

Assignment Project Exam Help

Machine learning lecture slides

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Classification I: Linear classification

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

- ▶ Logistic regression and linear classifiers
- ▶ Example: text classification
- ▶ Maximum likelihood estimation and empirical risk minimization
- ▶ Linear separators
- ▶

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Logistic regression model

- ▶ Suppose x is given by d real-valued features, so $x \in \mathbb{R}^d$, while $y \in \{-1, +1\}$.

- ▶ *Logistic regression model* for (X, Y)
 - ▶ $Y|X=x$ is Bernoulli (but taking values in $\{-1, +1\}$ rather than $\{0, 1\}$) with parameter $\sigma(x^\top w) := \frac{1}{1+\exp(-x^\top w)}$.

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

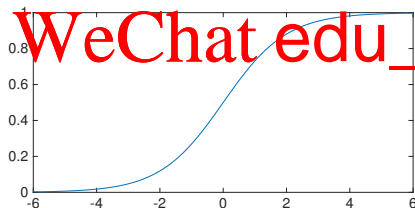


Figure 1: Logistic (sigmoid) function

Log-odds in logistic regression model

- Sigmoid function $\sigma(t) := 1/(1 + e^{-t})$

- Useful property: $1 - \sigma(t) = \sigma(-t)$

- $\Pr(Y = +1 | X = x) = \sigma(x^T w)$

- $\Pr(Y = -1 | X = x) = 1 - \sigma(x^T w) = \sigma(-x^T w)$

- Convenient formula: for each $y \in \{-1, +1\}$,

- <https://eduassistpro.github.io>

Add WeChat $\frac{\Pr(Y = +1 | X = x)}{\Pr(Y = -1 | X = x)}$ edu_assist_pro

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

Optimal classifier in logistic regression model

- Recall that Bayes classifier is

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



<https://eduassistpro.github.io>

Add WeChat $f(x) = \text{sign}(x^T w)$ edu_assist_pro

- This is a linear classifier
 - Compute linear combination of features, then check if above threshold (zero)
 - With affine feature expansion, threshold can be non-zero
- Many other statistical models for classification data lead to a linear (or affine) classifier, e.g., Naive Bayes

Geometry of linear classifiers

- ▶ Hyperplane specified by normal vector $w \in \mathbb{R}^d$:

- ▶ $H = \{x \in \mathbb{R}^d : x^\top w = 0\}$

- ▶ This is the decision boundary of a linear classifier

- ▶ Angle θ between x and w has

$$\cos \theta = \frac{x^\top w}{\|x\| \|w\|}$$

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro



Figure 3: Decision boundary of linear classifier

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

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

Figure 4: Decision boundary of linear classifier with quadratic feature expansion

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

Figure 5: Decision boundary of linear classifier with quadratic feature expansion (another one)

MLE for logistic regression

- ▶ Treat training examples as iid, same distribution as test example

- ▶ Log likelihood of w given data $(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R}^d \times \{-1, +1\}$:

<https://eduassistpro.github.io>

- ▶ No “closed form” expression for maximiz
- ▶ (Later, we'll discuss algorithms for finding maximizers using iterative methods like g

Example: Text classification (1)

- ▶ Data: articles posted to various internet message boards
- ▶ Label: -1 for articles on "religion", +1 for articles on "politics"
- ▶ Features
 - ▶ Vocabulary of $d = 61188$ words

$\{0, 1\}^d$,

- ▶ <https://eduassistpro.github.io>

$$\ln \frac{\Pr_w(Y = \text{politics} \mid \mathbf{x})}{\Pr_w(Y = \text{religion} \mid \mathbf{x})}$$

Add WeChat edu_assist_pro

- ▶ Each weight in weight vector w corresponds to a vocabulary word

Example: Text classification (2)

- ▶ Found \hat{w} that approximately maximizes likelihood given 3028 training examples
- ▶ Test error rate on 2017 examples is about 8.5%
- ▶ Vocabulary words with 10 highest (most positive) coefficients:

▶ <https://eduassistpro.github.io>

christ, athos

Add WeChat edu_assist_pr

450



Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Figure 6: Histogram of $\Pr_{\hat{w}}(Y = \text{politics} \mid X = x)$ values on test data

Example: Text classification (4)

- ▶ Article with $\Pr_{\hat{w}}(Y = \text{politics} \mid X = x) \approx 0.0$:

Rick, I think we can safely say, 1) Robert is not the only person who understands the Bible, and 2) the leadership of the LDS church historicly never has. Let's consider

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Example: Text classification (5)

- ▶ Article with $\Pr_{\hat{w}}(Y = \text{politics} \mid X = x) \approx 0.5$:

Does anyone know where I can access an online copy of the proposed 'jobs' or 'stimulus' legislation? Please E-mail me directly and if anyone else is interested, I can post this

- ▶ <https://eduassistpro.github.io>

titled "The Enemy Within" about the Anti-League.

Add WeChat edu_assist_pr

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

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

where loss function is zero-one loss.



<https://eduassistpro.github.io>

- Just like for linear regression, can apply plug
derive ERM, but now for linear classifiers

► Find $w \in \mathbb{R}^d$ to minimize

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

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

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

- <https://eduassistpro.github.io>

- Add WeChat [edu_assist_pro](#)
Unfortunately, solving this optimization problem for linear classifiers, is computationally intractable
 - (Sharp contrast to ERM optimization problem for linear regression!)

Linearly separable data

- ▶ Training data is linearly separable if there exists a linear classifier with training error rate zero.
- ▶ (Special case where FPM optimization problem is tractable.)
- ▶ There exists $w \in \mathbb{R}^d$ such that $\text{sign}(x_i^\top w) = y_i$ for all $i = 1, \dots, n$

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

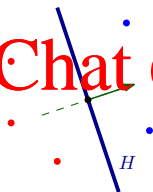


Figure 7: Linearly separable data

x

o

Assignment Project Exam Help

Figure 8: Data that is not linearly separable

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Finding a linear separator I

- ▶ Suppose training data $(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R}^d \times \{-1, +1\}$ is linearly separable.
- ▶ How to find a linear separator (assuming one exists)?
- ▶ Method 1: solve linear feasibility problem

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

- ▶ Method 2: approximately solve logistic regression MLE

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Surrogate loss functions I

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

Assignment Project Exam Help

- ▶ C.f. Zero-one loss: $\ell_{0/1}(s) := \mathbf{1}_{s \neq 0}$

- ▶ loss:

<https://eduassistpro.github.io>

- ▶ ℓ_{logistic} $\ell_{0/1\text{-risk}}$

Add WeChat edu_assist_pro

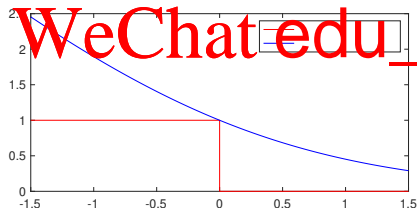


Figure 9: Comparing zero-one loss and (scaled) logistic loss

Surrogate loss functions II

- ▶ Another example: squared loss

- ▶ $\ell_{\text{sq}}(s) = (1 - s)^2$

- ▶ Note: $(1 - y_i x_i^T w)^2 = (y_i - x_i^T w)^2$ since $y_i \in \{-1, +1\}$

- ▶ Weild $\ell_{\text{sq}}(s) \rightarrow \infty$ as $s \rightarrow \infty$.

- ▶ Minimizing $\mathcal{R}_{\ell_{\text{sq}}}$ does not necessarily give a linear separator,

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

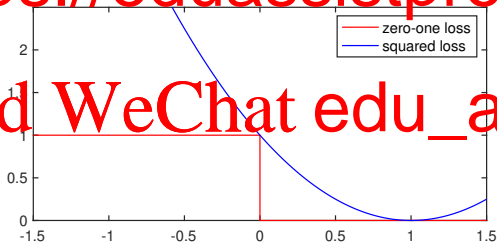
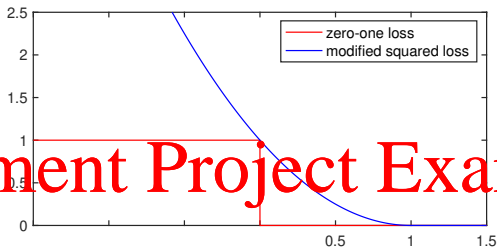


Figure 10: Comparing zero-one loss and squared loss



Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

(Regularized) empirical risk minimization for classification with surrogate losses

- ▶ We can combine these surrogate losses with regularizers, just as when we discussed linear regression
- ▶ This leads to regularized ERM objectives:

<https://eduassistpro.github.io>

where

- ▶ ℓ is a (surrogate) loss function
- ▶ Φ is a regularizer (e.g. $\Phi(w) = \|w\|^2$)