

## The general paradigm of supervised learning

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

- ▶ The goal of a large family of machine learning is to minimize the prediction errors of the model
  - ▶ Ideally we want to predict the true errors, error model when it is used in realistic scenarios
  - ▶ That is hard to do, so the common practice is to minimize the errors in a training set
  - ▶ To do that we need a metric, which is a metric that measures the errors in the prediction. I.
    - ▶ Cross-entropy Loss, Squared Error
- ▶ In other cases it is more natural to think of the problem is to optimize an *objective function*, e.g., Maximum Likelihood
- ▶ Whether to call it a loss function or objective function, there is no difference in how they are optimized

## Commonly used loss and objective functions in NLP

- ▶ Naïve Bayes:  $m$  independent and identically distributed samples of labeled samples

Assignment Project Exam Help

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} p(\mathbf{x}^{1:N}, \mathbf{y}^{1:N}; \theta)$$

- ▶ Logistic Regression: The weights are estimated by maximizing the log-likelihood function. **Maximum Condition**

Add WeChat edu\_assist\_pro

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

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \log$$

Add WeChat edu\_assist\_pro

- ▶ SVM: The weights are estimated by minimizing marginal loss

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^N \left( 1 - \gamma(\theta; \mathbf{x}^{(i)}, y^{(i)}) \right)_+$$

Note: Letters in bold indicates vector:  $\boldsymbol{\theta}$ ,  $\mathbf{x}$ ,  $\mathbf{f}$ . Alternative notations:  $\vec{\theta}$ ,  $\vec{x}$ ,  $\vec{f}$

## Naïve Bayes Objective

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

- Naïve Bayes. Maximize the joint probability of a training set of labeled samples, in a process called **likelihood Estimation**

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

$$\begin{aligned} &= \operatorname{argmax}_{\theta} \prod_{i=1}^N \\ &= \operatorname{argmax}_{\theta} \sum_{i=1}^N \log P(\mathbf{x}^i, y^i; \theta) \end{aligned}$$

## Logistic Regression Objective

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

- Logistic Regression: The weights are estimated by **Maximum Conditional Likelihood**

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \log p(y^{1:N}, x^{1:N} | \theta)$$
$$= \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^N \log \left( \exp(\theta \cdot \mathbf{f}(\mathbf{x}^{(i)}, y)) \right)$$

or by minimizing the **logistic loss**

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} - \sum_{i=1}^N \left( \theta \cdot \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) - \log \sum_{y \in \mathcal{Y}} \exp(\theta \cdot \mathbf{f}(\mathbf{x}^{(i)}, y)) \right)$$

## Support Vector Machine Objective

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

### Assignment Project Exam Help

- SVM: The weights are estimated by minimiz

$$\begin{aligned}\hat{\theta} &= \operatorname{argmin}_{\theta} \\ &= \operatorname{argmin}_{\theta} \sum_{i=1}^N (\max_{y \in \mathcal{Y}} (\theta \cdot \mathbf{f}(\mathbf{x}^{(i)}) - \theta \cdot \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}))_+)\end{aligned}$$

These look rather daunting, don't they?

How do we minimize a function?

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

## Assignment Project Exam Help

In order to minimize a function, we need to be able to compute the derivative, or rate of change of the function.

Let's start with a much simpler function and its derivative is:

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

$$\frac{d}{dx} f(x) = \frac{d}{dx} (x^2)$$

Add WeChat edu\_assist\_pro

"The derivative of the function  $f(x)$  with respect to (w.r.t.)  $x$ "  
This looks like magic, but it's really just calculus.

How do we find the minimum of a function with the derivative?

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

## Assignment Project Exam Help

- ▶ The derivative of a function can be interpreted at a certain point of the function.
- ▶ At the point that is the minimum (or maximum) of the function, the derivative is zero. We can find the derivative to zero:  
 $2x = 0, x = 0$
- ▶ For this particular function, there is a closed form solution. Most models in NLP don't have a closed form solution, but some do, e.g., Naïve Bayes.

Plot the function

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

Assignment Project Exam Help

Assignment Project Add WeChat edu\_assist\_pro

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

Add WeChat edu\_assist\_pro



## Finding the minimum iteratively

For functions that don't have a closed-form minimum, we can find the minimum iteratively.

from the input  $x$  so that the value of the function will decrease.

Suppose we start at the point where  $x = -1$ , and set the fraction  $\eta \triangleq 0.1$ , and  $\Delta x = \eta \frac{d}{dx} f(x)$ . So:

Assignment Project Exam Help  
Add WeChat: edu\_assist\_pro

$$x = x - \Delta x = -1 - 0.1 \times (-1.28) = -0.8$$

$$f(x) = (-0.8)^2 + 1 = 1.64$$

$$x = x - \Delta x = -0.8 - 0.1 \times (-1.28) = -0.64$$

$$f(x) = (-0.64)^2 + 1 = 1.4096$$

$$x = x - \Delta x = -0.64 - 0.1 \times (-1.28) = -0.512$$

$$f(x) = (-0.512)^2 + 1 = 1.262144$$

As  $x$  approaches 0,  $f(x)$  reaches the minimum, which is 1.

Finding the minimum iteratively

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

Assignment Project Exam Help

Assignment Project Add WeChat edu\_assist\_pro

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

Add WeChat edu\_assist\_pro

Finding the minimum iteratively

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

Assignment Project Exam Help

Assignment Project Add WeChat edu\_assist\_pro

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

Add WeChat edu\_assist\_pro

What if we try to learn fast using a larger learning rate?

Let's still start at  $x = 0$  the learning rate  $\eta \triangleq 1$  instead and see what happens

Assignment Project Exam Help

$$x = x - \Delta x = -1 - 1$$

Assignment Project Exam Help Add WeChat edu\_assist\_pro

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

$$x = x - \Delta x = -1 -$$

Assignment Project Exam Help Add WeChat edu\_assist\_pro

$$f(x) = (1)^2 + 1 = 2$$

So the  $x$  will just swing back and forth without ever reaching the minimum.

Setting the right learning rate is thus very important. If set improperly, we'll never reach the minimum, or at least take much longer than necessary.

Trying to learn fast with a larger learning rate

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

Assignment Project Exam Help

Assignment Project Add WeChat edu\_assist\_pro

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

Add WeChat edu\_assist\_pro

Trying to learn fast with a larger learning rate

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

Assignment Project Exam Help

Assignment Project Add WeChat edu\_assist\_pro

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

Add WeChat edu\_assist\_pro

## Derivative Rules

### Common Derivatives

$$\frac{d}{dx}(C) = 0$$

$$\text{e.g., } \frac{d}{dx}(91) = 0$$

$$\frac{d}{dx}(x) = 1$$

$$\frac{d}{dx}(x^n)$$

$$\frac{d}{dx}(a^x)$$

$$\frac{d}{dx}(e^x)$$

Assignment Project Exam Help

Add WeChat: edu\_assist\_pro

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

Add WeChat: edu\_assist\_pro

Note: ln: "Natural logarithm", logarithm to base of the mathematic constant e, where  $e = 2.71882 \dots$

## Derivative rules

More common derivatives

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

$$\frac{d}{dx}(\ln(x)) = \frac{1}{x}, x > 0$$

$$\frac{d}{dx}(\ln(|x|)) = \frac{1}{x}, x \neq 0$$

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

$$\frac{d}{dx}(\sin(x)) = \cos(x)$$

Add WeChat edu\_assist\_pro

$$\frac{d}{dx}(\cos(x)) = -\sin(x)$$

$$\frac{d}{dx}(\tan(x)) = \sec^2(x)$$

Note: When  $x \leq 0$ ,  $\ln(x)$  is unspecified. That is, you can't raise the constant  $e$  to any value to get a zero or a negative number.



## Derivatives of functions

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

“The derivative of the function with respect to  $x$ ”

Assignment Project Exam Help

$$\frac{d}{dx}(cf(x)) = c \frac{d}{dx}f(x)$$

$$\frac{d}{dx}(f(x) \pm g(x)) = \frac{d}{dx}f(x) \pm \frac{d}{dx}g(x)$$

$$\frac{d}{dx}(f(x)g(x)) = f(x) \frac{d}{dx}g(x) + g(x) \frac{d}{dx}f(x) \quad (\text{Product rule})$$

$$\frac{d}{dx} \frac{f(x)}{g(x)} = \frac{g(x) \frac{d}{dx}f(x) - f(x) \frac{d}{dx}g(x)}{g^2} \quad (\text{Quotient rule})$$

$$\frac{d}{dx}f(g(x)) = \frac{d}{dg(x)}f(g(x)) \frac{d}{dx}g(x) \quad (\text{Chain rule})$$

Breaking down the derivative of complex functions

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

Using these fundamental derivative rules, and particularly the chain rule, you can break down more complicated functions

Assignment Project Exam Help  
Add WeChat edu\_assist\_pro

$\frac{d}{dx} e^{f(x)} = e^{f(x)} \frac{d}{dx} f(x)$   
<https://eduassistpro.github.io/>

$\frac{d}{dx} \ln(f(x)) = \frac{1}{f(x)} \frac{d}{dx} f(x)$   
Add WeChat edu\_assist\_pro

## Partial Derivatives

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

### Assignment Project Exam Help

- ▶ We don't normally deal with single variable functions. A typical NLP model (function) has tens of thousands of variables (features). So we need to compute *partial derivatives*.
- ▶ Fortunately, it's quite simple. You just need to hold all other variables constant, and take the derivative with respect to the variable.  $\frac{\partial}{\partial x} f(x, y)$ ,  $\frac{\partial}{\partial y} f(x, y)$

More on partial derivatives

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

Assignment Project Exam Help

Assignment Project Exam Help  
Add WeChat edu\_assist\_pro

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

$\frac{\partial}{\partial y} f(x, y)$   
Add WeChat edu\_assist\_pro

More on partial derivatives

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

Assignment Project Exam Help

$f(x, y) = \min(x, y) = \begin{cases} x & \text{if } x \leq y \\ y & \text{if } x > y \end{cases}$

$\frac{\partial}{\partial x} f(x, y) = \begin{cases} 1, & \text{if } x < y \\ 0, & \text{if } x > y \end{cases}$

$\frac{\partial}{\partial y} f(x, y) = \begin{cases} 0, & \text{if } x < y \\ 1, & \text{if } x > y \end{cases}$

The function is not differentiable when  $x = y$

Plot multi-variable functions

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

Assignment Project Exam Help

Assignment Project Add WeChat edu\_assist\_pro

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

Add WeChat edu\_assist\_pro

## Gradient

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

## Assignment Project Exam Help

The gradient of a function  $\nabla f$  is the set of partial derivatives of the function

[Add WeChat edu\\_assist\\_pro](#)

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

$$\nabla f = \left[ \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right]$$
[Add WeChat edu\\_assist\\_pro](#)

## Properties of Logarithms

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

$$\log(xy) = \log(x) + \log(y) \quad \ln(e^x) = x$$

$$\log\left(\frac{x}{y}\right) = \log(x) - \log(y) \quad \ln(x) > 0$$

$$\log(x^y) = y \log(x)$$

$$\log\left(\prod_i x_i\right) = \sum_i \log(x_i)$$

Add WeChat edu\_assist\_pro

- ▶ It is common practice to map probabilities to logarithmic space to avoid *underflow* (when a value gets too close to zero for the computer to represent it).  
 $\ln(0.0001) = -9.2103403 \dots$
- ▶ You can map the log values back to probabilities using the exponent.  $e^{-9.2103403} = 0.0001$



## Convexity of functions

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

- Assignment Project Exam Help  
Add WeChat edu\_assist\_pro
- ▶ Intuitively, a convex (conclave) function is a continuous function in which there is a single minimum (m
  - ▶ A mathematical definition: A convex function is a function who domain does not have the ends of the interval.
  - ▶ How to decide a function is convex. If the first derivative in  $[a, b]$ , then a necessary condition for it to be convex on that interval is that the second derivative  $f''(x) \geq 0$  for all  $x$  in  $[a, b]$ .

Example convex functions

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

Assignment Project Exam Help

Assignment Project Add WeChat edu\_assist\_pro

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

Add WeChat edu\_assist\_pro

Example non-convex functions

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

Assignment Project Exam Help

Assignment Project Add WeChat edu\_assist\_pro

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

Add WeChat edu\_assist\_pro