

# Linear Models

Jonathan May

August 30, 2023(Prepared for Fall 2023)

**Discussion Question: How is your understanding of your own ability to learn like the machine learning paradigm and how is it different?**

We're going to start at section 9; the first sections should be review and I will refer to their contents. If you don't already know the contents here, we should talk!

## 1 Definitions

- Experiment: Some action that takes place in the world
- Outcomes = Sample Space =  $\Omega$  = the universe, every *basic outcome* that could happen
- Event =  $A \subseteq \Omega$  = something that happened (could be more than one basic outcome)
- Probability Distribution =  $P : \Omega \rightarrow [0, 1]$ ,  $\sum_{x \in \Omega} P(x) = 1$ , i.e. values sum to 1 and no value is negative

## 2 Example

- Experiment = “toss a coin three times”
- $\Omega = \{HHH, HHT, HTT, HTH, THH, THT, TTT, TTH\}$
- Event  $A$  = “exactly two heads” =  $\{HHT, HTH, THH\}$
- Event  $B$  = “first one was heads” =  $\{HHH, HHT, HTT, HTH\}$
- Distribution: assign a number between 0 and 1 (‘probability’) to each basic outcome<sup>1</sup>; sum of all such numbers = 1
- Uniform Distribution: define  $P(x) = c \forall x \in \Omega$ ...in this case?
- Probability of an event = sum of the probability of its basic outcomes
- So,  $P(A) = ?$  and  $P(B) = ?$

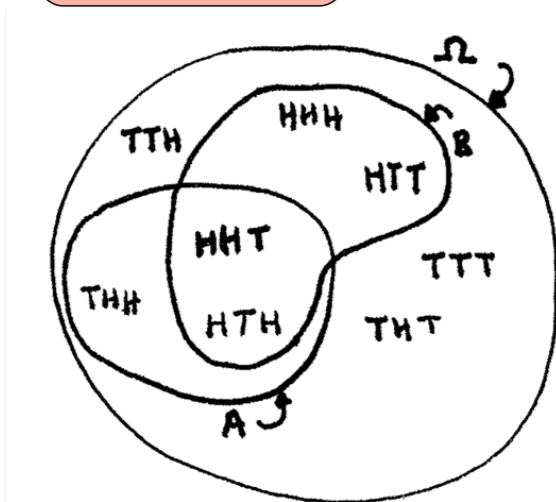
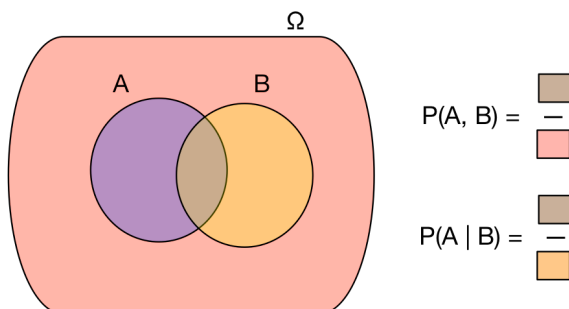
---

<sup>1</sup>actually to each event in a partition but we'll get back to that in a minute

### 3 Joint and Conditional Probability

$P(A, B) = P(A \cap B)$  = ‘joint probability of A and B’, i.e. probability of the event formed by the intersection operation (can think of it as probability of ‘the joint event’)

$P(A|B) = \frac{P(A \cap B)}{P(B)}$  = ‘conditional probability of A given B’, i.e. the joint event above, assuming that B is  $\Omega$



So  $P(A) = 3/8$ , and  $P(B) = 1/2$

$P(A, B) = P(A) + P(B)$ ? (No. What is it?)

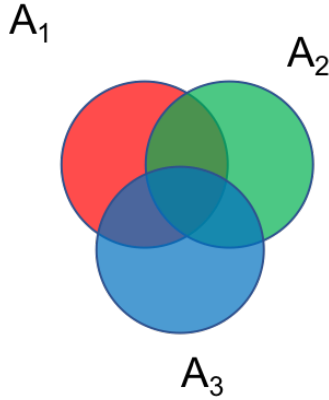
$P(B|A)$  = “if you’ve got two heads what’s the chance your first was heads” = ?

$P(A|B)$  = “if your first is heads what’s the chance you’ve got two” = ?

### 4 Chain Rule of Probability

(Not to be confused with the chain rule of Calculus)

Since  $P(A|B) = \frac{P(A, B)}{P(B)}$  (by definition), we can rewrite terms to get  $P(A, B) = P(A|B)P(B)$ .



Now consider three events,  $A_1, A_2, A_3$ . How can we define  $P(A_1, A_2, A_3) = P(A_1 \cap A_2 \cap A_3)$  in terms of conditional probabilities?

Recall that an event is just a set of basic outcomes. So let's define a new event

$$A_{23} = A_2 \cap A_3$$

Then we would write

$$P(A_1, A_{23}) = P(A_1|A_{23})P(A_{23})$$

Now, substitute back in the joint event that  $A_{23}$  represents

$$P(A_1, A_2, A_3) = P(A_1|A_2, A_3)P(A_2, A_3)$$

Now substitute the definition of joint probabilities in terms of conditional probabilities again

$$P(A_1, A_2, A_3) = P(A_1|A_2, A_3)P(A_2|A_3)P(A_3)$$

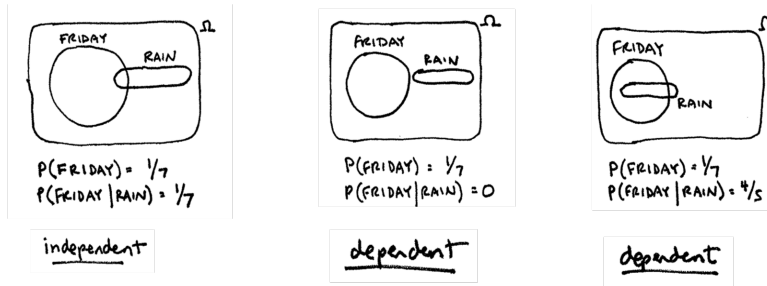
Of course  $P(A_1, A_2, A_3) = P(A_3, A_2, A_1)$  (set intersection is commutative) so you could write this instead as  $P(A_3|A_2, A_1)P(A_2|A_1)P(A_1)$

The general chain rule for probabilities is:

$$P(A_1, \dots, A_N) = P(A_1|A_2, \dots, A_N) \times \dots \times P(A_{N-1}|A_N) \times P(A_N)$$

## 5 Independence

$A$  and  $B$  are *independent* if the occurrence of one does not affect the occurrence of the other, i.e. if  $P(A|B) = P(A)$ . Corollary,  $P(B|A) = P(B)$



## 6 Bayes' Rule/Theorem/Law

$P(A|B) = \frac{P(A,B)}{P(B)}$ , by definition. Thus,  $P(A,B) = P(A|B)P(B)$ .

Because intersection is commutative (see above),  $P(A,B) = P(B|A)P(A)$ . This also explains the corollary noted in Section 5. This leads to Bayes' Rule/Theorem/Law:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

This can be very helpful when you have information about one conditional direction but you want info about the other direction.

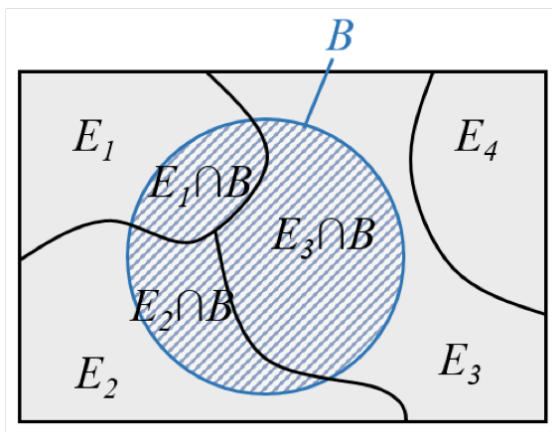
## 7 Law of Total Probability

We say events  $E_1, \dots, E_n$  *partition*  $\Omega$  if:

$$\forall i, j \in [1, n], E_i \cap E_j = \emptyset$$

and

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



The Law of Total Probability says, given partitioning events  $E_1 \dots E_n$  and event  $B$ :

$$P(B) = \sum_{i=1}^n P(B, E_i)$$

## 8 Example

Some people can read minds, but not many:  $P(MR) = 1/100,000 = .00001$ .

There is a test to read minds; if you are a mind reader I can detect this very well:  $P(T|MR) = 0.95$  and if you're not I can detect this even better:  $P(\neg T|\neg MR) = 0.995$ . Note:  $\{T, \neg T\}$  partition the event space, as do  $\{MR, \neg MR\}$ .

If Jill gets a positive result on the test, how likely is it she is the mind reader?

i.e.,  $P(MR|T) = ?$

By Bayes' Law,  $P(MR|T) = \frac{P(T|MR)P(MR)}{P(T)}$ .

We need to get  $P(T)$ . By law of total probability,  $P(T) = P(T, MR) + P(T, \neg MR)$ . By definition of conditional probability,  $P(T, MR) = P(T|MR)P(MR) = 0.95 \times .00001$ ;  $P(T, \neg MR) = P(T|\neg MR)P(\neg MR) = 0.005 \times .99999$ .  $P(T) = .00500945$  and  $P(MR|T) \approx .002$

## 9 A classifier model

Here's a simple framework for this problem:

```
def evaluate(sentence, option, model):
    # fill me in
    pass

def classify(sentence, options, model):
    scores = {}
    for option in options:
        scores[option] = evaluate(sentence, option, model)
    # key functions as argmax
    return max(scores, key=scores.get)
```

Let's consider methods for the 'evaluate' function:

## 10 Rule-Based (top-down?) Model

Positive reviews should have positive words, negative reviews should have negative words. Thankfully, people have compiled such *sentiment lexicons* for English. See <http://www.enchantedlearning.com/wordlist>.

Positive words example: {absolutely, adorable, bountiful, bounty, cheery}

Negative words example: {angry, abysmal, bemoan, callous}

Let's put the intuition and data together:

```

# probably shouldn't structure your code this way...
model = {}
model['good'] = set(['yay', 'love', ...])
model['bad'] = set(['terrible', 'boo', ...])

def evaluate(sentence, option, model):
    score = 0
    for word in sentence.split():
        if word in model[option]:
            score+=1
    return score

```

## 11 Empirical Model

Our intuitions about word sentiment aren't perfect and neither are those of the people who made the word list. But we do have many examples of reviews and their sentiments. So we can make our own list of good and bad words. Instead of hard-coding the model as I did above, we can create a *training* function that takes in sentences *with their labels* and then returns a model:

```

from collections import Counter, defaultdict
def train(labeled_sentences):
    scores = defaultdict(lambda: Counter()) # doubly nested structure
    for sentence, label in labeled_sentences:
        for word in sentence.split():
            scores[word][label]+=1
    model = defaultdict(lambda: set())
    for word, table in scores.items():
        # the most frequent label associated with the word
        # key functions as argmax
        label = max(table, key=table.get)
        model[label].add(word)
    return model

```

We now have our first *supervised* model. We should back up and consider what we are actually trying to model from a probabilistic view. For one thing let's consider the actual probability of each label, rather than our ad-hoc adding method above.

Our classifier should choose  $\operatorname{argmax}_y P(y|s)$  for sentence  $s$  where  $y$  is one of a fixed set of labels. But we didn't consider  $s$  as some monotone thing; we considered the occurrence of each word as an event.

So we'll make an assumption called the *bag of words assumption* which is that for sentence  $s = w_1 w_2 \dots w_n$ , we say  $P(y|s) = P(y|w_1, w_2, \dots, w_n)$ . Note this now doesn't depend on the order of the words.<sup>2</sup>

---

<sup>2</sup>This assumption isn't really part of Naive Bayes, it's an assumption about the feature set being used.



Figure 7.1 Intuition of the multinomial naive Bayes classifier applied to a movie review. The position of the words is ignored (the *bag of words* assumption) and we make use of the frequency of each word.

Figure from J&M 3rd ed. draft, sec 7.1

But this model will only be valid if we can calculate probabilities of a label for the exact multiset of each sentence's words. We need another assumption, the *Naive Bayes assumption* which is that the probability of a label given a word is conditionally independent of the other words.

Note from Bayes' rule:

$$P(y|w_1, w_2, \dots, w_n) = \frac{P(w_1, w_2, \dots, w_n|y)P(y)}{P(w_1, w_2, \dots, w_n)}$$

and since the word sequence itself is constant, we can say

$$\operatorname{argmax}_y P(y|w_1, w_2, \dots, w_n) = \operatorname{argmax}_y P(w_1, w_2, \dots, w_n|y)P(y)$$

(This means that we're optimizing the joint probability  $P(s, y)$ , but it's monotonic with  $P(s|y)$ .)

We can use the chain rule to break up the first term on the right:

$$P(w_1, w_2, \dots, w_n|y) = P(w_1|y, w_2, \dots, w_n)P(w_2|y, w_3, \dots, w_n) \dots P(w_n|y)$$

Then we apply the Naive Bayes assumption:

$$P(w_1|y, w_2, \dots, w_n) \approx P(w_1|y)$$

and so on for all the other terms. That gives us

$$P(w_1, w_2, \dots, w_n|y)P(y) \approx P(w_1|y)P(w_2|y), \dots, P(w_n|y)P(y)$$

Is this a good model? George Box, statistician: "All models are wrong, but some models are useful."

What are problems with the bag of words assumption?

What are problems with the naive bayes assumption?

Does it work? Yes, for many tasks actually. And it's very simple so it's usually worth trying.

Here's a new trainer:

It's probably better to say that for  $x, y$ , we calculate  $f(x, y) = f_1, \dots, f_n$  and then proceed from there. The above is a bit of a simplification. The reading avoids this simplification.

```

from collections import Counter, defaultdict
wprobdenom = "__ALL__" # assume this token doesn't appear
def train(labeled_sentences):
    wscores = defaultdict(lambda: Counter()) # doubly nested structure
    cscores = Counter()
    for sentence, label in labeled_sentences:
        cscores[label] += 1
        for word in sentence.split():
            wscores[label][word] += 1
            wscores[label][wprobdenom] += 1
    model = {'cprobs': {}, 'wprobs': {}}
    for label in cscores.keys():
        model['cprobs'][label] = cscores[label] / len(labeled_sentences)
        wprob = {}
        for word, score in wscores[label].items():
            wprob[word] = score / wscores[label][wprobdenom]
        model['wprobs'][label] = wprob
    return model

```

And a new, more appropriate classifier

```

def evaluate(sentence, option, model):
    score = model['cprobs'][option]
    for word in sentence.split():
        score *= model['wprobs'][option][word]
    return score

```

## 11.1 Practicalities: smoothing

`score *= model['wprobs'][option][word]` is going to be problematic if we have never seen a word with a particular class. Solution: smoothing!

Laplace (add-1) smoothing: assume you've seen every word (even words you haven't seen before) with every class!

Before (assume 10k words in training set of negative items):

$P(\text{amazing}|\text{negative}) = 0/10,000 = 0$  (seen the word but not with class 'negative')

$P(\text{blargh}|\text{negative}) = 0/10,000 = 0$  (never seen the word)

Introduce a new term, 'OOV', and if you haven't seen your test word during training, pretend your word is 'OOV'. Then, since you add 1 for each vocabulary word and the OOV with each class, if your vocabulary size was 500, you now get:

$P(\text{amazing}|\text{negative}) = 1/10,501$

and

$P(\text{blargh}|\text{negative}) = 1/10,501$

You may want to actually *introduce* some OOV into your training set, by replacing words that appear fewer than some  $k$  times with OOV. This is so that your model can learn how to behave with OOVs.



If you use a subword tokenizer like BPE, you can reduce the reliance on OOV and smoothing, but you can't quite eliminate either one. Why not?

## 11.2 Practicalities: underflow

Recall:

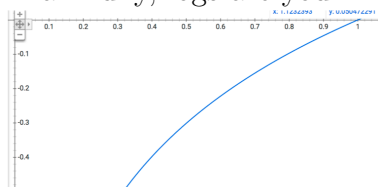
```
for word in sentence.split():
    score += model['wprobs'][option][word]
```

Sentence may be long! Probabilities may be small! It's very easy to run into underflow: try this:

```
a=1
for i in range(100):
    a*=.0001
    if i % 10 == 0:
        print(i, a)
```

0 0.0001  
10 1.0000000000000003e-44  
20 1.0000000000000007e-84  
30 1.0000000000000001e-124  
40 1.00000000000000015e-164  
50 1.00000000000000021e-204  
60 1.00000000000000029e-244  
70 1.00000000000000036e-284  
80 0.0  
90 0.0

Thankfully, logs are your friend, because of the following graph of  $\log(x)$ :



So instead rewrite the function as

```
from math import log
```

```
def evaluate(sentence, option, model):
    score = log(model['cprobs'][option])
    for word in sentence.split():
        score += log(model['wprobs'][option][word])
    return score
```

(Note the operator change)

## 12 problems with naive bayes assumption

We tolerated the naive bayes assumption:  $P(w_1, w_2, \dots, w_n|y)P(y) = P(w_1|y)P(w_2|y), \dots, P(w_n|y)P(y)$  because it was useful, but we know that it is incorrect. Let's take a closer look:

So, by chain rule,  $P(w_1, w_2, \dots, w_n|y) = P(w_1|y, w_2, \dots, w_n)P(w_2|y, w_3, \dots, w_n) \dots P(w_n|y)$  but by the Naive Bayes assumption we restrict the conditional to just be  $y$ .

Of course, this is wrong! Some words are clearly not conditionally independent of each other, i.e.  $P(\text{San}|y) \neq P(\text{San}|y, \text{Francisco})$ .

A more intuitive example: Imagine if 9/10 people recommended a movie to you. What if 8 of those 9 didn't actually see the movie but just repeated whatever the 9th person said to you?

Naive Bayes may in fact be a decent assumption for the specific case we've seen it in, and for most other words that occur alongside it (except for nearby words), but this is not the only *feature* we might care about. In particular we may care about overlapping features (e.g. the word, the word AND its predecessors/successors, prefix of the word, if the word is on certain lists of words, etc.

## 13 Representation of features and weights for multi-class classification (esp. in Eisenstein)

We're going to talk about features in general but more importantly we're going to talk about feature weights, and especially if you've taken ML before you may be a bit confused by the notation in Eisenstein that I will borrow:

Sentence  $x$ : 'it rocks'

Label possibilities: 'Positive (+), Negative (-), Neutral ( $\emptyset$ )'

Feature classes: "number of words" ( $nW$ ), "contains word ending in s" ( $*s$ ), "contains 'happy'" ( $ch$ )

Feature function  $f$  :

$$f(x, +) = [2, 1, 0, 0, 0, 0, 0, 0, 0]$$

$$f(x, -) = [0, 0, 0, 2, 1, 0, 0, 0, 0]$$

$$f(x, \emptyset) = [0, 0, 0, 0, 0, 0, 2, 1, 0]$$

For every feature we assign a weight. In the previous Naive Bayes discussion the features were all of the form 'contains  $x$ ' for word  $x$  (could actually be 'number of times we see  $x$ '), and the weights were  $P(x|y)$  (from the Naive Bayes assumption). As we see above, the features can be more arbitrary than that. We can arrange our weights in a weight vector, which by convention we'll call  $\theta$ :

$$\theta = [P(nW|+), P(*s|+), P(ch|+), P(nW|-), P(*s|-), P(ch|-), P(nW|\emptyset), P(*s|\emptyset), P(ch|\emptyset)]$$

So given  $f$  and  $\theta$  we can do 'inference' like so:

$$\operatorname{argmax}_{y \in \mathcal{Y}} \theta \cdot f(x, y)$$

Eisenstein uses  $\Psi(x, y) = \theta \cdot f(x, y)$ ; I call that the 'model score' or 'model cost'; it's the model's opinion of  $y$  being suitable for  $x$ .

What about  $P(y)$ ? Note this is the background probability of the class  $y$ . We can include this as a term that is always on for the set of features associated with class  $y$ . It's called the 'bias'. So revising the above,

$f(x, +) = [2, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0]$  and so on...  
 $\theta = [P(nW|+), P(*s|+), P(ch|+), P(+),$   
 $P(nW|-), P(*s|-), P(ch|-), P(-),$   
 $P(nW|\emptyset), P(*s|\emptyset), P(ch|\emptyset), P(\emptyset)]$

## 14 perceptron

I called  $\theta$  ‘weights’ for good reason – who says they have to be probabilities? We can adopt a ‘trial and error’ approach:

```
theta = random weights
```

```
for each sentence, label in data:
    if we would choose some other label over the correct one:
        modify theta so we don't do that
return theta
```

Here it is in a bit more gory detail using the framework from before:

```
import numpy as np
```

```
# note that evaluate needs to be rewritten to be more general;  
# left as an exercise to the reader, but should include a features() method
```

```
def train(labeled_sentences, options, featsize):
    # should be a better way to determine feat size  

    # probably want to initialize differently  

    model = {'theta': np.random(featsize)} # or zeroes  

    for i in range(iterations): # user-determined  

        for sentence, label in labeled_sentences:  

            hyp = classify(sentence, options, model)  

            if hyp != label:  

                model['theta'] += features(sentence, label) - features(sentence, hyp)  

    return model
```

Quick illustration: Some feature type  $f_1$  has value 2,  $f_2$  has value 1,  $f_3$  has value 0. There are three classes.

feats	$f(x, 1)$	$f(x, 2)$	$f(x, 3)$	$\theta$
$f1_1$	2	0	0	0.3
$f2_1$	1	0	0	0.7
$f3_1$	0	0	0	0.8
$f1_2$	0	2	0	-0.2
$f2_2$	0	1	0	2.2
$f3_2$	0	0	0	-.4
$f1_3$	0	0	2	-4
$f2_3$	0	0	1	-4
$f3_3$	0	0	0	-4

Let's say class 1 is correct. Given the table above, we get the following model costs:

$$\Psi(x, 1) = 1.3; \Psi(x, 2) = 1.8; \Psi(x, 3) = -12.$$

So update  $\theta = \theta + f(x, 1) - f(x, 2)$ :

feats	$f(x, 1)$	$f(x, 2)$	$f(x, 3)$	$\theta$
$f1_1$	2	0	0	2.3
$f2_1$	1	0	0	1.7
$f3_1$	0	0	0	0.8
$f1_2$	0	2	0	-2.2
$f2_2$	0	1	0	1.2
$f3_2$	0	0	0	-.4
$f1_3$	0	0	2	-4
$f2_3$	0	0	1	-4
$f3_3$	0	0	0	-4

Class 2's weights went down (if they affected the outcome) and class 1's went up. Class 3's weights weren't involved at all. Now we get the following model costs:

$$\Psi(x, 1) = 6.3; \Psi(x, 2) = -3.2; \Psi(x, 3) = -12.$$

So the item is correctly classified.

Some things to discuss in the context of machine learning: averaging all  $\theta$  at the end, learning rates, batch sizes.

## 15 loss function justification for perceptron

This seemed to work in the demo case above but has a very ad-hoc feel to it. (Side note: often times models *are* originally designed in an ad-hoc way and only later is the theory worked out. The original paper on perceptrons from 1957 has, AFAICT, nothing regarding the following) Why does this work?

It's worth considering the *loss* (often, though not always, written as  $\ell$ ) of the model. Remember, the model conveys an opinion about how to classify<sup>3</sup> that is to some degree untrue. Loss can be thought of as 'how wrong the model is.' For perceptron the loss for item  $i$  in a data set is (from Eisenstein):

$$\ell_{\text{perceptron}}(\theta; \mathbf{x}^{(i)}, y^{(i)}) = \max_{y \in \mathcal{Y}} \theta \cdot f(\mathbf{x}^{(i)}, y) - \theta \cdot f(\mathbf{x}^{(i)}, y^{(i)})$$

---

<sup>3</sup>or for regression models, the value associated with an input

The first term ( $\max_{y \in \mathcal{Y}} \theta \cdot f(\mathbf{x}^{(i)}, y)$ ) is exactly how we pick a label, and the second term is the model score of the right label. If we picked the right label, the loss is zero. Otherwise, the model score for the wrong label is higher, and the amount higher is how wrong we are.

Why do we have loss<sup>4</sup>? Because there's something wrong with our model, i.e.  $\theta$ . But we can change that. How much should we change it? To minimize the loss, of course!

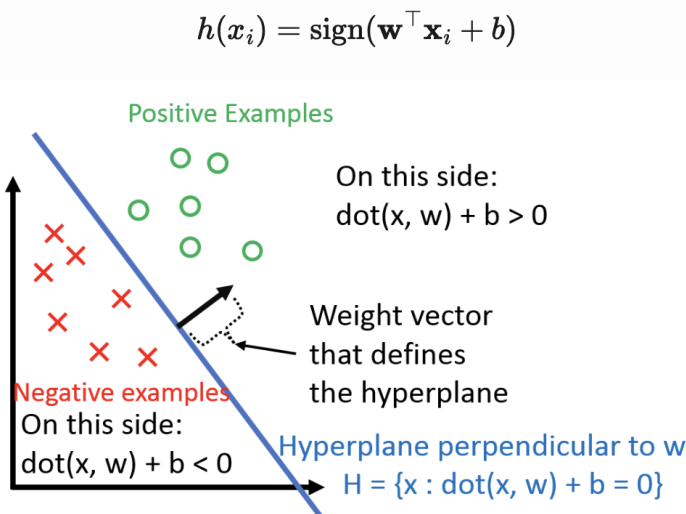
$\ell$  is an equation, and we can take its derivative with respect to  $\theta$ . By adjusting  $\theta$  in the negative direction of the gradient we will have a lower loss on our training data.

Let's assume  $\hat{y}$  is the hypothesis  $y$  and that it's not the true label; what's the derivative? Well what's the equation?

Let's call  $\hat{\mathbf{f}} f(\mathbf{x}, \hat{y})$  and  $\mathbf{f} f(\mathbf{x}, y)$  (I dropped the superscripts; too hard to type) and make them subscriptable. Then  $\ell = \theta_1 \hat{f}_1 - \theta_1 f_1 + \theta_2 \hat{f}_2 - \theta_2 f_2 + \dots + \theta_n \hat{f}_n - \theta_n f_n$ . Then  $\partial \ell / \partial \theta_1 = \hat{f}_1 - f_1$  or to be more vector wise about it,  $\partial \ell / \partial \theta = \hat{\mathbf{f}} - \mathbf{f}$ . So the negative of that is  $\mathbf{f} - \hat{\mathbf{f}}$ .

You might wonder 'isn't there a closed form of this?'<sup>5</sup> Yes, there is, in the sense of this being a solution for a general linear model. Let the features be matrix  $A$  of dimension  $N \times F$  for  $N$  training data and  $F$  features, and let numeric (integral?) labels be a vector  $B$  ( $N \times 1$ )<sup>6</sup>. Let the weights be vector  $W$  ( $F \times 1$ ). Then a linear model can be expressed  $A \times W = B$ . Then  $W = A^{-1}B$  but usually  $A^{-1}$  isn't invertible, so  $W = (A^T A)^{-1} A^T B$ .<sup>7</sup> But for large  $N$  that inverse operation is quite slow, so this isn't usually done.

BTW, another way to look at what the perceptron is doing is finding the *separating hyperplane* between differently labeled samples. This is easiest to visualize in the binary case:



The weights define a line/plane/hyperplane along which the score is zero; we want all the differently labeled examples to be divided by that separator.

<sup>4</sup>Assuming our data is *separable*, which see below

<sup>5</sup>i.e. a single equation that can be solved to give ideal weights, as in Naive Bayes

<sup>6</sup>I think any unique value per label should do but FYI the literature I'm seeing, which is all over the place on if this is even possible, at best shows examples with labels 1 and -1 only.

<sup>7</sup> $(A^T A)^{-1} A^T$  is a pseudo-inverse of  $A$ .

## 16 support vector machine

Very briefly, perceptron updates are as follows:

$$\hat{y} = \operatorname{argmax}_y \theta \cdot f(x^{(i)}, y)$$

If  $\hat{y} = y^{(i)}$  then no update, otherwise

$$\theta = \theta + f(x^{(i)}, y^{(i)}) - f(x^{(i)}, \hat{y})$$

Introduce cost function for making a mistake  $c(y^{(i)}, y)$ ; (for instance, let's just say cost is 1 if  $y$  is wrong, 0 otherwise). Change the updates as follows:

$$\hat{y} = \operatorname{argmax}_y \theta \cdot f(x^{(i)}, y) + c(y^{(i)}, y)$$

$$\theta = \lambda \theta + f(x^{(i)}, y^{(i)}) - f(x^{(i)}, \hat{y})$$

8

This is going to choose  $\hat{y}$  that are close to  $y^{(i)}$  instead of useless updates with the correct values, and it scales the previous weights to avoid overfitting. This is called *support vector machine* and has a nice derivation that we would spend time on in an ML class, but the upshot is it can work better than perceptron (see the Pang et al. paper in the reading, for instance).

## 17 logistic regression

One nice thing about the naive bayes model is that it's probabilistic, so if your classifier is one part of a pipeline, you can tell the rest of the pipeline your confidence in your output in a rational way.  $\Psi(x, y)$  has range  $(-\infty, \infty)$  which is less helpful. Another issue with perceptron/svm is that they are only concerned with making sure the correct answer is chosen/has large margin (in training).

Nobody forced us to keep  $\Psi(x, y)$  as the model score. A preferred model score would be  $P(y|x)$ , the conditional probability of the output given the input. How do we form a probability distribution from a set of scores?

We can't simply normalize:

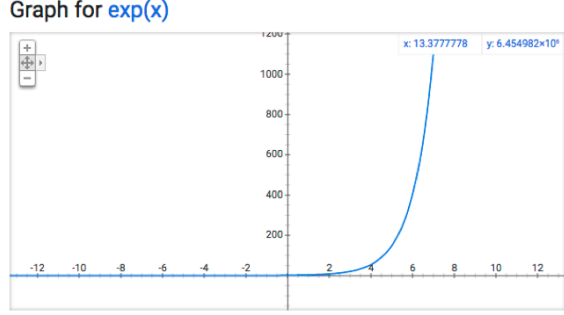
$$\frac{\Psi(x, y)}{\sum_{y' \in \mathcal{Y}} \Psi(x, y')}$$

will not work. Why?

Instead we can use  $e^\Psi$ . why?

---

<sup>8</sup>In Eisenstein  $(1 - \lambda)$  is used and this fits the derivation better.



So our new model cost will be

$$\frac{e^{\Psi(x,y)}}{\sum_{y' \in \mathcal{Y}} e^{\Psi(x,y')}} \quad (1)$$

This, by the way, is the ‘softmax’ function that comes up a lot in neural networks. It does exactly the same job we are trying to do here: convert a set of numbers in the range  $(-\infty, \infty)$  into a distribution while keeping the ordering the same. In fact, it generally ‘peaks’ the highest number; that’s why this can be thought of as a ‘soft’ form of the ‘max’ operator (hence the name).

Now that we’ve got a distribution we have a natural loss function we can use: cross-entropy!  $H(p, q) \triangleq -E_p \log q$ , where  $E_p$  is the *expected value* w/r/t  $p$ . Concretely:

$$H(p, q) = - \sum_{x \in \mathcal{X}} p(x) \log q(x)$$

Where here,  $q(x) = q(y|x) = \text{Equation 1}$ . From information theory this is the average number of bits<sup>9</sup> needed to identify an item if  $q$  is used to identify the item, but  $p$  is the true distribution. So intuitively you want as small a cross-entropy as possible, and that makes this a natural loss function. But what is  $p(y|x)$ , the true distribution? A good choice for supervised learning is to assume that the provided label is the truth and should always occur and that all other labels should never occur. That means that our calculation of cross entropy, which will in practice use training set  $\mathcal{X}$  and label space  $\mathcal{Y}$ :

$$- \frac{\sum_{x,y \in \mathcal{X}} \sum_{y' \in \mathcal{Y}} p(y'|x) \log q(y'|x)}{|\mathcal{X}|}$$

becomes

$$- \frac{\sum_{x,y \in \mathcal{X}} \log q(y|x)}{|\mathcal{X}|}$$

since  $p(y|x) = 1$  and  $\forall y' \in \mathcal{Y} \setminus y, p(y'|x) = 0$ .

Now considering a single item and substituting back in Equation 1:

$$- \log \frac{\exp(\theta \cdot f(\mathbf{x}, y))}{\sum_{y' \in \mathcal{Y}} \exp(\theta \cdot f(\mathbf{x}, y'))}$$

---

<sup>9</sup>really, nats

$$-\theta \cdot f(\mathbf{x}, y) + \log \sum_{y' \in \mathcal{Y}} \exp(\theta \cdot f(\mathbf{x}, y'))$$

Now, the gradient:

$$\partial \ell / \partial \theta = -f(\mathbf{x}, y) + \frac{1}{\sum_{y'' \in \mathcal{Y}} \exp(\theta \cdot f(\mathbf{x}, y''))} \sum_{y' \in \mathcal{Y}} \exp(\theta \cdot f(\mathbf{x}, y')) f(\mathbf{x}, y')$$

$$\partial \ell / \partial \theta = -f(\mathbf{x}, y) + \sum_{y' \in \mathcal{Y}} \frac{\exp(\theta \cdot f(\mathbf{x}, y'))}{\sum_{y'' \in \mathcal{Y}} \exp(\theta \cdot f(\mathbf{x}, y''))} f(\mathbf{x}, y')$$

$$\partial \ell / \partial \theta = -f(\mathbf{x}, y) + \sum_{y' \in \mathcal{Y}} q(y' | \mathbf{x}; \theta) f(\mathbf{x}, y')$$

That second term is the expectation of  $Y|X$ , which is to say, the probability (according to the model) of a value times that value, summed over all possible values. Naturally, we *reduce*  $\theta$  by (some fraction of) this value.

Note that we are now considering not only how far away each wrong answer is from the right answer, but how confident we are about each wrong answer. If we have low confidence about an answer it will affect the loss very little.