

## 1. Math and Probability Basics

### Q1.1 Definitions

[a] Give the definition of an orthogonal matrix.

<https://powcoder.com>

[b] Give the definition of an eigenvector and eigenvalue.

Assignment Project Exam Help

Assignment Project Exam Help

[c] How is the probability density function different from the cumulative probability distribution?

<https://powcoder.com>

Add WeChat powcoder

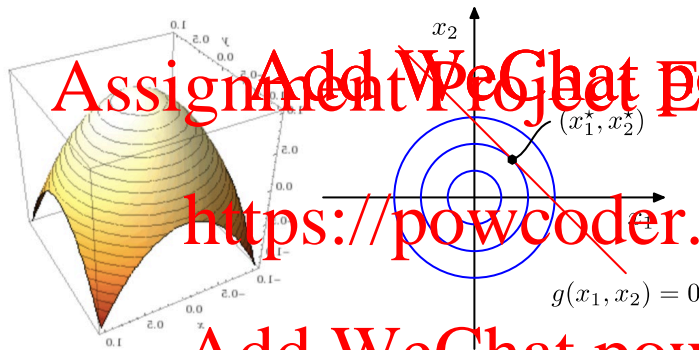
[d] What is a 'singular' matrix?

[e] Give the definition of Baye's Rule.

### Q1.2 Short questions

- a) Supervised classification models can be used to predict if an email attachment contains a computer virus [T/F]
- b) Suppose we use polynomial features for linear regression, then the hypothesis is linear in the original features [T/F]

- c) Find the maximum of the function  $f(x_1, x_2) = 1 - x_1^2 - x_2^2$  subject to the constraint  $g(x_1, x_2) = x_1 + x_2 - 1 = 0$ , using Lagrange Multipliers. The 3-D plot shows  $f$  and the 2-D plot shows its contours (blue circles) and the line corresponding to  $g = 0$  (red line).



### Q1.3 Covariance Matrix

Recall that in the derivation of normal equations in class, we used the fact that the data covariance matrix, i.e. a matrix whose element in the  $k, j$  position is the covariance between the  $k$ -th and  $j$ -th elements of the random input vector, is given by

$$\sum_{i=1}^m x^{(i)} x^{(i)T} = X^T X$$

where  $x^{(i)}$  is the  $i$ th  $n \times 1$  input vector,  $m$  is the number of input vectors in the dataset, and  $X$  is the  $m \times n$  design matrix. Show that the above equality is true. Show all your steps.

Assignment Project Exam Help

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

### Q1.4 Matrix Norm

The trace of a square matrix  $A \in \mathbb{R}^{n \times n}$  is defined as the sum of diagonal entries, or

$$\text{tr} A = \sum_{i=1}^n A_{ii}.$$

Prove the following fact

$$\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n A_{ij}^2} = \sqrt{\text{tr}(A^T A)}.$$

where  $\|A\|_F$  is the matrix Frobenius norm. Show all your steps.

Assignment Project Exam Help

Assignment Project Exam Help

### Q1.4 Flu Virus Test

After your yearly checkup, the doctor has bad news and good news. The bad news is that you tested positive for a flu virus, and that the test is 99% accurate (i.e., the probability of testing positive given that you have the virus is 0.99, as is the probability of testing negative given that you don't have the disease). The good news is that this is a rare virus, striking only one in 10,000 people ( $10^{-4}$ ). What are the chances that you actually have the disease? Show your calculations as well as giving the final result. *Hint: use Baye's Rule.*

## 2. Gradient Descent

### Q2.1 Gradient Descent for a Cost Function

Suppose we have a cost function

$$J(\theta) = \frac{1}{m} (\sum_{i=1}^m x_i^T \theta + b y_i) + \frac{1}{2} \theta^T A \theta,$$

where  $\theta \in \mathbb{R}^n$  is the parameter vector,  $x_i \in \mathbb{R}^n$ ,  $y_i \in \mathbb{R}$ ,  $\{x_i, y_i\}$  are  $m$  training data points,  $A \in \mathbb{R}^{n \times n}$  is a symmetric matrix, and  $b \in \mathbb{R}$ . We want to find parameters  $\theta$  using gradient descent.

- a) [3 points] Give the pseudocode for the gradient descent algorithm for a **generic** cost function  $J(\theta)$  (not the specific one above).

Assignment Project Exam Help

- b) [3 points] For the specific function above, what is the vector of partial gradients of the cost function, i.e. the vector with the  $j$ th element equal to  $\frac{\partial}{\partial \theta_j} J(\theta)$ ?

<https://powcoder.com>

Add WeChat powcoder

- c) [3 points] What is the design matrix? Describe its entries and give its dimensions.

- d) [3 points] Re-write the expression for the gradient without using the summation notation  $\sum$ .  
*Hint: use the design matrix  $X$ .*

- e) [3 points] Suppose we run gradient descent for two iterations. Give the expression for  $\theta$  after two updates, with step size  $\alpha = 1$  and initial value of  $\theta = \mathbf{0}$  (vector of zeros).



- f) [3 points] How do we know when the algorithm has converged?

<https://powcoder.com>

Assignment Project Exam Help

- g) [3 points] Give the closed-form solution for  $\theta$ . You do not need to prove it is the minimum of the cost.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

### 3. Regression and Classification

#### Q3.1 Linear Regression: Online Art Auction

Imagine you work for an online art auctioneer. You would like to estimate the price  $y$  that a piece of art will sell for in an auction (in dollars), based on the following features:

- $x_1$  = type of art (out of 15 types such as 1:painting, 2:sculpture, etc.),
- $x_2$  = artist popularity (rank out of 100 artists),
- $x_3$  = estimated value in dollars,
- $x_4$  = days in the auction,
- $x_5$  = previously owned (binary),
- $x_6$  = is abstract (binary),
- ...etc.

For example, a feature vector for the  $i^{th}$  item could be  $x^{(i)} = [1, 28, 1700, 5, 1, 0, \dots]$ . You have collected data points from previous auction sales,  $(x^{(i)}, y^{(i)})$ ,  $i = 1, \dots, m$ .

- a. [3 points] You decide to use a linear regression model,  $y = \sum_{j=0}^n \theta_j x_j$ . In what circumstances should you use gradient descent to find the parameters?

- b. [3 points] Suppose you decide to use gradient descent. How can you tell if it is converging?

- c. [3 points] Suppose you're monitoring convergence and find it is slow. Name two things you can try to speed it up (be specific).

- d. [3 points] You want to add new features to improve your predictor. You consider adding total minutes spent in the auction. Is this a good idea? Why or why not?

- e. [3 points] Your boss does not know how much to trust your prediction  $y^{test} = \$15,000$  for a certain watercolor painting. She asks you to estimate the probability of the painting selling for more than \$20,000. Give the equation for this probability, using a linear regression model that assumes the outputs have Gaussian noise with variance  $\beta^{-1}$ .

<https://powcoder.com>

## Assignment Project Exam Help

- f. [2 points] What is the probability the painting will sell for more than \$15,000?

[Add WeChat powcoder](https://powcoder.com)

<https://powcoder.com>

- g. [3 points] Suppose you now want to estimate the probability that a piece of art does not get any bids,  $p(no\_bids|x)$ , based on historic data. What sort of features and machine learning method should you use?



### Q3.2 Softmax Classifier

Typically, we use the *softmax* function to model the probability of one of several classes given the input. Instead, consider what would happen if a binary classifier with output labels  $j \in \{0,1\}$  used the softmax function to model the probability of the binary label  $y = j$  given the input:

$$P(y = j|x) = \frac{e^{w_j^T x}}{\sum_{j=0,1} e^{w_j^T x}}$$

Where  $w_j$  is the parameter vector of the  $j$ th class.

- a) [5 points] Show that the solution to this problem can be expressed as a logistic regression classifier with parameter  $w$ , and give the expression for  $w$ .

Assignment Project Exam Help

Assignment Project Exam Help

<https://powcoder.com>

- b) [5 points] Show that the posterior  $P(y = j|x)$  is invariant if a constant vector  $b$  is added to both weight vectors  $w_j$ .

Add WeChat powcoder

- c) [5 points] (b) implies that the solution  $w_j$  is not unique. We can guarantee a unique solution by making the objective function regularized, e.g. using the squared norm regularizer. Write down the objective function and say whether you should minimize or maximize it.

## 4. Overfitting and Regularization

### Q4.1 Bias-Variance and $\lambda$

Alice has a binary classification dataset of  $m$  points with  $n$ -dimensional inputs.

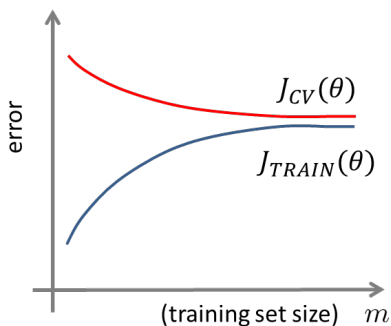
- a) [3 points] She has trained several regularized logistic regression models using regularization parameters  $\lambda = 10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}$ . She computed the cross-validation (CV) and training errors for each value of  $\lambda$ , shown in the table below, but the rows are out of order. Fill in the correct values of  $\lambda$  for each row.

Train error	CV error	$\lambda$
80%	85%	
40%	45%	
70%	76%	
35%	50%	

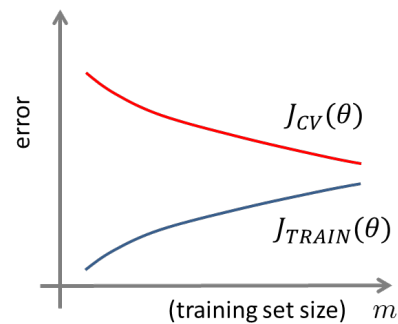
- b) [3 points] Based on these results, which  $\lambda$  should she choose, and why?

- c) [3 points] Which of the four models will have the highest error due to variance? Why?

- d) [3 points] Alice also plotted learning curves for the models with  $\lambda = 10^{-1}, 10^{-4}$ . Match each plot with the correct value, and explain why it matches.



$\lambda =$



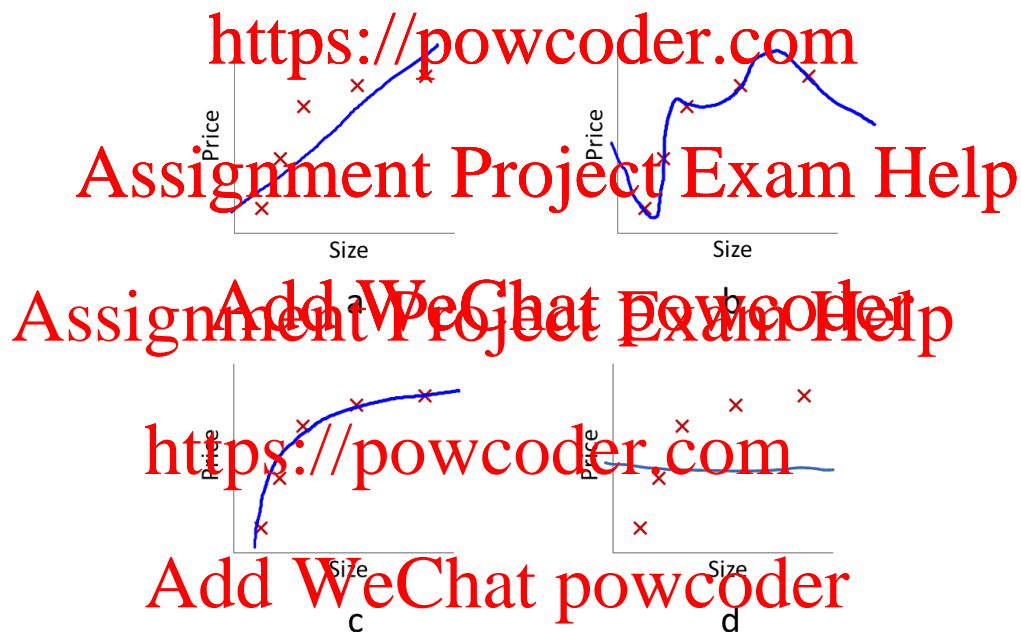
$\lambda =$

### Q4.2 Regularization for Linear Regression

Alice is trying to fit a linear regression model to predict house price based on size using polynomial features. Since her training dataset is very small, she is applying regularization. She fit several models by minimizing the cost function

$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

for  $\lambda = 10^0, 10^1, 10^2, 10^3$ . The following are sketches of the resulting models.



- [3 points] Which value of  $\lambda$  goes with each of the plots? (Write it next to the plot)
- [3 points] Alice tries her model on a test set. Which model will have the highest error due to bias?
- [3 points] Which model will have the highest error due to variance?
- [3 points] Which model, if any, will always have zero test error?

## 5. Maximum Likelihood Principle

### Q5.1 ML for Probabilistic Linear Regression

Recall that probabilistic linear regression defines the likelihood of observing outputs  $t^{(i)} \in \mathbb{R}$  given inputs  $x^{(i)} \in \mathbb{R}^p$ , where  $i = 1, \dots, m$  and  $m$  is the number of samples in the dataset, as

$$p(t_1, \dots, t_m | x_1, \dots, x_m, \theta, \beta) = \prod_{i=1}^m N(t^{(i)} | h(x^{(i)}), \beta^{-1})$$

where  $h(x)$  is the linear regression hypothesis,  $\theta, \beta$  are parameters and  $N(x | \mu, \sigma^2)$  is the normal (Gaussian) probability density with mean  $\mu$  and variance  $\sigma^2$ . Here  $\beta \in \mathbb{R}^+$  is the inverse variance of the Gaussian noise that we assume is added to the data.

(a) [8 points] Find  $\beta_{ML}$ , the maximum likelihood solution for  $\beta$ . *Hint: maximize log likelihood with respect to only  $\beta$ .*

(b) [2 points] What is the interpretation of the solution  $\beta_{ML}$ ? Explain in one sentence.

## Q5.2 ML for Linear Regression with Multivariate Outputs

Consider a probabilistic linear regression model for a multivariate  $p$ -dimensional target variable  $t = [t_1, \dots, t_p]^T$  that has a Gaussian distribution over  $t$  of the form

$$p(t|W, \Sigma) = N(t|y(\phi(x), W), \Sigma)$$

where  $\phi(x)$  is a basis function representation of the input, and

$$y(x, W) = W^T \phi(x)$$

$W$  is a matrix of parameters and  $\Sigma$  is the covariance parameter matrix. We are given a training dataset of input basis vectors  $x_n$  and corresponding target vectors  $t_n$ ,  $n = 1, \dots, N$ . Show that the Maximum Likelihood solution  $W_{ML}$  for the parameter matrix  $W$  has the property that each column is given by

$$W_{ML} = (\Phi^T \Phi)^{-1} \Phi^T t$$

where  $\Phi$  is the design matrix. You do not need to provide a solution for  $\Sigma$ . Show all your steps..

Hint: the  $p$ -dimensional multivariate normal distribution is given by

$$N(x; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left( -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right)$$

Hint: you may also find some of the matrix differentiation rules in the appendix helpful.

Add WeChat powcoder

### Q5.3 ML for Poisson Regression

This problem asks you to derive the maximum likelihood solution for a Poisson hypothesis.

- a) [3 points] Given a training set  $\{x_i, y_i\}$ , Poisson regression models the probability of observing an output given an input as  $p(y_i|\lambda) = \frac{1}{y_i!} \lambda^{y_i} e^{-\lambda}$  where  $\lambda = e^{\theta^T x_i}$  for some parameter vector  $\theta$ . Derive the cost function  $J(\theta)$  corresponding to maximizing the log likelihood for a **single training example**.

<https://powcoder.com>

Assignment Project Exam Help

- b) [3 points] The cost function in (b) has no closed form solution, so we must use an iterative method. Show the stochastic gradient descent update for this cost function.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

## 6. Unsupervised Learning

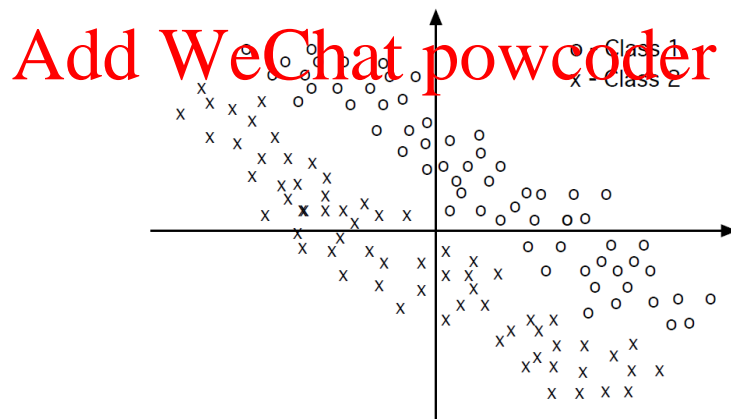
### Q6.1 Principle Component Analysis

- a) PCA assumes a specific relationship between the unobserved latent coordinates  $z$  and the observed data points  $x$ . Express this relationship as an equation. Clearly identify and name the parameters which are learned.

- b) Name one objective function which could be minimized to learn the parameters of PCA.

- c) For a dataset of arbitrary points  $x^{(1)}, \dots, x^{(m)}$ , specify the steps of the PCA algorithm.

- d) Suppose you are given 2D feature vectors for a classification task which are distributed according to the figure below. You apply PCA to the entire dataset. On the figure, draw all the PCA components.



- e) In (d) above, could you use PCA components directly to classify the data (without training a classifier)? Explain.

## Q6.2 Gaussian Mixture Models

- a) Describe in words the two main steps of the Expectation Maximization algorithm used to solve Gaussian Mixture Models.
- b) True or False: In the case of fully observed data, i.e. when all latent variables are observed, EM reduces to Maximum Likelihood.

<https://powcoder.com>

- c) True or False: Since the EM algorithm guarantees that the value of its objective function will increase after each iteration, it is guaranteed to eventually reach the global maximum.

Assignment Project Exam Help

Add WeChat powcoder

- d) Sketch a dataset on which K-Means would work poorly but a Gaussian Mixture Model with the same number of clusters would do well. Describe why K-Means wouldn't work well.

<https://powcoder.com>

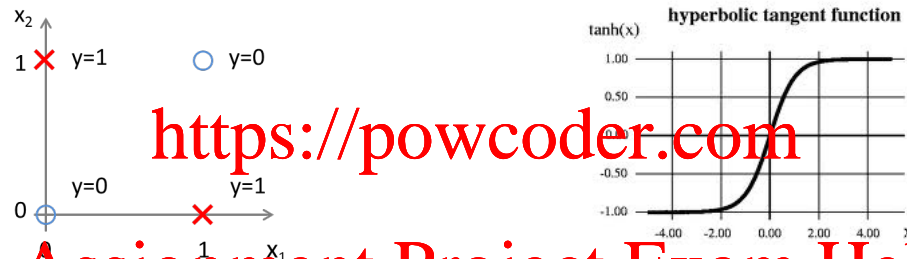
Add WeChat powcoder



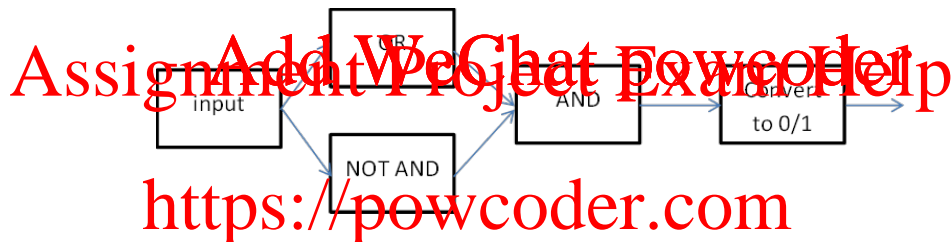
## 7. Neural Networks

### Q7.1 Neural Network for XOR

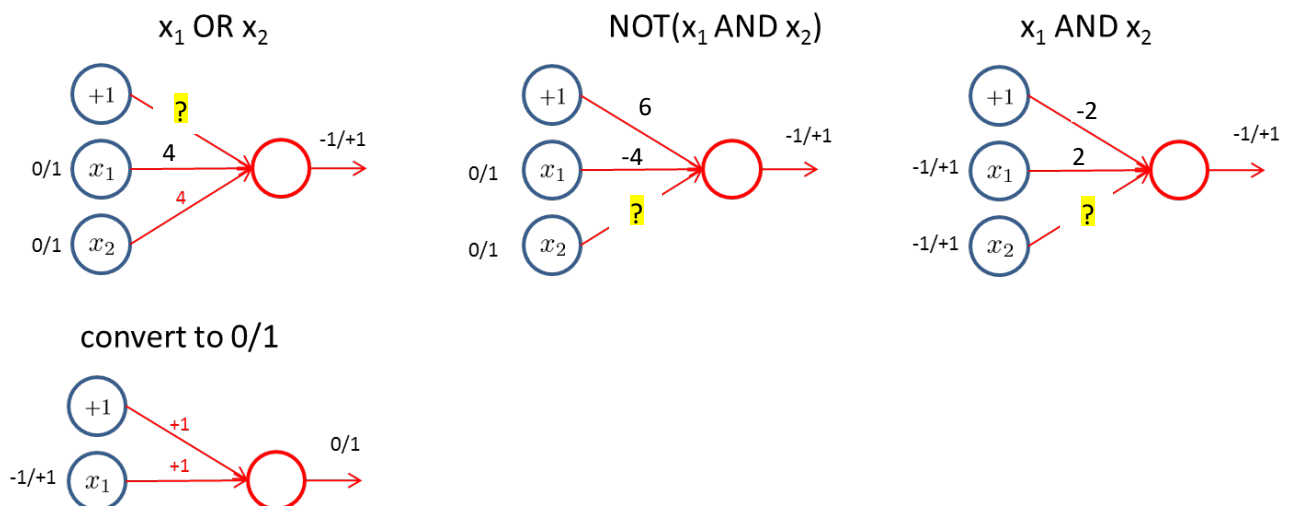
Design a neural network to solve the XOR problem, i.e. the network should output 1 if only one of the two binary input variables is 1, and 0 otherwise (see left figure). Use the hyperbolic tangent, or *tanh*, activation function in all nodes (right figure), which ranges in  $[-1, +1]$ .



Note that  $(A \text{ XOR } B)$  can be expressed as  $(A \text{ OR } B) \text{ AND NOT}(A \text{ AND } B)$ , as illustrated below:



In the diagrams below, we filled in most of the tanh units' parameters. Fill in the remaining parameters, keeping in mind that tanh outputs  $+1/-1$ , not 0/1. Note that we need to appropriately change the second layer (the AND node) to take  $+1/-1$  as inputs. Also, we must add an extra last layer to convert the final output from  $+1/-1$  to 0/1. Hint: assume tanh outputs  $-1$  for any input  $x \leq -2$ ,  $+1$  for any input  $x \geq +2$ , 0 for  $x = 0$ .



### Q7.2 Computation Graph and Backpropagation

In class, we learned how to take a complex function that consists of multiple nested functions and represent it with a computation graph, which allows us to write down the forward and backward pass used to compute the function gradient.

- a) Practice converting different functions  $f_{\theta}(x) = f_k(f_{k-1}(\dots f_1(x)))$  of input vector  $x$  parametrized by  $\theta$  to their computation graphs.

- b) For the computation graphs obtained in (a), write down the forward pass and the backward pass equations.

<https://powcoder.com>  
Assignment Project Exam Help

### Q7.3 Neural Network Architectures

- a) Draw a convolutional network with input  $x \in \mathbb{R}^2$ , one hidden layer with  $2 \times 1$  filters and 2 channels with stride 2, and a fully-connected output layer with one neuron. How many parameters does this network have?

<https://powcoder.com>

Add WeChat powcoder

- b) What algorithm is used for learning the parameters of a recurrent network? Name the algorithm and sketch out its main steps.

## Appendix: Useful Formulas

### Matrix Derivatives

For vectors  $x$ ,  $y$  and matrix  $A$ ,

$$y = Ax, \text{ then } \frac{\partial y}{\partial x} = A$$

$$\text{If } z = x^T A x, \text{ then } \frac{\partial z}{\partial x} = x^T (A + A^T). \text{ For the special case of a symmetric matrix } A, \frac{\partial z}{\partial x} = 2x^T A.$$

$$\text{Chain Rule: if } z \text{ is a function of } y, \text{ which is a function of } A, \text{ then } \frac{\partial z}{\partial A} = \frac{\partial y}{\partial A} \frac{\partial z}{\partial y} \text{ (note the order).}$$

<https://powcoder.com>

### Single-Dimension Normal Distribution

$$N(x; \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

### Multivariate Normal Distribution

The  $p$ -dimensional multivariate normal distribution with mean  $\mu$  and covariance matrix  $\Sigma$  is given by

$$N(x; \mu, \Sigma) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu)\right)$$

Add WeChat powcoder