

Assignment Project Exam Help

Logistic Regression and MaxEnt

<https://eduassistpro.github.io>

Add WeChat

April 9, 2020

- Generative models:

Assignment Project Exam Help

$$\Pr[y | \mathbf{x}] = \frac{\Pr[\mathbf{x} | y] \Pr[y]}{\Pr[\mathbf{x}]}$$

<https://eduassistpro.github.io>

- Example: Naive Bayes.
- Discriminative models:
 - Models $\Pr[y | \mathbf{x}]$ directly as $g(\mathbf{x})$
 - Example: Decision tree, Logistic Regre
- Instance-based Learning.
 - Example: k NN classifier.

Add WeChat

Assignment Project Exam Help

<https://eduassistpro.github.io>

Figure: Linear Regres

Task

Add WeChat

- Input: $(x^{(i)}; y^{(i)})$ pairs $(1 \leq i \leq n)$
- Preprocess: let $\mathbf{x}^{(i)} = [1 \ x^{(i)}]^T$
- Output: The best $\mathbf{w} = [w_0 \ w_1]^T$ such that $\hat{y} = \mathbf{w}^T \mathbf{x}$ best explains the observations

Assignment Project Exam Help

The cr

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

Find \mathbf{w} such that \hat{y} is minimized.

Add WeChat

Taylor Series of $f(x)$ at point a

$$f(x) = \sum_{n=0}^{\infty} \frac{f^{(n)}(a)}{n!} (x-a)^n$$

<https://eduassistpro.github.io>

- close to x .
- If $f(x)$ has local minimum x , then
 - $f'(x) = 0$, and
 - $f''(x) > 0$.

Minimum of the local minima is the global minimum if it is smaller than the function values at all the boundary points.

- Intuitively, $f(x)$ is almost $f(a) + \frac{f''(a)}{2}(x-a)^2$ if a is close to x .

Assignment Project Exam Help

<https://eduassistpro.github.io>

By setting the above to 0, this essentially requires, f

Add WeChat

$$\sum_{i=1}^n \hat{y}^{(i)} x_j^{(i)} = \sum_{i=1}^n y^{(i)} x_j^{(i)}$$

what the model predicts

what the data says

Find the Least Square Fit for Linear Regression

In the simple 1D case, we have only two parameters in $\mathbf{w} = \begin{matrix} w_0 \\ w_1 \end{matrix}$

$$\sum_{i=1}^n (w_0 + w_1 x_1^{(i)}) x_0^{(i)} = \sum_{i=1}^n y^{(i)} x_0^{(i)}$$

<https://eduassistpro.github.io>

Since $x_0^{(i)} = 1$, they are essentially

$$\sum_{i=1}^n (w_0 + w_1 x_1^{(i)}) \cdot 1 =$$

$$\sum_{i=1}^n (w_0 + w_1 x_1^{(i)}) x_1^{(i)} = \sum_{i=1}^n y^{(i)} x_1^{(i)}$$

Example

Using the same example in [https://en.wikipedia.org/wiki/Linear_least_squares_\(mathematics\)](https://en.wikipedia.org/wiki/Linear_least_squares_(mathematics))

$$\mathbf{X} = \begin{bmatrix} 1 & 1 \\ 2 & 3 \\ 6 & 5 \\ 4 & 7 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} 1 \\ 6 \\ 5 \\ 7 \end{bmatrix}$$

<https://eduassistpro.github.io>

$$\begin{bmatrix} 1 & 1 \\ 2 & 3 \\ 6 & 5 \\ 4 & 7 \end{bmatrix} \begin{bmatrix} 2 \\ 3 \end{bmatrix} = \begin{bmatrix} 4 \\ 10 \end{bmatrix} \quad \mathbf{y}_3 = 5 \quad \mathbf{y}_4 = 7$$

- Easily generalizes to more than 2-dim:

Assignment Project Exam Help

$$\begin{matrix} 1 & x_1^{(1)} & \dots & x_m^{(1)} & & y^{(1)} \\ & 1 & \dots & \dots & w_0 & \vdots \\ & & & & & y^{(i)} \\ & & & & & \vdots \\ & & & & & y^{(n)} \end{matrix}$$

<https://eduassistpro.github.io>

- How to perform polynomial regression for x ?

• $y = w_0 + w_1 x + w_2 x^2 + \dots + w_m x^m$

- Let $x_j^{(i)} = (x_1^{(i)})^j$ Polynomi
(<http://mathworld.wolfram.com/LeastSquaresFittingPolynomial.html>)

High-level idea:

- Observations (i.e., training data) are noisy
- $P(y^{(i)} | \hat{y}^{(i)}) = f_i(\mathbf{w})$
- Any \mathbf{w} is possible, but some \mathbf{w} is most likely.

<https://eduassistpro.github.io>

- Maximum likelihood estimation (MLE)

- If we also incorporate some prior on

Maximum Posterior Estimation (

Gaussian prior on \mathbf{w} , this will add

the objective function.

- Many models and their variants can be deemed as different ways of estimating $P(y^{(i)} | \hat{y}^{(i)})$

Find \mathbf{w} such that $\|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2$ is minimized.

Assignment Project Exam Help

- What is $\mathbf{X}\mathbf{w}$ when \mathbf{X} is fixed?
 - It is the hyperplane spanned by the d column vectors of \mathbf{X} .

• <https://eduassistpro.github.io>

column of \mathbf{X} as X_i)

Add WeChat

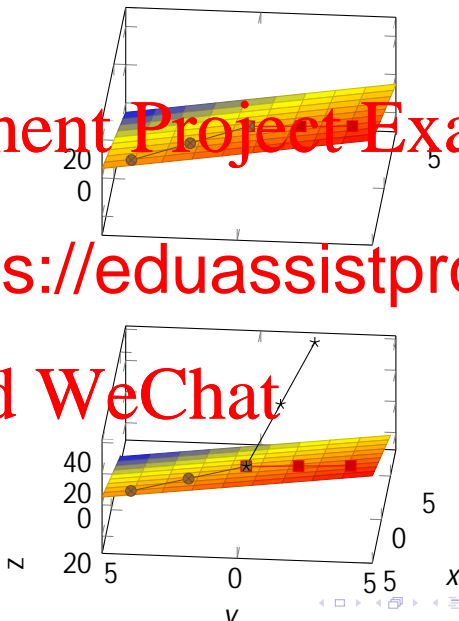
$$\begin{aligned} X_1^T(\mathbf{y} - \mathbf{X}\mathbf{w}) &= 0 \\ X_2^T(\mathbf{y} - \mathbf{X}\mathbf{w}) &= 0 \\ \vdots & \\ X_d^T(\mathbf{y} - \mathbf{X}\mathbf{w}) &= 0 \end{aligned} \Rightarrow \mathbf{X}^T(\mathbf{y} - \mathbf{X}\mathbf{w}) = \mathbf{0}$$

- $\mathbf{w} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y} = \mathbf{X}^+\mathbf{y}$ (\mathbf{X}^+ : pseudo inverse of \mathbf{X})

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat



Special case: $y^{(i)} \in \{0, 1\}$.

- Not appropriate to directly regress $y^{(i)}$.
- Rather, model $y^{(i)}$ as the observed outcome of a Bernoulli trial with an unknown parameter p_i .

<https://eduassistpro.github.io>

- What can we say about $p_{x^+} = \mathbf{w}^T \mathbf{x}$.
- Answer: we impose a linear relationship between $\mathbf{w}^T \mathbf{x}$ to $p_{\mathbf{x}}$.
 - What about a simple linear model
(Note: all points share the same parameter \mathbf{w})
 - Problem: mismatch of the domains: \mathbb{R} vs $[0, 1]$
 - Solution: mean function / inverse of link function:
 $g^{-1} : \mathbb{R} \rightarrow [0, 1]$ params

- Solution: Link function $g(\text{parameters}) = \logit(p)$

Assignment Project Exam Help

-

<https://eduassistpro.github.io>

Where $\sigma(z) = \frac{1}{1 + \exp(-z)}$.

Recall that $p_x = \mathbf{E}[y = 1 | \mathbf{x}]$.

- Decision boundary is $p = 0.5$.
 - Equivalent to whether $\mathbf{w}^T \mathbf{x} = 0$. Hence, LR is a linear classifier.

- Consider a training data point $\mathbf{x}^{(i)}$.
 - Recall that the conditional probability ($\Pr[y^{(i)} = 1 \mid \mathbf{x}^{(i)}]$) computed by the model is denoted by the shorthand notation p (which is a function of w and $\mathbf{x}^{(i)}$).
 - The likelihood of $\mathbf{x}^{(i)}$ is $p^{y^{(i)}}$, if $y^{(i)} = 1$, or equivalently,

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

$$L(w) = \prod_{i=1}^n p(\mathbf{x}^{(i)})^{y^{(i)}} (1 - p(\mathbf{x}^{(i)}))^{1 - y^{(i)}}$$

- Log-likelihood is (assume \log is \ln)

$$\ell(w) = \sum_{i=1}^n y^{(i)} \log p(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - p(\mathbf{x}^{(i)})) \quad (5)$$

- To maximize ℓ , notice that it is concave. So take its partial derivatives

Assignment Project Exam Help

$$\frac{\partial \ell(w)}{\partial w_j} = \sum_{i=1}^n y^{(i)} \frac{1}{1 + p(x^{(i)})} \frac{\partial p(x^{(i)})}{\partial w_j} + (1 - y^{(i)}) \frac{1}{1 + p(x^{(i)})} \frac{\partial (1 - p(x^{(i)}))}{\partial w_j}$$

<https://eduassistpro.github.io>

$i=1$

- and set them to 0 essentially means, for all

Add WeChat

$$\sum_{i=1}^n \hat{y}^{(i)} x^{(i)}_j = \sum_{i=1}^n p(x^{(i)}) x^{(i)}_j = \sum_{i=1}^n y^{(i)} x^{(i)}_j$$

what the model predicts

what the data says

- Consider one dimensional \mathbf{x} . The above condition is simplified

Assignment Project Exam Help

$$\sum_{i=1}^n x_{(i)} = \sum_{i=1}^n x_{(i)}$$

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

training data in class $Y = 1$.

- The LHS says: if we use our learned model to assign probability (of belonging to the class) for each training data, the LHS is the expected sum of
- If this is still abstract, think of an example.

Add WeChat

Assignment Project Exam Help

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

Add WeChat

- \mathbf{w} is initialized to some random value (e.g., $\mathbf{0}$).
- Since the gradient gives the *steepest* direction to increase a function's value, we move a small step towards that direction, i.e.,

Assignment Project Exam Help

<https://eduassistpro.github.io>

$i=1$

where η (learning rate) is usually a small value decreasing over the epochs.

- Stochastic version: using the gradient on a randomly selected training instance, i.e.,

$$\mathbf{w}_j \leftarrow \mathbf{w}_j + \eta (y^{(i)} - p(\mathbf{x}^{(i)})) \mathbf{x}_j^{(i)}$$

- Gradient Ascent moves to the "right" direction a tiny step a time. Can we find a good step size?

- Consider 1D case minimize $f(x)$ and the current point is a .

- $f(x) = f(a) + f'(a)(x - a) + \frac{f''(a)}{2}(x - a)^2$ for x near a .

<https://eduassistpro.github.io>

$$x = a - \frac{f'(a)}{f''(a)}$$

- Can be applied to multiple dimension cases too) need to use r (gradient) and Hess (Hessian).

- Regularization is another method to deal with **over fitting**.
 - It is designed to penalize large values of the model parameters.
 - Hence it *encourages* simpler models, which are less likely to over fit.
- Instead of optimizing for $\|w\|$, we optimize $\|w\| + R(w)$.

<https://eduassistpro.github.io>

- Grid search: http://scikit-learn.org/stable/tutorial/grid_search_parameter_tuning.html

- There are alternative methods.

- $R(w)$ quantifies the "size" of the model parameters. Choices are:

- L_2 regularization (**Ridge LR**) $R(w) = \|w\|_2^2$
- L_1 regularization (**Lasso LR**) $R(w) = \|w\|_1$
- L_1 regularization is more likely to result in sparse models.

Assignment Project Exam Help

- LR can be generalized to multiple classes \Rightarrow MaxEnt.

$$\frac{e^{c \cdot x}}{Z}$$

<https://eduassistpro.github.io>

- Z is the normalization constant.
- Let c be the last class in C , the
- Derive LR from MaxEnt. How?
- Both belong to *exponential* or *log*

Assignment Project Exam Help

- Andrew Ng's note:

1. pdf

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

- Tom Mitchell's book chapter: <http://www.cs.cmu.edu/~tom/mlbook>

Add WeChat