# GRADIENT DESCENT

David Kauchak
CS 158 – Spring 2022

1

## Admin

Assignment 3 graded

Assignment 5 out
- Course feedback

Midterm next week

Assignment 6 will also be next week

2

## Midterm details

Time limited take home exam (you'll have 2 hours to complete it)

Available on Monday (2/21)

Must finish by end of the day on Friday (2/25)

You may use your notes, the class notes, the class book(s), and your assignments

You may NOT use any other resources on the web or search for things on the web

3

## Midterm topics

| Date | Topic |
|------|-------|
| 1/18 | introduction (ppt) |
| 1/20 | decision trees (ppt) |
| 1/25 | geometric view of data (ppt) |
| 1/27 | perceptron (ppt) |
| 2/1 | features (ppt) |
| 2/3 | evaluation (ppt) |
| 2/8 | imbalanced data (ppt) |
| 2/10 | beyond binary classification (ppt) |
| 2/15 | gradient descent |
| 2/17 | regularization |

**(More details on Wednesday!)**

4

## Midterm topics

Machine learning basics
- different types of learning problems
- feature-based machine learning
- data assumptions/data generating distribution

Classification problem setup

Proper experimentation
- train/dev/test
- evaluation/accuracy/training error
- optimizing hyperparameters

5

## Midterm topics

Learning algorithms
- Decision trees
- K-NN
- Perceptron
- Gradient descent

Algorithm properties
- training/learning
  - rational/why it works
- classifying
- hyperparameters
- avoiding overfitting
- algorithm variants/improvements

6

## Midterm topics

Geometric view of data
- distances between examples
- decision boundaries

Features
- example features
- removing erroneous features/picking good features
- challenges with high-dimensional data
- feature normalization

Other pre-processing
- outlier detection

7

## Midterm topics

Comparing algorithms
- n-fold cross validation
  - leave one out validation
- bootstrap resampling
- t-test

imbalanced data
- evaluation
  - precision/recall, F1, AUC
- subsampling
- oversampling
- weighted binary classifiers

8

## Midterm topics

Multiclass classification
- Modifying existing approaches
- Using binary classifier
  - OVA
  - AVA
  - Tree-based
- micro- vs. macro-averaging

Ranking
- using binary classifier
- using weighted binary classifier

9

## Midterm topics

Gradient descent
- 0/1 loss
- Surrogate loss functions
- Convexity
- minimization algorithm
- regularization
  - different regularizers
  - p-norms

Misc
- good coding habits
- JavaDoc

10

## Midterm general advice

2 hours goes by fast!
- Don't plan on looking everything up
  - Lookup equations, algorithms, random details
- Make sure you understand the key concepts
- Don't spend too much time on any one question
  - Skip questions you're stuck on and come back to them
- Watch the time as you go

Be careful on the T/F questions

For written questions
- think before you write
- make your argument/analysis clear and concise

11

## An aside: text classification

Raw data          labels

Chardonnay

Pinot Grigio

Zinfandel

12

## Text: raw data

| Raw data | labels | Features? |
|---|---|---|
| | Chardonnay | |
| | Pinot Grigio | |
| | Zinfandel | |

13

## Feature examples

| Raw data | labels | Features |
|---|---|---|
| | Chardonnay | Clinton said pinot repeatedly last week on tv, "pinot, pinot, pinot" |
| | Pinot Grigio | (1, 1, 1, 0, 0, 1, 0, 0, ...) |
| | Zinfandel | Occurrence of words |

(words: pinot, clinton, said, california, across, tv, wrong, capital)

14

## Feature examples

| Raw data | labels | Features |
|---|---|---|
| | Chardonnay | Clinton said pinot repeatedly last week on tv, "pinot, pinot, pinot" |
| | Pinot Grigio | (4, 1, 1, 0, 0, 1, 0, 0, ...) |
| | Zinfandel | Frequency of word occurrences |

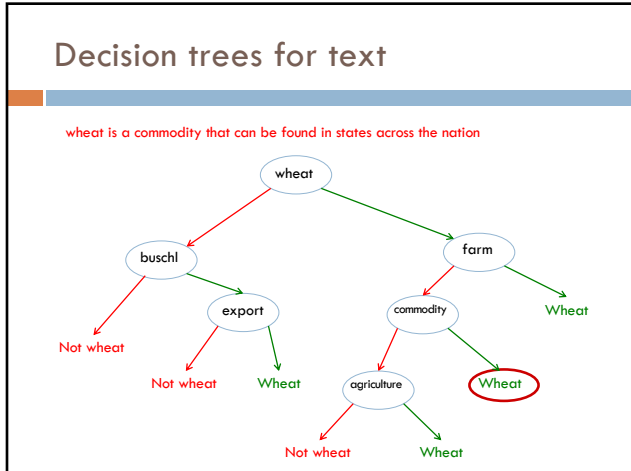(words: pinot, clinton, said, california, across, tv, wrong, capital)

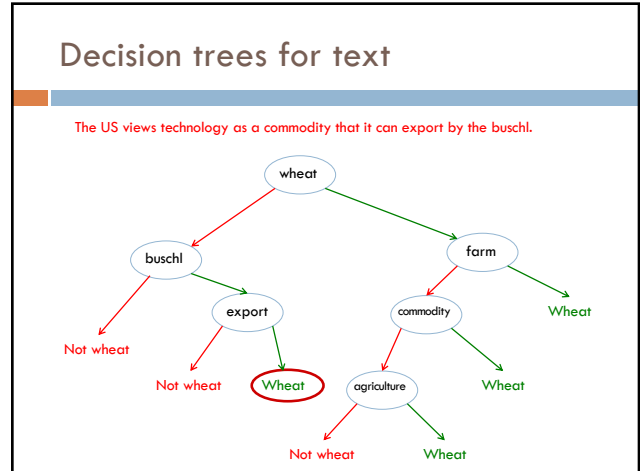This is the representation we're using for assignment 5

15

## Decision trees for text

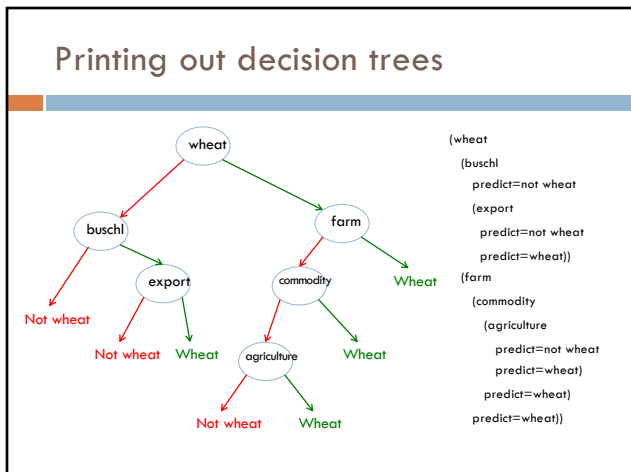Each internal node represents whether or not the text has a particular word
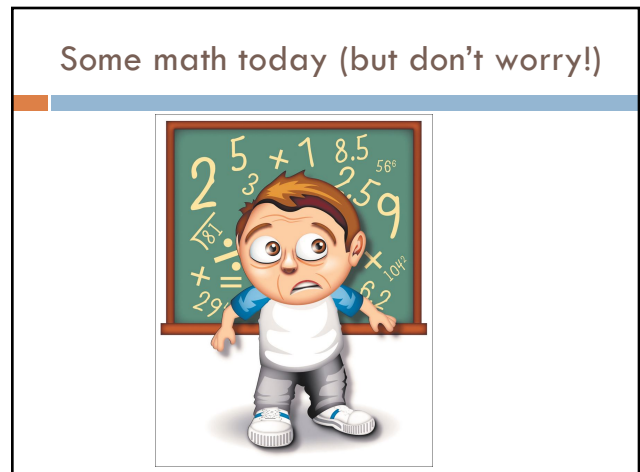


16

## Decision trees for text

wheat is a commodity that can be found in states across the nation



17

## Decision trees for text

The US views technology as a commodity that it can export by the buschl.



18

## Printing out decision trees



```
(wheat
  (buschl
    predict=not wheat
    (export
      predict=not wheat
      predict=wheat))
  (farm
    (commodity
      (agriculture
        predict=not wheat
        predict=wheat)
      predict=wheat)
    predict=wheat))
```

19

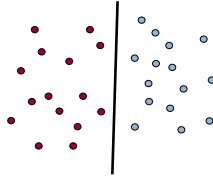## Some math today (but don't worry!)



20

2/15/22

5

## Linear models

A high-bias assumption is *linear separability*:

- in 2 dimensions, can separate classes by a line
- in higher dimensions, need hyperplanes

A *linear model* is a model that assumes the data is linearly separable



21

## Linear models

A linear model in *n*-dimensional space (i.e. *n* features) is define by *n*+1 weights:

In two dimensions, a line:
$$0 = w_1 f_1 + w_2 f_2 + b$$  (where b = -a)

In three dimensions, a plane:
$$0 = w_1 f_1 + w_2 f_2 + w_3 f_3 + b$$

In *m*-dimensions, a *hyperplane*
$$0 = b + \sum_{j=1}^{m} w_j f_j$$

22

## Perceptron learning algorithm

repeat until convergence (or for some # of iterations):

for each training example ($f_1, f_2, \ldots, f_m$, label):

$$prediction = b + \sum_{j=1}^{m} w_j f_j$$

if *prediction* * label ≤ 0:  // they don't agree
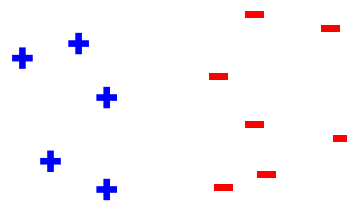
for each $w_i$:

$w_i = w_i + f_i$*label

$b = b + $label

23

## Which line will it find?



24

## Which line will it find?



Only guaranteed to find **some**
line that separates the data

25

## Linear models

Perceptron algorithm is one example of a linear classifier

Many, many other algorithms learn a line (i.e. a setting of a linear combination of weights)

Goals:

- Explore a number of linear training algorithms
- Understand *why these algorithms work*

26

## Perceptron learning algorithm

repeat until convergence (or for some # of iterations):

   for each training example ($f_1, f_2, ..., f_m$, label):

$$prediction = b + \sum_{j=1}^{m} w_j f_j$$

   if *prediction * label* $\leq$ 0:  // they don't agree

     for each $w_i$:

       $w_i = w_i + f_i$*label

     $b = b + label$

27

## A closer look at why we got it wrong

$w_1$     $w_2$            (-1, -1, positive)

$$0 * f_1 + 1 * f_2 =$$

$$0 * -1 + 1 * -1 = -1 \leftarrow$$

We'd like this value to be positive since it's a positive value

didn't contribute, but could have

contributed in the wrong direction

Intuitively these make sense
Why change by 1?
Any other way of doing it?

decrease              decrease

0 -> -1                1 -> 0

28

## Model-based machine learning

1. pick a model
   - e.g. a hyperplane, a decision tree,…
   - A model is defined by a collection of parameters

   What are the parameters for DT? Perceptron?

29

## Model-based machine learning

1. pick a model
   - e.g. a hyperplane, a decision tree,…
   - A model is defined by a collection of parameters

   DT: the structure of the tree, which features each node splits on, the predictions at the leaves

   perceptron: the weights and the b value

30

## Model-based machine learning

1. pick a model
   - e.g. a hyperplane, a decision tree,…
   - A model is defined by a collection of parameters

2. pick a criterion to optimize (aka objective function)

   What criteria do decision tree learning and perceptron learning optimizing?

31

## Model-based machine learning

1. pick a model
   - e.g. a hyperplane, a decision tree,…
   - A model is defined by a collection of parameters

2. pick a criterion to optimize (aka objective function)
   - e.g. training error

3. develop a learning algorithm
   - the algorithm should try and minimize the criteria
   - sometimes in a heuristic way (i.e. non-optimally)
   - sometimes exactly

32

## Linear models in general

1. pick a model

$$0 = \boxed{b} + \sum_{j=1}^{m} \boxed{w_j} f_j$$

These are the parameters we want to learn

2. pick a criterion to optimize (aka objective function)

33

## Some notation: indicator function

$$1[x] = \begin{cases} 1 & if \ x = True \\ 0 & if \ x = False \end{cases}$$

Convenient notation for turning T/F answers into numbers/counts:

$$beers\_to\_bring\_for\_class = \sum_{age \in class} 1[age >= 21]$$

34

## Some notation: dot-product

Sometimes it is convenient to use vector notation

We represent an example $f_1, f_2, \ldots, f_m$ as a single vector, $x$
- j subscript will indicate feature indexing, i.e., $x_j$
- i subscript will indicate examples indexing over a dataset, i.e., $x_i$ or sometimes $x_{ij}$

Similarly, we can represent the weight vector $w_1, w_2, \ldots, w_m$ as a single vector, $w$

The dot-product between two vectors $a$ and $b$ is defined as:

$$a \cdot b = \sum_{j=1}^{m} a_j b_j$$

35

## Linear models

1. pick a model

$$0 = \boxed{b} + \sum_{j=1}^{n} \boxed{w_j} f_j$$

These are the parameters we want to learn

2. pick a criterion to optimize (aka objective function)

$$\sum_{i=1}^{n} 1[y_i(w \cdot x_i + b) \leq 0]$$

What does this equation say?

36

## 0/1 loss function

$$\sum_{i=1}^{n} 1\left[y_i(w \cdot x_i + b) \le 0\right]$$

- distance from hyperplane
- sign is prediction

whether or not the prediction and label agree, true if **they don't**

total number of mistakes, aka 0/1 loss

37

## Model-based machine learning

1. pick a model

$$0 = b + \sum_{j=1}^{m} w_j f_j$$

2. pick a criteria to optimize (aka objective function)

$$\sum_{i=1}^{n} 1\left[y_i(w \cdot x_i + b) \le 0\right]$$

3. develop a learning algorithm

$$\operatorname{argmin}_{w,b} \sum_{i=1}^{n} 1\left[y_i(w \cdot x_i + b) \le 0\right]$$

Find w and b that minimize the 0/1 loss (i.e. training error)

38

## Minimizing 0/1 loss

$$\operatorname{argmin}_{w,b} \sum_{i=1}^{n} 1\left[y_i(w \cdot x_i + b) \le 0\right]$$
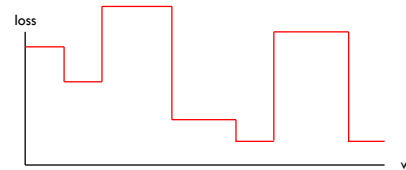
Find w and b that minimize the 0/1 loss

How do we do this?
How do we *minimize* a function?
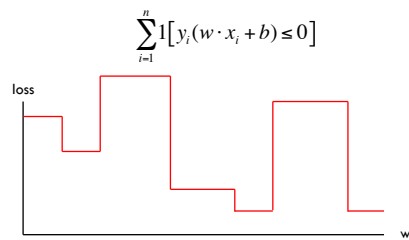Why is it hard for this function?

39

## Minimizing 0/1 in one dimension

$$\sum_{i=1}^{n} 1\left[y_i(w \cdot x_i + b) \le 0\right]$$

loss

w

Each time we change w such that the example is right/wrong the loss will increase/decrease

40

## Minimizing 0/1 over all w

$$\sum_{i=1}^{n} 1\left[y_i(w \cdot x_i + b) \le 0\right]$$

loss



w

Each new feature we add (i.e. weights) adds another dimension to this space!

41

## Minimizing 0/1 loss

$$\text{argmin}_{w,b} \sum_{i=1}^{n} 1\left[y_i(w \cdot x_i + b) \le 0\right]$$

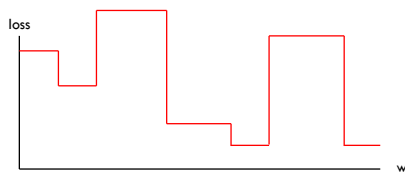Find w and b that minimize the 0/1 loss

This turns out to be hard (in fact, NP-HARD ☹)

Challenge:
- small changes in any w can have large changes in the loss (the change isn't continuous)
- there can be many, many local minima
- at any given point, we don't have much information to direct us towards any minima
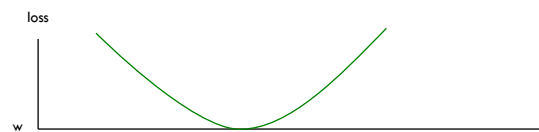
42

## More manageable loss functions

loss



w

What property/properties do we want from our loss function?

43

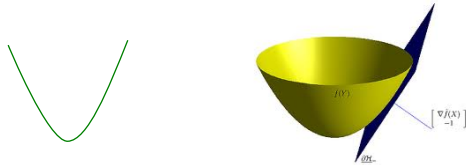## More manageable loss functions

loss



w

- Ideally, continuous (i.e. differentiable) so we get an indication of direction of minimization
- Only one minima

44

11

## Convex functions

Convex functions look something like:



One definition: The line segment between any two points on the function is *above* the function

45

## Surrogate loss functions

For many applications, we really would like to minimize the 0/1 loss

A surrogate loss function is a loss function that provides an upper bound on the actual loss function (in this case, 0/1)

We'd like to identify a convex surrogate loss functions to make them easier to minimize

Key to a loss function: how it scores the difference between the actual label *y* and the predicted label *y'*

46

## Surrogate loss functions

0/1 loss:     $l(y,y') = 1[yy' \leq 0]$

Ideas?
Some function that is a proxy for error, but is continuous and convex

47

## Surrogate loss functions
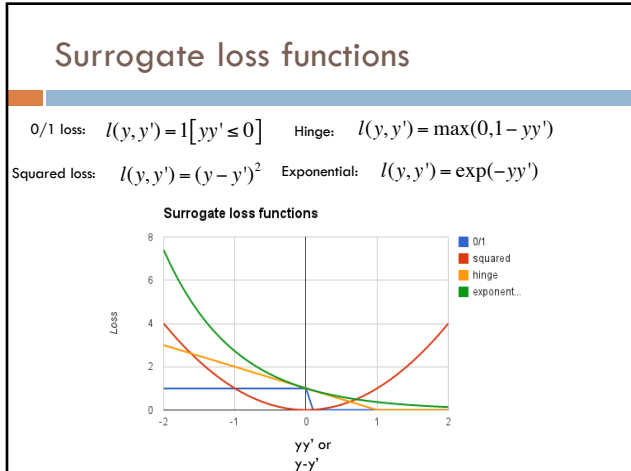
0/1 loss:     $l(y,y') = 1[yy' \leq 0]$

Hinge:     $l(y,y') = \max(0, 1 - yy')$

Exponential:     $l(y,y') = \exp(-yy')$

Squared loss:     $l(y,y') = (y - y')^2$

Why do these work?  What do they penalize?

48

## Surrogate loss functions

0/1 loss: $l(y, y') = 1[yy' \le 0]$  Hinge: $l(y, y') = \max(0, 1 - yy')$

Squared loss: $l(y, y') = (y - y')^2$  Exponential: $l(y, y') = \exp(-yy')$



Surrogate loss functions

49

## Model-based machine learning

1. pick a model

$$0 = b + \sum_{j=1}^{m} w_j f_j$$

2. pick a criteria to optimize (aka objective function)

$$\sum_{i=1}^{n} \exp(-y_i(w \cdot x_i + b))$$

use a convex surrogate loss function

3. develop a learning algorithm

$$\operatorname{argmin}_{w,b} \sum_{i=1}^{n} \exp(-y_i(w \cdot x_i + b))$$

Find w and b that minimize the surrogate loss

50

## Finding the minimum



You're blindfolded, but you can see out of the bottom of the blindfold to the ground right by your feet.  I drop you off somewhere and tell you that you're in a convex shaped valley and escape is at the bottom/minimum.  How do you get out?
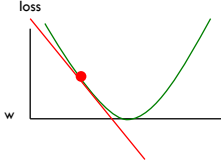
51

## Finding the minimum



loss

w

How do we do this for a function?

52

## One approach: gradient descent

Partial derivatives give us the slope (i.e. direction to move) in that dimension



53

## One approach: gradient descent

Partial derivatives give us the slope (i.e. direction to move) in that dimension

Approach:
- ☐ pick a starting point (w)
- ☐ repeat:
  - ■ pick a dimension
  - ■ move a small amount in that dimension towards decreasing loss (using the derivative)
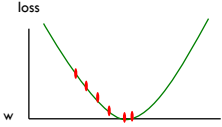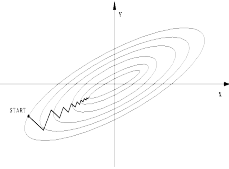


54

## One approach: gradient descent

Partial derivatives give us the slope (i.e. direction to move) in that dimension

Approach:
- ☐ pick a starting point (w)
- ☐ repeat:
  - ■ pick a dimension
  - ■ move a small amount in that dimension towards decreasing loss (using the derivative)



55

## Gradient descent

- ☐ pick a starting point (w)
- ☐ repeat until loss doesn't decrease in any dimension:
  - ■ pick a dimension
  - ■ move a small amount in that dimension towards decreasing loss (using the derivative)

$$w_j = w_j - \frac{d}{dw_j} loss(w)$$

Why negative?

56

14

## Gradient descent

- pick a starting point (w)
- repeat until loss doesn't decrease in any dimension:
  - pick a dimension
  - move a small amount in that dimension towards decreasing loss (using the derivative)

$$w_j = w_j - \eta \; \frac{d}{dw_j} loss(w)$$

What does this do?

57

## Gradient descent

- pick a starting point (w)
- repeat until loss doesn't decrease in any dimension:
  - pick a dimension
  - move a small amount in that dimension towards decreasing loss (using the derivative)

$$w_j = w_j - \eta \; \frac{d}{dw_j} loss(w)$$

learning rate (how much we want to move in the error direction, often this will change over time)

58

## Some math

$$\frac{d}{dw_j} loss = \frac{d}{dw_j} \sum_{i=1}^{n} \exp(-y_i(w \cdot x_i + b))$$

$$= \sum_{i=1}^{n} \exp(-y_i(w \cdot x_i + b)) \frac{d}{dw_j} - y_i(w \cdot x_i + b)$$

59

## Some math

$$-\frac{d}{dw_j} y_i(w \cdot x_i + b) = -\frac{d}{dw_j} y_i(\sum_{j=1}^{m} w_j x_{ij} + b)$$

$$= -\frac{d}{dw_j} y_i(w_1 x_{i1} + w_2 x_{i2} + \ldots + w_m x_{im} + b)$$

$$= -\frac{d}{dw_j} y_i w_1 x_{i1} + y_i w_2 x_{i2} + \ldots + y_i w_m x_{im} + y_i b)$$

$$= -y_i x_{ij}$$

60

15

## Some math

$$\frac{d}{dw_j} loss = \frac{d}{dw_j} \sum_{i=1}^{n} \exp(-y_i(w \cdot x_i + b))$$

$$= \sum_{i=1}^{n} \exp(-y_i(w \cdot x_i + b)) \frac{d}{dw_j} - y_i(w \cdot x_i + b)$$

$$= \sum_{i=1}^{n} -y_i x_{ij} \exp(-y_i(w \cdot x_i + b))$$

61

## Gradient descent

- ☐ pick a starting point (w)
- ☐ repeat until loss doesn't decrease in any dimension:
  - ■ pick a dimension
  - ■ move a small amount in that dimension towards decreasing loss (using the derivative)

$$w_j = w_j + \eta \sum_{i=1}^{n} y_i x_{ij} \exp(-y_i(w \cdot x_i + b))$$

**What is this doing?**

62

## Exponential update rule

$$w_j = w_j + \eta \sum_{i=1}^{n} y_i x_{ij} \exp(-y_i(w \cdot x_i + b))$$

for each example $x_i$:

$$w_j = w_j + \eta y_i x_{ij} \exp(-y_i(w \cdot x_i + b))$$

**Does this look familiar?**

63

## Perceptron learning algorithm!

repeat until convergence (or for some # of iterations):

for each training example ($f_1$, $f_2$, …, $f_m$, label):

$$prediction = b + \sum_{j=1}^{m} w_j f_j$$

if *prediction* * *label* ≤ 0:  // they don't agree

for each $w_i$:

$w_i = w_i + f_i$*label

$b = b$ + label

$$w_j = w_j + \eta y_i x_{ij} \exp(-y_i(w \cdot x_i + b))$$

or

$$w_j = w_j + x_{ij} y_i c \quad \text{where} \quad c = \eta \exp(-y_i(w \cdot x_i + b))$$

64

## The constant

$$c = \eta \exp(-y_i(w \cdot x_i + b))$$

learning rate   label          prediction

When is this large/small?

65

## The constant

$$c = \eta \exp(-y_i(w \cdot x_i + b))$$

label          prediction

If they're the same sign, as the predicted gets larger the update gets smaller

If they're different, the more different they are, the bigger the update

66

## Perceptron learning algorithm!

repeat until convergence (or for some # of iterations):

for each training example ($f_1$, $f_2$, …, $f_m$, label):

$$prediction = b + \sum_{j=1}^{m} w_j f_j$$

~~if prediction * label ≤ 0:~~ // they don't agree

for each $w_i$:                  Note: for gradient descent, we always update

$w_i = w_i + f_i$*label

$b = b + $ label

$$w_j = w_j + \eta y_i x_{ij} \exp(-y_i(w \cdot x_i + b))$$

or

$$w_j = w_j + x_{ij} y_i c \quad \text{where} \quad c = \eta \exp(-y_i(w \cdot x_i + b))$$

67

## One concern

$$\text{argmin}_{w,b} \sum_{i=1}^{n} \exp(-y_i(w \cdot x_i + b))$$

loss

w

We're calculating this on the **training set**

We still need to be careful about overfitting!

The min w,b on the training set is generally NOT the min for the test set

How did we deal with this for the perceptron algorithm?

68

# Summary

Model-based machine learning:
- define a model, objective function (i.e. loss function), minimization algorithm

Gradient descent minimization algorithm
- require that our loss function is convex
- make small updates towards lower losses

Perceptron learning algorithm:
- gradient descent
- exponential loss function (modulo a learning rate)

69