbf{w}_t\| \geq \mathbf{w}_t^\top \mathbf{u} = \sum_{i \in \mathcal{M}} y_i \mathbf{x}_i^\top \mathbf{u} \geq \gamma |\mathcal{M}|.$$

On the other hand, by induction,

$$\begin{aligned} \|\mathbf{w}_t\|^2 &= \|\mathbf{w}_{t-1}\|^2 + y_t \mathbf{x}_t^\top \mathbf{w}_{t-1} \mathbb{1}[t-1 \in \mathcal{M}] + \|y_t \mathbf{x}_t \mathbb{1}[t-1 \in \mathcal{M}]\|^2 \\ &\leq \|\mathbf{w}_{t-1}\|^2 + \mathbb{1}[t-1 \in \mathcal{M}] \leq \dots \leq |\mathcal{M}|. \end{aligned}$$

Together, $\gamma |\mathcal{M}| \leq \|\mathbf{w}_t\| \leq \sqrt{|\mathcal{M}|}$, and $|\mathcal{M}| \leq 1/\gamma^2$. □