

Linear classification

COMS 4771 Fall 2019

Overview

- ▶ Logistic regression model
- ▶ Linear classifiers
- ▶ Gradient descent and SGD
- ▶ Multinomial logistic regression model

0 / 27

1 / 27

Logistic regression model I

- ▶ Suppose \mathbf{x} is given by d real-valued features, so $\mathbf{x} \in \mathbb{R}^d$, while $y \in \{-1, +1\}$.
- ▶ Logistic regression model for (\mathbf{X}, Y) :
 - ▶ $Y \mid \mathbf{X} = \mathbf{x}$ is Bernoulli (but taking values in $\{-1, +1\}$ rather than $\{0, 1\}$), with “heads probability” parameter

$$\frac{1}{1 + \exp(-\mathbf{x}^\top \mathbf{w})}.$$

- ▶ $\mathbf{w} \in \mathbb{R}^d$ is parameter vector of interest
- ▶ \mathbf{w} not involved in marginal distribution of \mathbf{X} (which we don’t care much about)

2 / 27

Logistic regression model II

- ▶ Sigmoid function $\sigma(t) := 1/(1 + e^{-t})$
 - ▶ Useful property: $1 - \sigma(t) = \sigma(-t)$
 - ▶ $\Pr(Y = +1 \mid \mathbf{X} = \mathbf{x}) = \sigma(\mathbf{x}^\top \mathbf{w})$
 - ▶ $\Pr(Y = -1 \mid \mathbf{X} = \mathbf{x}) = 1 - \sigma(\mathbf{x}^\top \mathbf{w}) = \sigma(-\mathbf{x}^\top \mathbf{w})$
- ▶ Convenient formula: for each $y \in \{-1, +1\}$,

$$\Pr(Y = y \mid \mathbf{X} = \mathbf{x}) = \sigma(y\mathbf{x}^\top \mathbf{w}).$$

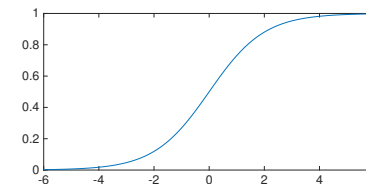


Figure 1: Sigmoid function

3 / 27

Log-odds in logistic regression model

- ▶ Log-odds in the model is given by a linear function:

$$\ln \frac{\Pr(Y = +1 \mid \mathbf{X} = \mathbf{x})}{\Pr(Y = -1 \mid \mathbf{X} = \mathbf{x})} = \mathbf{x}^\top \mathbf{w}.$$

- ▶ Just like in linear regression, common to use feature expansion!
 - ▶ E.g., affine feature expansion $\varphi(\mathbf{x}) = (1, \mathbf{x}) \in \mathbb{R}^{d+1}$

4 / 27

Optimal classifier in logistic regression model

- ▶ Recall that [Bayes classifier](#) is

$$f^*(x) = \begin{cases} +1 & \text{if } \Pr(Y = +1 \mid X = x) > 1/2 \\ -1 & \text{otherwise.} \end{cases}$$

- ▶ If distribution of (\mathbf{X}, Y) comes from logistic regression model with parameter \mathbf{w} , then Bayes classifier is

$$\begin{aligned} f^*(\mathbf{x}) &= \begin{cases} +1 & \text{if } \mathbf{x}^\top \mathbf{w} > 0 \\ -1 & \text{otherwise.} \end{cases} \\ &= \text{sign}(\mathbf{x}^\top \mathbf{w}). \end{aligned}$$

- ▶ This is a [linear classifier](#)
- ▶ Many other statistical models for classification data lead to a linear (or affine) classifier, e.g., Naive Bayes

5 / 27

Linear classifiers, operationally

- ▶ Compute linear combination of features, then check if above threshold (zero)

$$\text{sign}(\mathbf{x}^\top \mathbf{w}) = \begin{cases} +1 & \text{if } \sum_{i=1}^d w_i x_i > 0 \\ -1 & \text{otherwise.} \end{cases}$$

- ▶ With affine feature expansion, threshold can be non-zero

```
1: if  $0.335 \cdot x_1 + 2.5 \cdot x_2 + \dots + 6.35 \cdot x_{10^6} > 4.3$  then
2:   return spam
3: else
4:   return not spam
5: end if
```

Figure 2: Example of an affine classifier

6 / 27

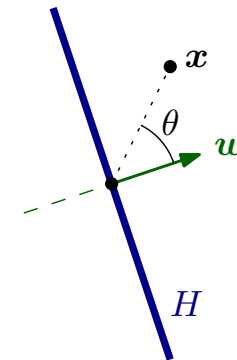
Geometry of linear classifiers I

- ▶ [Hyperplane](#) specified by [normal vector](#) $\mathbf{w} \in \mathbb{R}^d$:

- ▶ $H = \{\mathbf{x} \in \mathbb{R}^d : \mathbf{x}^\top \mathbf{w} = 0\}$
- ▶ This is the [decision boundary](#) of a linear classifier
- ▶ Angle θ between \mathbf{x} and \mathbf{w} has

$$\cos(\theta) = \frac{\mathbf{x}^\top \mathbf{w}}{\|\mathbf{x}\|_2 \|\mathbf{w}\|_2}$$

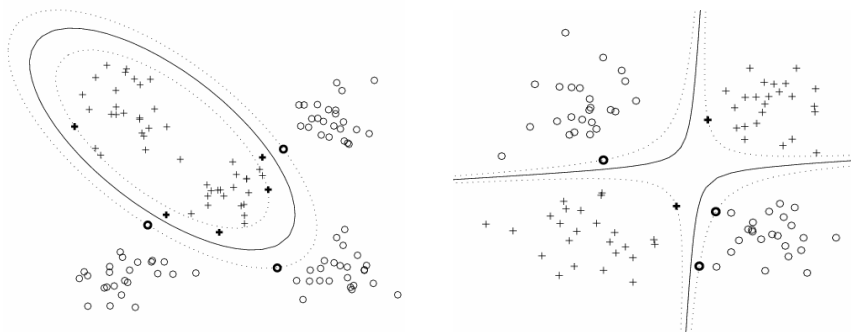
- ▶ Distance to hyperplane given by $\|\mathbf{x}\|_2 \cdot \cos(\theta)$
- ▶ \mathbf{x} is on same side of H as \mathbf{w} iff $\mathbf{x}^\top \mathbf{w} > 0$



7 / 27

Geometry of linear classifiers II

- ▶ With feature expansion, can obtain other types of decision boundaries



8 / 27

MLE for logistic regression

- ▶ Treat training examples as iid, same distribution as test example
- ▶ Log-likelihood of \mathbf{w} given data $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n) \in \mathbb{R}^d \times \{-1, +1\}$:

$$-\sum_{i=1}^n \ln(1 + \exp(-y_i \mathbf{x}_i^T \mathbf{w})) + \{\text{terms not involving } \mathbf{w}\}$$

- ▶ No “closed form” expression for maximizer
- ▶ (Later, we’ll discuss algorithms for finding approximate maximizers using iterative methods like gradient descent.)
- ▶ What about ERM perspective?

9 / 27

Zero-one loss and ERM for linear classifiers

- ▶ Recall: error rate of classifier f can also be written as risk:

$$\mathcal{R}(f) = \mathbb{E}[\mathbb{1}_{\{f(X) \neq Y\}}] = \Pr(f(X) \neq Y),$$

where loss function is zero-one loss.

- ▶ For classification, we are ultimately interested in classifiers with small error rate
 - ▶ I.e., small (zero-one loss) risk
- ▶ Just like for linear regression, can apply plug-in principle to derive empirical risk minimization (ERM), but now for linear classifiers.
 - ▶ Find $\mathbf{w} \in \mathbb{R}^d$ to minimize

$$\hat{\mathcal{R}}(\mathbf{w}) := \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{\text{sign}(\mathbf{x}_i^T \mathbf{w}) \neq y_i\}}.$$

- ▶ Very different from MLE for logistic regression

10 / 27

ERM for linear classifiers

- ▶ **Theorem:** In IID model, ERM solution $\hat{\mathbf{w}}$ satisfies

$$\mathbb{E}[\mathcal{R}(\hat{\mathbf{w}})] \leq \min_{\mathbf{w} \in \mathbb{R}^d} \mathcal{R}(\mathbf{w}) + O\left(\sqrt{\frac{d}{n}}\right).$$

- ▶ Unfortunately, solving this optimization problem, even for linear classifiers, is computationally intractable.
 - ▶ (Sharp contrast to ERM optimization problem for linear regression!)

11 / 27

Linearly separable data I

- ▶ Training data is linearly separable if there exists a linear classifier with training error rate zero.

- ▶ (Special case where ERM optimization problem is tractable.)
- ▶ There exists $\mathbf{w} \in \mathbb{R}^d$ such that $\text{sign}(\mathbf{x}_i^T \mathbf{w}) = y_i$ for all $i = 1, \dots, n$.
- ▶ Equivalent:

$$y_i \mathbf{x}_i^T \mathbf{w} > 0 \quad \text{for all } i = 1, \dots, n$$

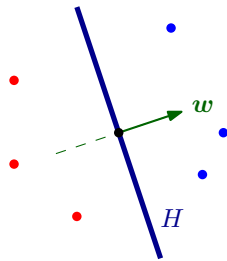


Figure 4: Linearly separable

12 / 27

Linearly separable data II

- ▶ Training data is linearly separable if there exists a linear classifier with training error rate zero.

- ▶ (Special case where ERM optimization problem is tractable.)
- ▶ There exists $\mathbf{w} \in \mathbb{R}^d$ such that $\text{sign}(\mathbf{x}_i^T \mathbf{w}) = y_i$ for all $i = 1, \dots, n$.
- ▶ Equivalent:

$$y_i \mathbf{x}_i^T \mathbf{w} > 0 \quad \text{for all } i = 1, \dots, n$$

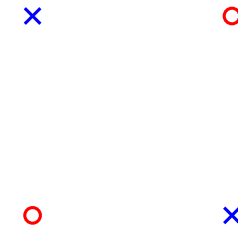


Figure 5: Not linearly separable

13 / 27

Finding a linear separator I

- ▶ Suppose training data $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n) \in \mathbb{R}^d \times \{-1, +1\}$ is linearly separable.
- ▶ How to find a linear separator (assuming one exists)?
- ▶ Method 1: solve linear feasibility problem

14 / 27

Finding a linear separator II

- ▶ Method 2: (approximately) solve logistic regression MLE

15 / 27

Surrogate loss functions I

- ▶ Often, a linear separator will not exist.
- ▶ Regard each term in negative log-likelihood as a “loss”

$$\ell_{\log}(s) := \ln(1 + \exp(-s))$$

- ▶ C.f. Zero-one loss:

$$\ell_{zo}(s) := \mathbb{1}_{\{s \leq 0\}}$$

- ▶ ℓ_{\log} (up to scaling) is upper-bound on ℓ_{zo} : a surrogate loss:

$$\ell_{zo}(s) \leq \frac{1}{\ln 2} \ell_{\log}(s) = \ell_{\log_2}(s).$$

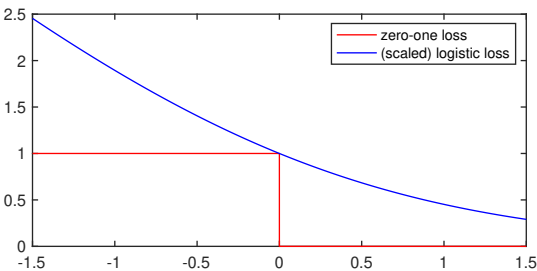


Figure 6: Logistic loss vs zero-one loss

16 / 27

Surrogate loss functions II

- ▶ Another example: squared loss

$$\ell_{sq}(s) = (1 - s)^2$$

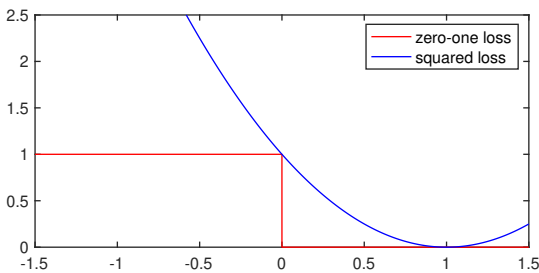


Figure 7: Squared loss vs zero-one loss

17 / 27

Surrogate loss functions III

- ▶ Modified squared loss:
 - ▶ $\ell_{msq}(s) := \max\{0, 1 - s\}^2$.

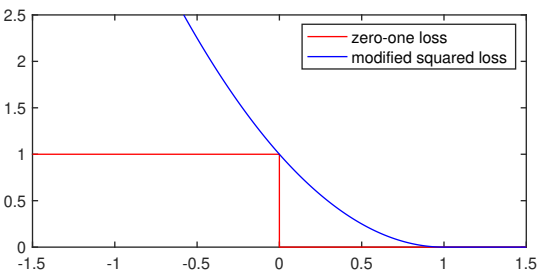


Figure 8: Modified squared loss vs zero-one loss

18 / 27

Gradient descent for logistic regression

- ▶ No “closed form” expression for logistic regression MLE
- ▶ Instead, use iterative algorithms like gradient descent.
- ▶ Gradient of empirical risk: by linearity of derivative operator,

$$\nabla \hat{\mathcal{R}}_{\ell_{\log}}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \nabla \ell_{\log}(y_i \mathbf{x}_i^T \mathbf{w}).$$

- ▶ Algorithm: start with some $\mathbf{w}^{(0)} \in \mathbb{R}^d$ and $\eta > 0$.
 - ▶ For $t = 1, 2, \dots$:

$$\mathbf{w}^{(t)} := \mathbf{w}^{(t-1)} - \eta \nabla \hat{\mathcal{R}}_{\ell_{\log}}(\mathbf{w}^{(t-1)})$$

19 / 27

Gradient of logistic loss

- ▶ Gradient of logistic loss on i -th training example: using chain rule,

$$\begin{aligned}\nabla \ell_{\log}(y_i \mathbf{x}_i^\top \mathbf{w}) &= \ell'_{\log}(y_i \mathbf{x}_i^\top \mathbf{w}) y_i \mathbf{x}_i \\ &= - \left(1 - \frac{1}{1 + \exp(-y_i \mathbf{x}_i^\top \mathbf{w})} \right) y_i \mathbf{x}_i \\ &= - (1 - \Pr_{\mathbf{w}}(Y = y_i \mid \mathbf{X} = \mathbf{x}_i)) y_i \mathbf{x}_i\end{aligned}$$

Here, $\Pr_{\mathbf{w}}$ is probability distribution of (\mathbf{X}, Y) from logistic regression model with parameter \mathbf{w} .

20 / 27

Interpretation of gradient descent for logistic regression

- ▶ Interpretation of gradient descent:

$$\nabla \ell_{\log}(y_i \mathbf{x}_i^\top \mathbf{w}) = - (1 - \Pr_{\mathbf{w}}(Y = y_i \mid \mathbf{X} = \mathbf{x}_i)) y_i \mathbf{x}_i$$

- ▶ Trying to make $\Pr_{\mathbf{w}}(Y = y_i \mid \mathbf{X} = \mathbf{x}_i)$ as close to 1 as possible.
- ▶ (Achieved by making \mathbf{w} infinitely far in direction of $y_i \mathbf{x}_i$.)
- ▶ How much of $y_i \mathbf{x}_i$ to add to \mathbf{w} is scaled by how far the $\Pr_{\mathbf{w}}(\cdots)$ currently is from 1.

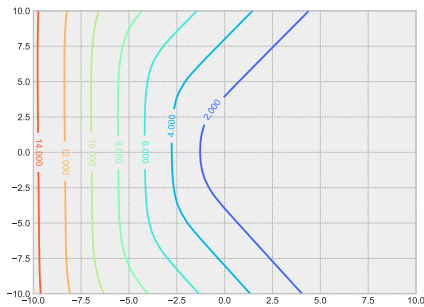
21 / 27

Behavior of gradient descent for logistic regression

- ▶ Analysis of gradient descent for logistic regression MLE much more complicated than for linear regression
 - ▶ Solution could be at infinity!
- ▶ **Theorem:** for appropriate choice of step size $\eta > 0$,

$$\hat{\mathcal{R}}(\mathbf{w}^{(t)}) \rightarrow \inf_{\mathbf{w} \in \mathbb{R}^d} \hat{\mathcal{R}}(\mathbf{w})$$

as $t \rightarrow \infty$ (even if the infimum is never attained).



Stochastic gradient method II

- ▶ Minibatch
 - ▶ To reduce variance of estimate, use several random examples J_1, \dots, J_B and average—called [minibatch gradient](#).

$$\frac{1}{B} \sum_{b=1}^B \nabla \ell(y_{J_b} \mathbf{x}_{J_b}^\top \mathbf{w}^{(t)}).$$

- ▶ Rule of thumb: larger batch size $B \rightarrow$ larger step size η .
- ▶ Alternative: instead of picking example uniformly at random, shuffle order of training examples, and take next example in this order.
 - ▶ Verify that expected value is same!
 - ▶ Seems to reduce variance as well, but not fully understood.

24 / 27

Example: SGD for logistic regression

- ▶ Logistic regression MLE for data $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n) \in \mathbb{R}^d \times \{-1, +1\}$.
- ▶ Start with $\mathbf{w}^{(0)} \in \mathbb{R}^d$, $\eta > 0$, $t = 1$
- ▶ For epoch $p = 1, 2, \dots$:
 - ▶ For each training example (\mathbf{x}, y) in a random order:

$$\mathbf{w}^{(t)} := \mathbf{w}^{(t-1)} + \eta \left(1 - \frac{1}{1 + \exp(-y \mathbf{x}^\top \mathbf{w}^{(t-1)})} \right) y \mathbf{x}; \quad t := t + 1.$$

- ▶ (If $\mathbf{w}^{(0)} = \mathbf{0}$, then final solution is in span of \mathbf{x}_i .)

25 / 27

Multinomial logistic regression

- ▶ How to handle $K > 2$ classes?
- ▶ [Multinomial logistic regression model](#)
 - ▶ $Y \mid \mathbf{X} = \mathbf{x}$ has a categorical distribution over $\{1, \dots, K\}$

$$\Pr(Y = k \mid \mathbf{X} = \mathbf{x}) = \text{softmax}(\mathbf{W} \mathbf{x})_k$$

- ▶ [Softmax function](#) (vector-valued function):

$$\text{softmax}(\mathbf{v})_k := \frac{\exp(v_k)}{\sum_{l=1}^K \exp(v_l)}$$

- ▶ $\mathbf{W} = [\mathbf{w}_1 \mid \dots \mid \mathbf{w}_K]^\top \in \mathbb{R}^{K \times d}$ is parameter matrix of interest

26 / 27

MLE for multinomial logistic regression

- ▶ Treat training examples as iid, same distribution as test example
- ▶ Encode label y_i as a [one-hot](#) vector $\tilde{\mathbf{y}}_i \in \{\mathbf{e}_1, \dots, \mathbf{e}_K\}$
 - ▶ $\tilde{\mathbf{y}}_i = (\tilde{y}_{i,1}, \dots, \tilde{y}_{i,K})$, where $\tilde{y}_{i,k} = \mathbb{1}_{\{y_i=k\}}$
- ▶ MLE equivalent to minimizing empirical risk with [cross-entropy loss](#)

$$\frac{1}{n} \sum_{i=1}^n \ell_{\text{ce}}(\tilde{\mathbf{y}}_i, \text{softmax}(\mathbf{W} \mathbf{x}_i))$$

where $\ell_{\text{ce}}(\mathbf{p}, \mathbf{q}) = - \sum_{k=1}^K p_k \ln q_k$.

27 / 27