

# CS 181 Spring 2021 Section 4 Notes Solution

## 1 Bayesian Regression

### 1.1 Motivations

The Bayesian frame of reference helps us answer three types of questions regarding data  $X$  and labels  $Y$ :

- The **posterior** over models:  $p(\theta|X, Y) \propto p(Y|X, \theta)P(\theta|X)$   
This tells us how likely different values of the model  $\theta$  are after updating the prior distribution with the observed data.
- The **posterior predictive** for new data:  $p(y^*|x^*, X, Y) = \int p(y^*|x^*, \theta)p(\theta|X, Y)d\theta$   
This tells us how to predict the label of a new data point according to the posterior over models obtained by updating the prior with the observed data.
- The **marginal likelihood** of data:  $p(Y|X) = \int p(Y|X, \theta)p(\theta)d\theta$   
This tells us how likely the data is, marginalizing over possible models  $\theta$ . In contrast to the likelihoods we've seen before which are computed given a specific setting of weights  $\theta$ , the marginal likelihood accounts for a distribution over models  $\theta$ . The marginal likelihood allows us to compare different priors or even different model classes in terms of how well they fit the data (this is model selection).

### 1.2 Conjugate Pairs

Focusing in on the concept of posteriors, note that  $p(\theta)$  represents our prior beliefs regarding the optimal model  $\theta$ , and  $p(\theta|X, Y)$  represents our updated beliefs after observing some data  $X, Y$ . There is no guarantee that  $p(\theta)$  and  $p(\theta|X, Y)$  share a nice, clean relationship, but sometimes they do, thanks to some special structure in the distribution of labels  $p(Y|X, \theta)$ . When a certain model can represent  $p(\theta)$  and  $p(\theta|X, Y)$  with the same distribution, we call this a **conjugacy pair**.

#### 1.2.1 Beta-Binomial Conjugacy

One example, seen in lecture, occurs when we model the prior distribution  $p(\theta)$  using a Beta distribution, and we model the data  $p(Y|X, \theta)$  according to a Bernoulli (or Binomial) distribution. In this case, the posterior distribution  $p(\theta|X, Y)$  follows a Beta distribution, just like the prior. Note that the parameter values may have changed, but the class of the distribution is the same!

### 1.3 Bayesian Linear Regression

Let  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ ,  $\mathbf{x}_i \in \mathbb{R}^m$ ,  $y_i \in \mathbb{R}$ . Consider the generative model:

$$y_i \sim \mathcal{N}(\mathbf{w}^\top \mathbf{x}_i, \beta^{-1}) \quad (1)$$

The likelihood of the data has the form:

$$p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = \mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \beta^{-1}\mathbf{I}) \quad (2)$$

Put a conjugate prior on the weights (assume covariance  $\mathbf{S}_0$  known):

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mu_0, \mathbf{S}_0) \quad (3)$$

We want a posterior distribution on  $\mathbf{w}$ . Using Bayes' Theorem:

$$p(\mathbf{w}|\mathcal{D}) \propto p(\mathcal{D}|\mathbf{w})p(\mathbf{w}) \quad (4)$$

It turns out that our posterior after  $N$  examples is also Gaussian:

$$p(\mathbf{w}|\mathcal{D}) = \mathcal{N}(\mathbf{w}|\mu_N, \mathbf{S}_N) \quad (5)$$

where

$$\mathbf{S}_N = \left( \mathbf{S}_0^{-1} + \beta \mathbf{X}^\top \mathbf{X} \right)^{-1} \quad (6)$$

$$\mu_N = \mathbf{S}_N (\mathbf{S}_0^{-1} \mu_0 + \beta \mathbf{X}^\top \mathbf{y}) \quad (7)$$

This tells us that Gaussian-Gaussian is yet another example of a conjugacy pair.

## 1.4 Concept Question

Why do we care about conjugacy pairs? What makes them convenient to study?

When the prior distribution and posterior distribution are different, it becomes very difficult to model updates in our beliefs as we see new data. Conjugacy pairs are remarkably helpful because they allow us to safely assume that the class of distributions doesn't change as we see new data. We can instead focus on adjusting the parameters of such a distribution and thus achieve much deeper insights regarding a model.

## 1.5 Posterior Predictive Distributions

We have seen how to obtain a posterior distribution over  $\mathbf{w}$ . But, given this posterior and a new data point  $\mathbf{x}^*$ , how do we actually make a prediction  $y^*$ ? How do we deal with *uncertainty* about  $\mathbf{w}$ ? We can expand the predictive distribution over  $y^*$  given  $\mathbf{x}^*$  using the Law of Total Probability:

$$p(y^*|\mathbf{x}^*, \mathcal{D}) = \int_{\mathbf{w}} p(y^*|\mathbf{x}^*, \mathbf{w})p(\mathbf{w}|\mathcal{D})d\mathbf{w} \quad (8)$$

$$= \int_{\mathbf{w}} \mathcal{N}(y^*|\mathbf{w}^\top \mathbf{x}^*, \beta^{-1})\mathcal{N}(\mathbf{w}|\mu_N, \mathbf{S}_N)d\mathbf{w} \quad (9)$$

This is the **posterior predictive** distribution over  $y^*$ . This can be interpreted as a weighted average of many predictors, one for each choice of  $\mathbf{w}$ , weighted by how likely  $\mathbf{w}$  is according to the posterior. Since each of the terms on the right hand side follows a normal distribution, we can use some math to find that:

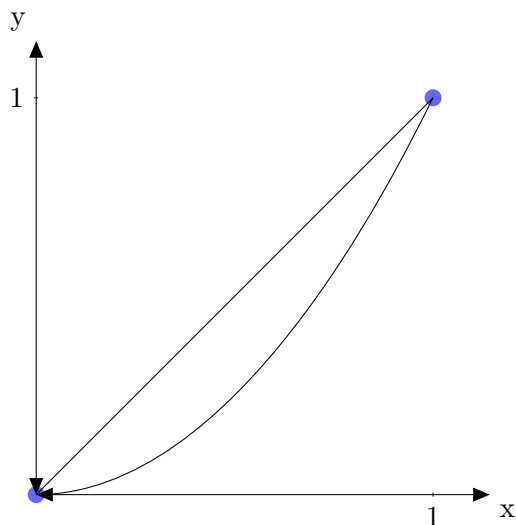
$$p(y^*|\mathbf{x}^*, \mathcal{D}) = \mathcal{N}(y^*|\mu_N^\top \mathbf{x}^*, \mathbf{x}^{*\top} \mathbf{S}_N \mathbf{x}^* + \beta^{-1}) \quad (10)$$

## 1.6 Exercise: A Simple Bayesian Model

Say you are tasked with fitting a parabolic regression  $\hat{y} = a_0 + a_1x + a_2x^2$  on data of the form  $(x, y)$  for feature  $x$  and label  $y$ . That is, you want to find the best possible  $a_0, a_1, a_2$  to fit the data you are given.

1. Before seeing any data, what are reasonable prior distributions for the parameters  $a_0, a_1, a_2$ ?

You are now presented with two data points,  $(0,0)$  and  $(1,1)$ . You are told that the data was generated from some  $y = f(x)$  with negligible error, that is  $\epsilon \sim \mathcal{N}(0, \sigma^2)$  where  $\sigma^2 \approx 0$ .



2. Given  $\sigma^2 \approx 0$ , does the actual function  $y = f(x)$  go through the points  $(0,0)$  and  $(1,1)$ ? Using this information, what relationships must exist between  $a_0, a_1, a_2$ ? Update your prior distributions on  $a_0, a_1, a_2$  to become posterior distributions.
3. You are now told that  $a_1 \sim \text{Unif}(-1, 1)$ . Given this and the data you've observed, if a new data point  $x^* = 1/2$  is presented, what is the posterior predictive on  $y^*$ ?
4. What would the posterior predictive on  $y^*$  have been for a linear regression  $\hat{y} = a_0 + a_1x$  instead? Compare the strengths of the two models' predictions for  $y^*$ . How does this relate to the expressivity of the different model classes?

---

**Solution**

---

1. There are many possible answers, but a natural one is that  $a_0, a_1, a_2$  are uniformly distributed across the real numbers, since we don't know anything about them yet. One could also argue that the weights are more likely to be close to zero than not, so perhaps a Gaussian distribution centered at zero makes sense as well.
2. Yes, with zero error, we must have  $f(0) = 0$  and  $f(1) = 1$ . Thus, we can calculate the posterior relationships  $\hat{y}(0) = a_0 = 0$  and  $\hat{y}(1) = a_0 + a_1 + a_2 = 1 \Rightarrow a_1 + a_2 = 1$ . Our posterior distributions are thus  $a_0 = 0$ ,  $a_1$  is still uniform over real numbers, and  $a_2 = 1 - a_1$ .
3. We see that  $y^*(1/2) = a_0 + a_1(1/2) + a_2(1/2)^2 = \frac{1}{2}a_1 + \frac{1}{4}(1 - a_1) = \frac{1}{4} + \frac{1}{4}a_1 \sim \text{Unif}(0, 1/2)$ . Thus, the posterior predictive for  $y^*$  is a uniform distribution with support  $[0, 1/2]$ .
4. For a linear regression, we must have  $\hat{y}(0) = a_0 = 0$  and  $\hat{y}(1) = a_0 + a_1 = 1 \Rightarrow a_1 = 1$ . Thus,  $y^*(1/2) = a_0 + a_1(1/2) = 1/2$ , and the posterior predictive for  $y^*$  is a dirac delta distribution centered at  $1/2$ , which means  $y^* = 1/2$  with probability 1. If we compare the two models' posterior distributions, we see that the linear regression has a much spikier distribution, all concentrated at  $y^* = 1/2$ , which is stronger than the parabolic regression's distribution spread evenly across  $[0, 1/2]$ . Thus, we can reason about the models in the following way: a less expressive (linear) model has a smaller range of possible predictives but conveys a stronger conviction for the predictions it does issue, while a more expressive (parabolic) model has a larger range of possible predictives but conversely a weaker convictions for all the predictions it can issue.

---

**End Solution**

---

## 1.7 Exercise: Posterior Distribution By Completing the Square (Bishop 3.7)

We know from (3.10) in Bishop that the likelihood can be written as

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}, \mathbf{w}) &= \prod_{i=1}^n \mathcal{N}(y_i | \mathbf{w}^\top \mathbf{x}_i, \beta^{-1}) \\ &\propto \exp\left(-\frac{\beta}{2}(\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w})\right) \end{aligned}$$

where precision  $\beta = \frac{1}{\sigma^2}$  and in the second line above we have ignored the Gaussian normalization constants. By completing the square, show that with a prior distribution on  $\mathbf{w}$  given by  $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mu_0, \mathbf{S}_0)$  where  $\mathbf{S}_0$  is the covariance matrix, the posterior distribution  $p(\mathbf{w}|\mathcal{D})$  is given by

$$p(\mathbf{w}|\mathcal{D}) = \mathcal{N}(\mathbf{w}|\mu_N, \mathbf{S}_N)$$

where

$$\begin{aligned} \mu_N &= \mathbf{S}_N(\mathbf{S}_0^{-1}\mu_0 + \beta\mathbf{X}^\top \mathbf{y}) \\ \mathbf{S}_N &= (\mathbf{S}_0^{-1} + \beta\mathbf{X}^\top \mathbf{X})^{-1} \end{aligned}$$

Here's the first step. Take  $\ln[(\text{likelihood})(\text{prior})]$  and collect normalization terms that don't depend on  $\mathbf{w}$ :

$$\begin{aligned} \ln p(\mathbf{w}|\mathcal{D}) &\propto \ln p(\mathbf{y}|\mathbf{X}, \mathbf{w}) + \ln p(\mathbf{w}) \\ &= \text{const} - \frac{\beta}{2}(\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w}) - \frac{1}{2}(\mathbf{w} - \mu_0)^\top \mathbf{S}_0^{-1}(\mathbf{w} - \mu_0) \end{aligned}$$

Hint: Remember, you already know what the posterior should look like. Once you simplify your expression enough, try foiling the posterior in terms of  $\mu_N$  and  $\mathbf{S}_N^{-1}$  and see if you can see the relationship between your expression and this posterior.

---

**Solution**

---

As the problem statement suggestions, the first step is to take  $\ln[(\text{likelihood})(\text{prior})]$  and collect normalization terms that don't depend on  $\mathbf{w}$ :

$$\ln p(\mathbf{w}|D) \propto \ln p(\mathbf{y}|\mathbf{X}, \mathbf{w}) + \ln p(\mathbf{w}) \quad (11)$$

$$= \text{const} - \frac{\beta}{2}(\mathbf{y} - \mathbf{X}\mathbf{w})^\top(\mathbf{y} - \mathbf{X}\mathbf{w}) - \frac{1}{2}(\mathbf{w} - \mu_0)^\top \mathbf{S}_0^{-1}(\mathbf{w} - \mu_0) \quad (12)$$

Expanding, we have:

$$\text{const} - \frac{1}{2} \left( \beta \mathbf{y}^\top \mathbf{y} - \beta \mathbf{y}^\top \mathbf{X}\mathbf{w} - \beta \mathbf{w}^\top \mathbf{X}^\top \mathbf{y} + \beta \mathbf{w}^\top \mathbf{X}^\top \mathbf{X}\mathbf{w} \right) \quad (13)$$

$$+ \mathbf{w}^\top \mathbf{S}_0^{-1} \mathbf{w} - \mathbf{w}^\top \mathbf{S}_0^{-1} \mu_0 - \mu_0^\top \mathbf{S}_0^{-1} \mathbf{w} + \mu_0^\top \mathbf{S}_0^{-1} \mu_0 \quad (14)$$

Important: Remember that terms like  $\mathbf{y}^\top \mathbf{X}\mathbf{w}$  are the same scalar as  $\mathbf{w}^\top \mathbf{X}^\top \mathbf{y}$ . Collecting together quadratic and linear terms, factoring the  $\mathbf{w}$ s out, and moving terms that don't depend on  $\mathbf{w}$  into the constant, we have

$$\text{const} - \frac{1}{2} \left( \mathbf{w}^\top (\mathbf{S}_0^{-1} + \beta \mathbf{X}^\top \mathbf{X}) \mathbf{w} - 2\mathbf{w}^\top (\mathbf{S}_0^{-1} \mu_0 + \beta \mathbf{X}^\top \mathbf{y}) \right) \quad (15)$$

Put aside what we have done so far. Recall that our target looks like:

$$-\frac{1}{2} \left( (\mathbf{w} - \mu_N)^\top \mathbf{S}_N^{-1} (\mathbf{w} - \mu_N) \right), \quad (16)$$

When expanded, this looks like

$$-\frac{1}{2} \left( \mathbf{w}^\top \mathbf{S}_N^{-1} \mathbf{w} - \mu_N^\top \mathbf{S}_N^{-1} \mathbf{w} - \mathbf{w}^\top \mathbf{S}_N^{-1} \mu_N + \mu_N^\top \mathbf{S}_N^{-1} \mu_N \right) \quad (17)$$

Drop the term that doesn't have  $\mathbf{w}$  and combine the two middle terms

$$-\frac{1}{2} \left( \mathbf{w}^\top \mathbf{S}_N^{-1} \mathbf{w} - 2\mathbf{w}^\top \mathbf{S}_N^{-1} \mu_N \right) \quad (18)$$

This looks like what we ended up with in (15) where

$$\mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta \mathbf{X}^\top \mathbf{X}$$

and

$$\mathbf{S}_N^{-1} \mu_N = \mathbf{S}_0^{-1} \mu_0 + \beta \mathbf{X}^\top \mathbf{y},$$

as required. This has the desired form of an (unnormalized) Gaussian.

---

**End Solution**

---

## 2 Neural Networks

### 2.1 Takeaways

Recall that in the case of binary classification, we can think about a neural network as being equivalent to logistic regression with parameterized, adaptive basis functions. Adaptive means that you don't need to specify a feature basis. We can train a matrix that linearly transforms the data, run the resulting vector through an element-wise non-linearity, potentially repeat this process, and then finally run logistic regression on the resulting vector, where the logistic regression itself has weights that need to be trained.

### 2.2 Activation Functions

There a wide range of possible non-linear activation functions to choose from when designing neural networks. Among these, the most popular are ReLU, which you will encounter below, and the tanh activation function, which is exactly what it sounds like:  $\tanh(z)$ . There are many others, including sigmoid and softmax for the final output layers, that researchers decide between when constructing their models.

**Note:** Activation functions are vital to constructing neural networks because they introduce non-linear relationships between the inputs and outputs of various layers. If there were no special activation functions, every step of a feed-forward neural network would just reduce to matrix multiplication (try this out yourself!), and all the relationships from the initial inputs to the final outputs would just be linear combinations of variables. Only linear relationships would come out of the model!

### 2.3 Concept Question

What is the difference between the activation functions described here (ReLU, tanh, sigmoid, softmax) and the usage of sigmoid and softmax earlier in logistic regression?

The purpose of these activation functions is to be non-linear, while earlier the sigmoid and softmax functions were used to transform outputs into probabilities. Thus, sigmoid and softmax both necessarily have range  $[0, 1]$ , while for instance tanh has range  $[-1, 1]$ . For neural networks, having a range centered around 0, like tanh, is an advantage because it means there is no systemic bias being introduced to the output values, whereas functions like ReLU, sigmoid and softmax result in all non-negative output values which can limit the expressiveness of the neural network.

## 2.4 Exercise: A Simple NN Classifier

Let's think about a neural network binary classifier with  $\mathbf{x} \in \mathbb{R}^2$  ( $D = 2$ ) and with a single two-dimensional hidden layer ( $M = 2$ ) followed by the one-dimensional output layer.

For the non-linear activation function, we use the *ReLU* function defined by the following:

$$\text{ReLU}(z) = \begin{cases} z & \text{if } z > 0 \\ 0 & \text{otherwise} \end{cases} \quad (19)$$

Consider a function  $h$  defined by the following:

$$h(\mathbf{x}) = \sigma(\mathbf{w}^\top \phi(\mathbf{x}) + w_0) \quad (20)$$

We can decompose  $h$  in terms of our neural network classifier's weights where the first (hidden) layer composes the "adaptive basis" and the second (output) layer composes the logistic regression:

$$h(\mathbf{x}) = \sigma(\mathbf{w}^\top \mathbf{ReLU}(\mathbf{W}^{(hid)\top} \mathbf{x} + \mathbf{w}_0^{(hid)})) + w_0 \quad (21)$$

$$= \sigma\left(\sum_{m=1}^M w_{1m}^{(2)} \mathbf{ReLU}\left(\sum_{d=1}^D w_{md}^{(1)} x_d + w_{m0}^{(1)}\right) + w_{10}^{(2)}\right), \quad (22)$$

where  $\sigma(z)$  is the sigmoid function and in the last equation above, we explicitly separately notate the weights of the output layer

$$\mathbf{w} = \begin{pmatrix} w_{11}^{(2)} \\ w_{12}^{(2)} \end{pmatrix} \in \mathbb{R}^2, w_0 = w_{10}^{(2)} \in \mathbb{R}$$

and of the hidden layer

$$\mathbf{W}^{(hid)} = \begin{pmatrix} w_{11}^{(1)} & w_{21}^{(1)} \\ w_{12}^{(1)} & w_{22}^{(1)} \end{pmatrix} \in \mathbb{R}^{2 \times 2}, \mathbf{w}_0^{(hid)} = \begin{pmatrix} w_{10}^{(1)} \\ w_{20}^{(1)} \end{pmatrix} \in \mathbb{R}^{2 \times 2}.$$

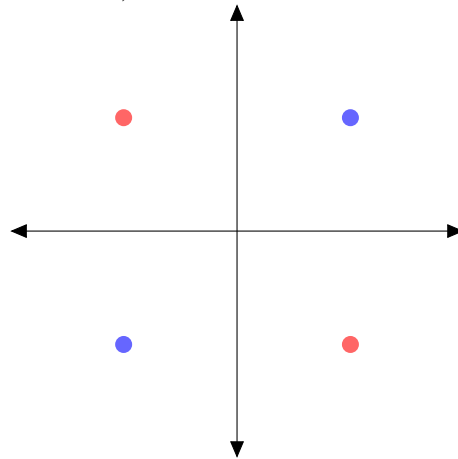
Then, classifications are made according to  $\mathbb{I}_{h(\mathbf{x}) > .5}$ . Suppose we want to fit the following data:

$$\mathbf{x}_1 = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad y_1 = 1, \quad \mathbf{x}_2 = \begin{pmatrix} 1 \\ -1 \end{pmatrix} \quad y_2 = 0$$

$$\mathbf{x}_3 = \begin{pmatrix} -1 \\ 1 \end{pmatrix} \quad y_3 = 0, \quad \mathbf{x}_4 = \begin{pmatrix} -1 \\ -1 \end{pmatrix} \quad y_4 = 1$$



This looks as follows (blue = 1, red = 0):



Why can't we solve this problem with a linear classifier? What values of parameters  $\mathbf{W}^{(hid)}$ ,  $\mathbf{w}_0^{(hid)}$ ,  $\mathbf{w}$ , and  $w_0$  will allow the neural network to solve the problem? Show that your choice of parameters allows the network to correctly classify the data.

**Note:** You do not need to learn to find these weights systematically by hand. This is not very possible to do manually in general, but you might see something in this low-dimensional case.

*Hint:* Think carefully about what the ReLU activation function can do for us. What does it do to various kinds of vectors? Think geometrically.

---

### Solution

---

We can solve the problem with

$$\mathbf{W}^{(hid)} = \begin{pmatrix} 1 & -1 \\ 1 & -1 \end{pmatrix}, \quad \mathbf{w}_0^{(hid)} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$\mathbf{w} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \quad w_0 = -1$$

How does this solve the problem? Note that

$$h(\mathbf{x}_1) = (1 \ 1) \mathbf{ReLU} \begin{pmatrix} 2 \\ -2 \end{pmatrix} - 1 = 1 \Rightarrow \mathbb{I}_{h(\mathbf{x}_1) > .5} = 1$$

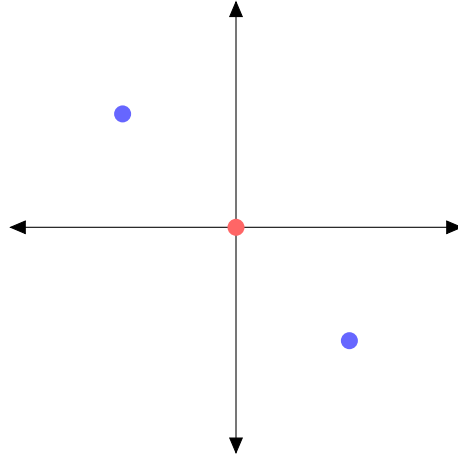
$$h(\mathbf{x}_2) = (1 \ 1) \mathbf{ReLU} \begin{pmatrix} 0 \\ 0 \end{pmatrix} - 1 = -1 \Rightarrow \mathbb{I}_{h(\mathbf{x}_2) > .5} = 0$$

$$h(\mathbf{x}_3) = (1 \ 1) \mathbf{ReLU} \begin{pmatrix} 0 \\ 0 \end{pmatrix} - 1 = -1 \Rightarrow \mathbb{I}_{h(\mathbf{x}_3) > .5} = 0$$

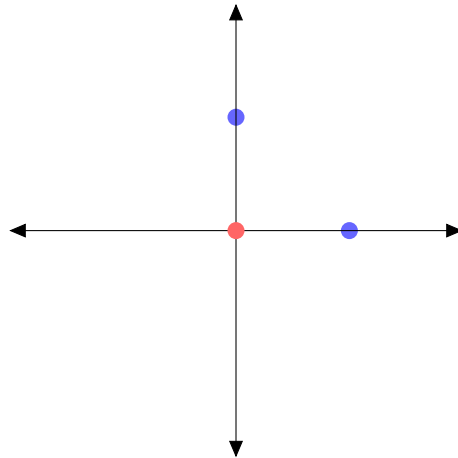
$$h(\mathbf{x}_4) = (1 \ 1) \mathbf{ReLU} \begin{pmatrix} -2 \\ 2 \end{pmatrix} - 1 = 1 \Rightarrow \mathbb{I}_{h(\mathbf{x}_4) > .5} = 1$$

This works because we have a ReLU! Having a nonlinearity in the activation function allows us to find a linear classifier for these examples in the new basis.

What did we do? With  $\mathbf{W}^{(hid)}$ , we map our data linearly so it looks as follows:



This is great, but still not linearly separable. Applying the ReLU, however, maps this picture to the following one:



Now our data is linearly separable and we can classify the points correctly. Note that the parameters we chose here are not unique (far from it). The important thing is that we transformed our input space well enough to apply a linear classifier.

---

**End Solution**

---