

Problem 6

a)

The basic function of activation function in neural network is to introduce nonlinearity. There are many choices for the specific nonlinear form. The advantage of sigmoid is that the output range is limited, so the data is not easy to diverge in the process of transmission. Of course, there is also a corresponding disadvantage, that is, the gradient is too small at saturation. Another advantage of sigmoid is that the output range is $(0, 1)$, so it can be used as an output layer, and the output represents probability. Besides, another advantage of sigmoid is easy derivation.

b)

i)

Sigmoid and tanh are "saturated activation functions", while relu and its variants are "unsaturated activation functions". The advantages of using "unsaturated activation function" lie in two points:

(1) "unsaturated activation function" can solve the so-called "gradient disappearance" problem.

(2) It can accelerate the convergence speed.

Sigmoid compresses the real value output in the range of $[0,1]$, and the tanh function compresses the real value output in the range of $[-1,1]$.

Sigmoid function was used very often in history, and the output value range is real number between $[0,1]$. But now it's not very popular and it's rarely used in practice. The reasons are as follows:

(1) The saturation of sigmoid function makes the gradients disappear.

(2) The sigmoid function output is not "zero centered".

(3) The calculation of exponential function consumes more computing resources.

tanh function

tanh function is very similar to sigmoid in shape. In fact, tanh is a deformation of sigmoid,

Unlike sigmoid, tanh is "zero centered.". Therefore, in practical application, tanh is better than sigmoid. However, in the case of saturated neurons, tanh did not solve the problem of gradient disappearance.

Advantages: (1) tanh solves the problem that the output of sigmoid is not "Zero Center"

Disadvantages: (1) there is still the problem of supersaturation of sigmoid function. (2) It's still exponential.

ReLU

In recent years, relu functions have become more and more popular. Its full name is rectified linear unit, and its Chinese name is modified linear unit.

Advantages: (1) relu solves the problem of gradient disappearance, at least x is in the positive range, neurons will not be saturated;

(2) Due to the linear and unsaturated form of relu, it can converge quickly in SGD;

(3) The calculation speed is much faster. Relu function has only linear relationship, and does not need exponential calculation. It is faster than sigmoid and tanh in both forward and backward propagation.

Disadvantages: (1) the output of relu is not "Zero Center";

(2) With the training going on, the neuron may die and the weight cannot be updated. This neuronal death is irreversible.

ii)

Sigmoid functions and their combination usually work better in classifiers. Due to the problem of gradient collapse, sigmoid and tanh activation functions need to be avoided in some cases. Relu function is a common activation function, which is most used at present. If we encounter some dead neurons, we can use the leaky relu function, which is always used only in hidden layers. According to experience, we can generally start with the relu

activation function, but if relu can't solve the problem well, we can try other activation functions

c)

The idea of hinge loss is to make the distance between those who fail to classify correctly and correctly classify sufficiently. If the difference reaches a threshold value of Δ Delta Δ , the error of incorrect classification can be considered as 0, otherwise, the calculation error will be accumulated. Hinge loss function is mainly used in support vector machine (SVM), Advantages: stable classification surface, convex function. The classification interval can be maximized.

The 0-1 loss function is mainly used for perceptron

The quadratic loss function is mainly used in least squares (OLS)

Logarithmic loss function (logarithmic loss function, cross entropy loss function, softmax loss function) is mainly used in logistic regression and softmax classification

Exponential loss function is mainly used in AdaBoost ensemble learning algorithm

d)

Jacobian

Sometimes we need to find all of the partial derivatives of a function whose input and output are both vectors. The matrix containing all such partial derivatives is the Jacobian.

Given:

$$\vec{f}: \mathbb{R}^m \rightarrow \mathbb{R}^n$$

The Jacobian matrix **J** is given by:

$$\mathbf{J} \in \mathbb{R}^{n \times m}$$

$$J_{i,j} = \frac{\partial}{\partial x_j} f(\mathbf{x})_i$$

Example of Jacobian Matrix

Let's say:

$$\vec{f}: \mathbb{R}^2 \rightarrow \mathbb{R}^2$$

The function f , takes in a vector of size 2 and outputs a vector of size 2. The operation can be explained by 2 functions in a vector:

$$\begin{bmatrix} x^2y \\ 5x + \sin(y) \end{bmatrix}$$

Where:

$$f_1(x, y) = x^2y$$

$$f_2(x, y) = 5x + \sin(y)$$

And the Jacobian Matrix of \mathbf{f} is:

$$J_f(x, y) = \begin{bmatrix} \frac{\partial f_1}{\partial x} & \frac{\partial f_1}{\partial y} \\ \frac{\partial f_2}{\partial x} & \frac{\partial f_2}{\partial y} \end{bmatrix}$$

Hessian

Hessian is a square matrix of second order partial derivatives of a scalar-valued function or scalar field.

It describes the local curvature of a function of many variables.

The loss functions of neural nets are very high dimensional and computing and storing the full Hessian matrix takes n^2 memory, which makes it intractable.

The directional second derivative tell us how well we can expect a gradient descent step to perform. We can make a second order Taylor Series approximation to the function $f(\mathbf{x})$ around the current point \mathbf{x}^0 :

$$f(\mathbf{x}) \approx f(\mathbf{x}^0) + \nabla f(\mathbf{x}^0)^T (\mathbf{x} - \mathbf{x}^0) + \frac{1}{2} (\mathbf{x} - \mathbf{x}^0)^T \mathbf{H}(\mathbf{x}^0) (\mathbf{x} - \mathbf{x}^0)$$

... where \vec{g} is the gradient and H is the Hessian at \mathbf{x}^0 .

If we use a learning rate of epsilon, then the new point \mathbf{x} , will be given by:

$$\mathbf{x}^0 - \epsilon \vec{g}$$

... substituting this into our equation, we get:

$$f(\mathbf{x}^0 - \epsilon \vec{g}) \approx f(\mathbf{x}^0) - \epsilon \vec{g}^T \vec{g} + \frac{1}{2} \epsilon^2 \vec{g}^T H \vec{g}$$

You can think of each part of the equation like this: there's the original value of the function, the expected improvement due to the slope of the function and then the curvature correction from the second derivative.

Let's call

$$\vec{g}^T H \vec{g} = \alpha$$

When α is 0 or negative, the Taylor series predicts that increasing epsilon forever will result in the function decreasing forever. This is not true, in practice Taylor expansion is not accurate for large epsilon.

When α is positive, solving for the optimal step size yields:

$$\epsilon^* = \frac{\vec{g}^T \vec{g}}{\vec{g}^T H \vec{g}}$$

Taylor Series

Taylor Series is a series expansion of a function about a point. The one-dimensional Taylor Series is an expansion of a real function $f(x)$ about a point $x = a$:

$$f(x) \approx f(a) + f'(a)(x-a) + \frac{1}{2} f''(a)(x-a)^2 + \frac{1}{6} f'''(a)(x-a)^3 + \dots$$

In order for this to be true, the function must meet certain requirements, but given that it meets them, Taylor Series will always work.

Taylor Series also works for complex variables and multivariable expansion.



If you remember your Calculus, there's the whole concept of optimization, where you can start to graph a polynomial from just the equation. You know that if the first derivative at a point is positive that it must be sloping upwards and if the second derivative is positive then it must be accelerating or curving up.

You also know that when a point has a first derivative of 0, that it's a critical point and it must be either a minima or maxima. Which brings you to the point that there can be many minima and maxima. Which are referred to as local/global minima/maxima.

And finally you also know that at a critical point you can determine if it's a minima or maxima by finding the second derivative which tells you that if it's negative then it must be a maximum.

So similarly in multivariable state, you can use the Jacobian and the Gradient as the first derivative and the Hessian as the second derivative. And roughly apply the same principles to graph in 3D space.

Remember also that you can introduce artificial constraints to your graphs as the problem suggests a real physical phenomena. After all we live in a constrained physical universe. Here we treat these constraint points as critical points. This greatly helps us in finding the global minima and maxima as it reduces our search space.

e)

Generally speaking, mathematicians like to generalise concepts and results up to the maximal point that they can, *to the limits of their usefulness*. That is, when mathematicians develop a concept, and find that one or more useful theorems apply to that concept, they will generally seek to generalise the concept and results more and more, until they get to the point where further generalisation would render the results inapplicable or no longer useful. As can be seen from your list, the exponential family has

a number of useful theorems attached to it, and it encompasses a wide class of distributions. This is sufficient to make it a worthy object of study, and a useful mathematical class in practice.

By exponential family, I mean the distributions which are given as

$$f(x|\theta)=h(x)\exp\{\eta(\theta)T(x)-B(\theta)\}$$

whose support doesn't depend on the parameter θ . Here are some advantages:

- (a) It incorporates a wide variety of distributions.
- (b) It offers a natural sufficient statistics $T(x)$ according to the Neyman-Fisher theorem.
- (c) It makes possible to provide a nice formula for the moment generating function of $T(x)$.
- (d) It makes it easy to decouple the relationship between the response and predictor from the conditional distribution of the response (via link functions).

In particular, the exponential family distributions always have conjugate priors, and the resulting posterior predictive distribution has a simple form. This makes it an extremely useful class of distributions in Bayesian statistics. Indeed, it allows you to undertake Bayesian analysis using conjugate priors at an extremely high level of generality, encompassing all the distributional families in the exponential family.

examples:

Bernoulli distribution and uniform distribution.

f)

The common feature of SVM and RVM is that they have sparse solutions, so the prediction of new data only depends on the kernel function calculated on a subset of training data. This subset is support vector for SVM and correlation vector for RVM.

The important property of SVM is that the determination of its model parameters corresponds to a convex optimization problem, so many local solutions are global optimal solutions. However, SVM does not provide posterior probability, and the important property of RVM is that the Bayesian method is introduced into RVM to provide the output of posterior probability, and can often produce more sparse solutions (the prediction speed is faster on the test set). SVM often needs to use cross validation method to determine the model complexity parameter C . For RVM, another advantage of introducing Bayesian method is to omit the step of model selection. But RVM often needs more training time because of the operation of inverse matrix.

g)

The *Kullback-Leibler divergence* (hereafter written as KL divergence) is a measure of how a probability distribution differs from another probability distribution. Classically, in Bayesian theory, there is some *true distribution* $P(X)$; we'd like to estimate with an *approximate distribution* $Q(X)$. In this context, the KL divergence measures the distance from the approximate distribution Q to the true distribution P .

Mathematically, consider two probability distributions P, Q on some space X . The Kullback-Leibler divergence from Q to P (written as $DKL(P \parallel Q)$)

$$DKL(P \parallel Q) = \mathbb{E}_{X \sim P} [\log \frac{P(X)}{Q(X)}] \quad DKL(P \parallel Q) = \mathbb{E}_{X \sim P} [\log \frac{P(X)}{Q(X)}]$$

Properties of KL Divergence

There are some immediate notes that are worth pointing out about this definition.

The KL Divergence is **not symmetric**: that is $DKL(P \parallel Q) \neq DKL(Q \parallel P)$. As a result, it is also **not a distance metric**.

The KL Divergence can take on values in $[0, \infty]$. Particularly, if P and Q are the exact same distribution ($P \text{ a.e. } = Q$), then $DKL(P \parallel Q) = 0$, and by symmetry $DKL(Q \parallel P) = 0$. In fact, with a little bit of math, a stronger statement can be proven: if $DKL(P \parallel Q) = 0$, then $P \text{ a.e. } = Q$.

In order for the KL divergence to be finite, the support of P needs to be contained in the support of Q . If a point x exists with $Q(x) = 0$ but $P(x) > 0$, then $DKL(P \parallel Q) = \infty$.

Independently of the interpretation, the KL divergence is always defined as a *specific* function of the [cross-entropy](#) (which you should be familiar with before attempting to understand the KL divergence) between two distributions (in this case, probability mass functions)

$$\begin{aligned} \text{DKL}(P \parallel Q) &= -\sum_{x \in X} p(x) \log q(x) + \sum_{x \in X} p(x) \log p(x) = H(P, Q) \\ &\quad - H(P) \end{aligned}$$

where $H(P, Q)$ is the cross-entropy of the distribution P and Q and $H(P) = H(P, P)$. The KL is not a metric, given that it does not obey the triangle inequality. In other words, in general, $\text{DKL}(P \parallel Q) \neq \text{DKL}(Q \parallel P)$.

Given that a neural network is trained to output the mean (which can be a scalar or a vector) and the variance (which can be a scalar, a vector or a matrix), why don't we use a metric like the MSE to compare means and variances? When you use the KL divergence, you don't want to compare just numbers (or matrices), but probability distributions (more precisely, probability densities or mass functions), so you will not compare just the mean and the variance of two different distributions, but you will actually compare the distributions

examples

Cross entropy and relative entropy

h)

Regularization (traditionally in the context of shrinkage) adds prior knowledge to a model; a prior, literally, is specified for the parameters. Augmentation is also a form of adding prior knowledge to a model; e.g. images are rotated, which **you** know does not change the class label. Increasing training data (as with augmentation) decreases a model's variance. Regularization also decreases a model's variance. They do so in different ways, but ultimately both decrease regularization error.

Section 5.2.2 of Goodfellow et al's [Deep Learning](#) proposes a much broader definition:

Regularization is any modification we make to a learning algorithm that is intended to reduce its generalization error but not its training error.

There is a tendency to associate regularization with shrinkage because of the term "l-p norm regularization"...perhaps "augmentation regularization" is equally valid, although it doesn't roll off the tongue. Regularization (traditionally in the context of shrinkage) adds prior knowledge to a model; a prior, literally, is specified for the parameters. Augmentation is also a form of adding prior knowledge to a model; e.g. images are rotated, which **you** know does not change the class label. Increasing training data (as with augmentation) decreases a model's variance. Regularization also decreases a model's variance. They do so in different ways, but ultimately both decrease regularization error.

Section 5.2.2 of Goodfellow et al's [Deep Learning](#) proposes a much broader definition:

Regularization is any modification we make to a learning algorithm that is intended to reduce its generalization error but not its training error.

There is a tendency to associate regularization with shrinkage because of the term "l-p norm regularization"...perhaps "augmentation regularization" is equally valid, although it doesn't roll off the tongue.

Problem 7

Supervised learning method can be divided into generation method and discriminant method. The models learned are called generative model and discriminant model respectively.

Generative model and discriminant model

definition

Discriminative vs. Generative

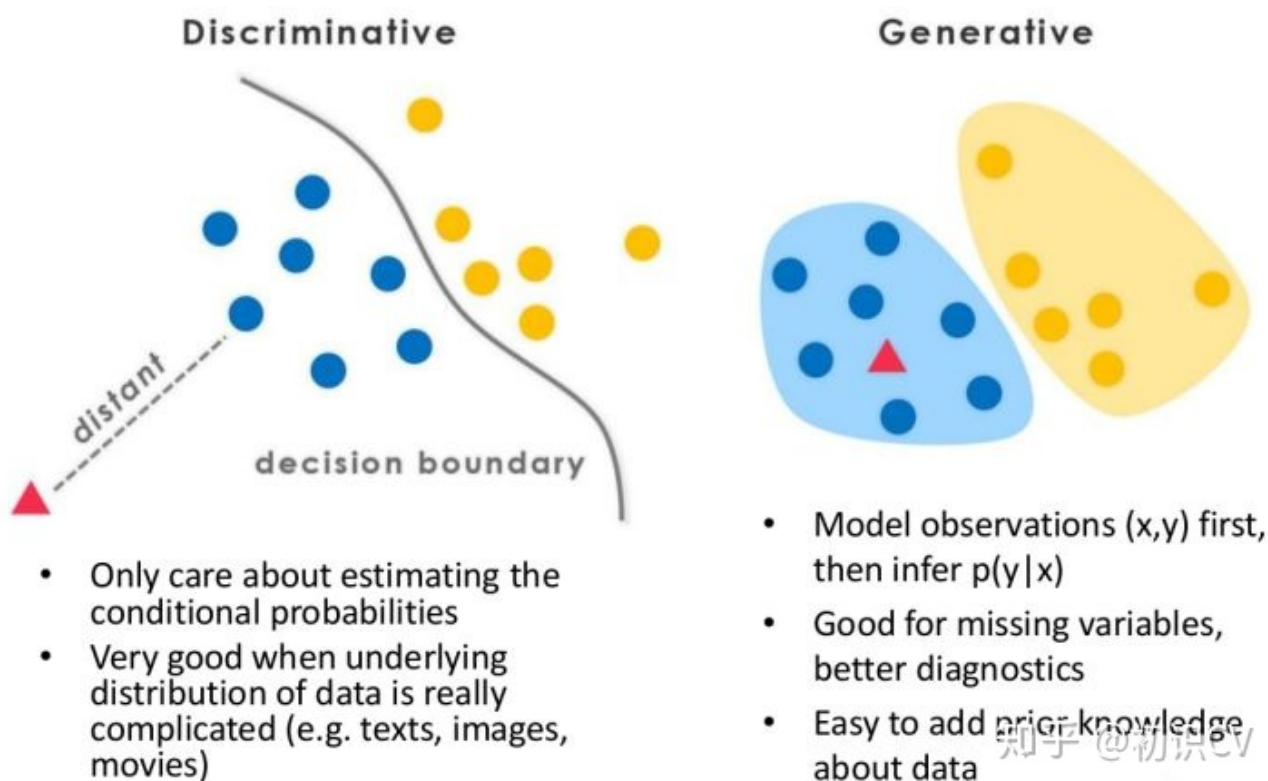


Fig. 1: left discriminant method, right generation method

The generation method (left side of Figure 1) learns the joint probability distribution $P(X,Y)$ from the data, and then obtains the conditional probability distribution $P(Y|X)$ as the prediction model, that is, the generation model: $P(Y|X)=P(X,Y)/P(X)$

It is called generation method because the model represents the generation relationship of a given input [Formula] and an output [Formula]. The typical generating models are naive Bayes method and Markov model.

For *example*: to determine whether a sheep is a goat or a sheep, the generative model is to first learn a goat model according to the characteristics of the goat, then learn a sheep model according to the characteristics of the sheep, and then extract features from the sheep, put it into the goat model to see what the probability is, and put it into the sheep model to see what the probability is.

The decision function $f(X)$ or conditional probability distribution $P(Y|X)$ is directly learned from the data as the prediction model, i.e. the discriminant model. The relationship between the discriminant methods is the given input X , what kind of output Y should be predicted. Typical discriminant models include k-nearest neighbor, perceptron, logistic regression model, maximum entropy model, support vector machine, lifting method and conditional random field.

For *example*: to determine whether a sheep is a goat or a sheep, the method of discriminant model is to learn the model from the historical data, and then to predict the probability that the sheep is a goat and a sheep by extracting the characteristics of the sheep.

difference

The generating method can restore the joint probability distribution $P(X,Y)$, but the discriminant method cannot.

The learning convergence speed of the generation method is faster, that is, when the sample size increases, the learned model can converge to the real model faster.

When there are hidden variables, the generation method can still be used to learn, but the discriminant method can not be used.

The discriminant method directly learns the conditional probability $P(Y|X)$ or decision function $f(X)$, and when it directly faces the prediction, it often has higher learning accuracy.

Because the discriminant method can directly learn the conditional probability $P(Y|X)$ or the decision function $f(X)$, it can abstract the data in various degrees, define the characteristics and use the features, so the discrimination method can simplify the learning problem.