# FEATURE PRE-PROCESSING

David Kauchak
CS 451 – Fall 2013

## Admin

Assignment 3

## So far…

1. Throw out outlier examples

## Feature pruning/selection

Good features provide us information that helps us distinguish between labels. However, not all features are good

Feature pruning is the process of removing "bad" features

Feature selection is the process of selecting "good" features

What makes a bad feature and why would we have them in our data?

## Bad features

Each of you are going to generate a feature for our data set: pick 5 random binary numbers

$f_1$ $f_2$ ...        label

☐

☐

☐   I've already labeled these examples and I have two features

☐

☐

## Bad features

label

1

0

1

1

0

If we have a "random" feature, i.e. a feature with random binary values, what is the probability that our feature perfectly predicts the label?

## Bad features

| label | $f_i$ | probability |
|-------|-------|-------------|
| 1 | 1 | 0.5 |
| 0 | 0 | 0.5 |
| 1 | 1 | 0.5 |
| 1 | 1 | 0.5 |
| 0 | 0 | 0.5 |

$0.5^5 = 0.03125 = 1/32$

Is that the only way to get perfect prediction?

## Bad features

| label | $f_i$ | probability |
|-------|-------|-------------|
| 1 | 0 | 0.5 |
| 0 | 1 | 0.5 |
| 1 | 0 | 0.5 |
| 1 | 0 | 0.5 |
| 0 | 1 | 0.5 |

$0.5^5 = 0.03125 = 1/32$

Total = 1/32+1/32 = 1/16

Why is this a problem?

Although these features perfectly correlate/predict the training data, they will generally NOT have any predictive power on the test set!

2

## Bad features

| label | $f_i$ | probability |
|-------|-------|-------------|
| 1 | 0 | 0.5 |
| 0 | 1 | 0.5 |
| 1 | 0 | 0.5 |
| 1 | 0 | 0.5 |
| 0 | 1 | 0.5 |

Total = 1/32+1/32 = 1/16

$0.5^5 = 0.03125 = 1/32$

Is perfect correlation the only thing we need to worry about for random features?

## Bad features

| label | $f_i$ |
|-------|-------|
| 1 | 1 |
| 0 | 0 |
| 1 | 1 |
| 1 | 0 |
| 0 | 0 |

Any correlation (particularly any strong correlation) can affect performance!

## Noisy features

Adding features *can* give us more information, but not always

Determining if a feature is useful can be challenging

| Terrain | Unicycle-type | Weather | Jacket | ML grade | Go-For-Ride? |
|---------|---------------|---------|--------|----------|--------------|
| Trail | Mountain | Rainy | Heavy | D | YES |
| Trail | Mountain | Sunny | Light | C- | YES |
| Road | Mountain | Snowy | Light | B | YES |
| Road | Mountain | Sunny | Heavy | A | YES |
| Trail | Normal | Snowy | Light | D+ | NO |
| Trail | Normal | Rainy | Heavy | B- | NO |
| Road | Normal | Snowy | Heavy | C+ | YES |
| Road | Normal | Sunny | Light | A- | NO |
| Trail | Normal | Sunny | Heavy | B+ | NO |
| Trail | Normal | Snowy | Light | F | NO |
| Trail | Normal | Rainy | Light | C | YES |

## Noisy features

These can be particularly problematic in problem areas where we automatically generate features

Features

Clinton said banana repeatedly last week on tv, "banana, banana, banana"

(1, 1, 1, 0, 0, 1, 0, 0, …)

clinton said
said banana
across the
california schools
tv banana
wrong way
capital city

## Noisy features

Ideas for removing noisy/random features?

| Terrain | Unicycle-type | Weather | Jacket | ML grade | Go-For-Ride? |
|---------|---------------|---------|--------|----------|--------------|
| Trail | Mountain | Rainy | Heavy | D | YES |
| Trail | Mountain | Sunny | Light | C- | YES |
| Road | Mountain | Snowy | Light | B | YES |
| Road | Mountain | Sunny | Heavy | A | YES |
| Trail | Normal | Snowy | Light | D+ | NO |
| Trail | Normal | Rainy | Heavy | B- | NO |
| Road | Normal | Snowy | Heavy | C+ | YES |
| Road | Normal | Sunny | Light | A- | NO |
| Trail | Normal | Sunny | Heavy | B+ | NO |
| Trail | Normal | Snowy | Light | F | NO |
| Trail | Normal | Rainy | Light | C | YES |

## Removing noisy features

The expensive way:
- Split training data into train/dev
- Train a model on all features
- for each feature f:
  - Train a model on all features – f
  - Compare performance of all vs. all-f on dev set

- Remove all features where decrease in performance between all and all-f is less than some constant

Feature ablation study          Issues/concerns?

## Removing noisy features

Binary features:
remove "rare" features, i.e. features that only occur (or don't occur) a very small number of times

Real-valued features:
remove features that have low variance

In both cases, can either use thresholds, throw away lowest x%, use development data, etc.

Why?

## Some rules of thumb for the number of features

Be very careful in domains where:
- the number of features > number of examples
- the number of features ≈ number of examples
- the features are generated automatically
- there is a chance of "random" features

In most of these cases, features should be removed based on some domain knowledge (i.e. problem-specific knowledge)

## So far…

1. Throw out outlier examples
2. Remove noisy features
3. Pick "good" features

## Feature selection

Let's look at the problem from the other direction, that is, selecting good features.

What are good features?

How can we pick/select them?

## Good features

A good feature correlates well with the label

label

| | | | |
|---|---|---|---|
| 1 | 1 | 0 | 1 |
| 0 | 0 | 1 | 1 |
| 1 | 1 | 0 | 1 |
| 1 | 1 | 0 | 1 |
| 0 | 0 | 1 | 0 |

…

How can we identify this?
- training error (like for DT)
- correlation model
- statistical test
- probabilistic test
- …

## Training error feature selection

- for each feature f:
  - calculate the training error if only feature f were used to pick the label

- rank each feature by this value
- pick top k, top x%, etc.
  - can use a development set to help pick k or x

## So far…

1. Throw out outlier examples
2. Remove noisy features
3. Pick "good" features

## Feature normalization

| Length | Weight | Color | Label |
|--------|--------|-------|-------|
| 4 | 4 | 0 | Apple |
| 5 | 5 | 1 | Apple |
| 7 | 6 | 1 | Banana |
| 4 | 3 | 0 | Apple |
| 6 | 7 | 1 | Banana |
| 5 | 8 | 1 | Banana |
| 5 | 6 | 1 | Apple |

| Length | Weight | Color | Label |
|--------|--------|-------|-------|
| 40 | 4 | 0 | Apple |
| 50 | 5 | 1 | Apple |
| 70 | 6 | 1 | Banana |
| 40 | 3 | 0 | Apple |
| 60 | 7 | 1 | Banana |
| 50 | 8 | 1 | Banana |
| 50 | 6 | 1 | Apple |

Would our three classifiers (DT, k-NN and perceptron) learn the same models on these two data sets?

## Feature normalization

| Length | Weight | Color | Label |
|--------|--------|-------|-------|
| 4 | 4 | 0 | Apple |
| 5 | 5 | 1 | Apple |
| 7 | 6 | 1 | Banana |
| 4 | 3 | 0 | Apple |
| 6 | 7 | 1 | Banana |
| 5 | 8 | 1 | Banana |
| 5 | 6 | 1 | Apple |

| Length | Weight | Color | Label |
|--------|--------|-------|-------|
| 40 | 4 | 0 | Apple |
| 50 | 5 | 1 | Apple |
| 70 | 6 | 1 | Banana |
| 40 | 3 | 0 | Apple |
| 60 | 7 | 1 | Banana |
| 50 | 8 | 1 | Banana |
| 50 | 6 | 1 | Apple |

Decision trees don't care about scale, so they'd learn the same tree

## Feature normalization

| Length | Weight | Color | Label |
|--------|--------|-------|-------|
| 4 | 4 | 0 | Apple |
| 5 | 5 | 1 | Apple |
| 7 | 6 | 1 | Banana |
| 4 | 3 | 0 | Apple |
| 6 | 7 | 1 | Banana |
| 5 | 8 | 1 | Banana |
| 5 | 6 | 1 | Apple |

| Length | Weight | Color | Label |
|--------|--------|-------|-------|
| 40 | 4 | 0 | Apple |
| 50 | 5 | 1 | Apple |
| 70 | 6 | 1 | Banana |
| 40 | 3 | 0 | Apple |
| 60 | 7 | 1 | Banana |
| 50 | 8 | 1 | Banana |
| 50 | 6 | 1 | Apple |

k-NN: NO! The distances are biased based on feature magnitude.

$$D(a,b) = \sqrt{(a_1 - b_1)^2 + (a_2 - b_2)^2 + ... + (a_n - b_n)^2}$$

## Feature normalization

| Length | Weight | Label |
|---|---|---|
| 4 | 4 | Apple |
| 7 | 5 | Apple |
| 5 | 8 | Banana |

Which of the two examples are closest to the first?

| Length | Weight | Label |
|---|---|---|
| 40 | 4 | Apple |
| 70 | 5 | Apple |
| 50 | 8 | Banana |

$$D(a,b) = \sqrt{(a_1-b_1)^2 + (a_2-b_2)^2 + ... + (a_n-b_n)^2}$$

## Feature normalization

| Length | Weight | Label |
|---|---|---|
| 4 | 4 | Apple |
| 7 | 5 | Apple |
| 5 | 8 | Banana |

$D = \sqrt{(7-4)^2+(5-4)^2} = \sqrt{10}$
$D = \sqrt{(5-4)^2+(8-4)^2} = \sqrt{17}$

| Length | Weight | Label |
|---|---|---|
| 40 | 4 | Apple |
| 70 | 5 | Apple |
| 50 | 8 | Banana |

$D = \sqrt{(70-40)^2+(5-4)^2} = \sqrt{901}$
$D = \sqrt{(70-50)^2+(8-4)^2} = \sqrt{416}$

$$D(a,b) = \sqrt{(a_1-b_1)^2 + (a_2-b_2)^2 + ... + (a_n-b_n)^2}$$

## Feature normalization

| Length | Weight | Color | Label |
|---|---|---|---|
| 4 | 4 | 0 | Apple |
| 5 | 5 | 1 | Apple |
| 7 | 6 | 1 | Banana |
| 4 | 3 | 0 | Apple |
| 6 | 7 | 1 | Banana |
| 5 | 8 | 1 | Banana |
| 5 | 6 | 1 | Apple |

| Length | Weight | Color | Label |
|---|---|---|---|
| 40 | 4 | 0 | Apple |
| 50 | 5 | 1 | Apple |
| 70 | 6 | 1 | Banana |
| 40 | 3 | 0 | Apple |
| 60 | 7 | 1 | Banana |
| 50 | 8 | 1 | Banana |
| 50 | 6 | 1 | Apple |

perceptron: NO!
The classification and weight update are based on the magnitude of the feature value

## Geometric view of perceptron update

for each $w_i$:
$w_i = w_i + f_i *$label

Geometrically, the perceptron update rule is equivalent to "adding" the weight vector and the feature vector

example weights

7

## Geometric view of perceptron update

for each $w_i$:
$$w_i = w_i + f_i*label$$

Geometrically, the perceptron update rule is equivalent to "adding" the weight vector and the feature vector

new weights    example

weights

## Geometric view of perceptron update

If the features dimensions differ in scale, it can bias the update

example    weights

example    weights

same f1 value, but larger f2

## Geometric view of perceptron update

If the features dimensions differ in scale, it can bias the update

new weights    example

weights

new weights    example

weights

- different separating hyperplanes
- the larger dimension becomes much more important

## Feature normalization

| Length | Weight | Color | Label |
|--------|--------|-------|-------|
| 4 | 4 | 0 | Apple |
| 5 | 5 | 1 | Apple |
| 7 | 6 | 1 | Banana |
| 4 | 3 | 0 | Apple |
| 6 | 7 | 1 | Banana |
| 5 | 8 | 1 | Banana |
| 5 | 6 | 1 | Apple |

| Length | Weight | Color | Label |
|--------|--------|-------|-------|
| 40 | 4 | 0 | Apple |
| 50 | 5 | 1 | Apple |
| 70 | 6 | 1 | Banana |
| 40 | 3 | 0 | Apple |
| 60 | 7 | 1 | Banana |
| 50 | 8 | 1 | Banana |
| 50 | 6 | 1 | Apple |

How do we fix this?

## Feature normalization

| Length | Weight | Color | Label |
|---|---|---|---|
| 40 | 4 | 0 | Apple |
| 50 | 5 | 1 | Apple |
| 70 | 6 | 1 | Banana |
| 40 | 3 | 0 | Apple |
| 60 | 7 | 1 | Banana |
| 50 | 8 | 1 | Banana |
| 50 | 6 | 1 | Apple |

Modify all values for a given feature

## Normalize each feature

For each feature (over all examples):

Center: adjust the values so that the mean of that feature is 0. How do we do this?

## Normalize each feature

For each feature (over all examples):

Center: adjust the values so that the mean of that feature is 0: subtract the mean from all values

Rescale/adjust feature values to avoid magnitude bias. Ideas?

## Normalize each feature

For each feature (over all examples):

Center: adjust the values so that the mean of that feature is 0: subtract the mean from all values

Rescale/adjust feature values to avoid magnitude bias:
- Variance scaling: divide each value by the std dev
- Absolute scaling: divide each value by the largest value

Pros/cons of either scaling technique?

## So far…

1. Throw out outlier examples
2. Remove noisy features
3. Pick "good" features
4. Normalize feature values
   1. center data
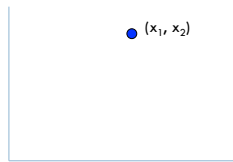   2. scale data (either variance or absolute)

## Example normalization

| Length | Weight | Color | Label |
|--------|--------|-------|-------|
| 4 | 4 | 0 | Apple |
| 5 | 5 | 1 | Apple |
| 7 | 6 | 1 | Banana |
| 4 | 3 | 0 | Apple |
| 6 | 7 | 1 | Banana |
| 5 | 8 | 1 | Banana |
| 5 | 6 | 1 | Apple |

| Length | Weight | Color | Label |
|--------|--------|-------|-------|
| 4 | 4 | 0 | Apple |
| 5 | 5 | 1 | Apple |
| 70 | 60 | 1 | Banana |
| 4 | 3 | 0 | Apple |
| 6 | 7 | 1 | Banana |
| 5 | 8 | 1 | Banana |
| 5 | 6 | 1 | Apple |

Any problem with this?
Solutions?

## Example length normalization

Make all examples roughly the same scale, e.g. make all have length = 1

What is the length of this example/vector?

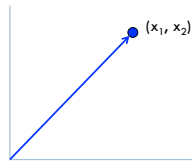$(x_1, x_2)$

## Example length normalization

Make all examples roughly the same scale, e.g. make all have length = 1

What is the length of this example/vector?

$(x_1, x_2)$

$$length(x) = \|x\| = \sqrt{x_1^2 + x_2^2}$$

## Example length normalization

Make all examples roughly the same scale, e.g. make all have length = 1

What is the length of this example/vector?

$(x_1, x_2)$

$$length(x) = \|x\| = \sqrt{x_1^2 + x_2^2 + ... + x_n^2}$$

## Example length normalization

Make all examples have length = 1

Divide each feature value by ||x||

- Prevents a single example from being too impactful
- Equivalent to projecting each example onto a unit sphere

$$length(x) = \|x\| = \sqrt{x_1^2 + x_2^2 + ... + x_n^2}$$

## So far…

1. Throw out outlier examples
2. Remove noisy features
3. Pick "good" features
4. Normalize feature values
   1. center data
   2. scale data (either variance or absolute)
5. Normalize example length
6. Finally, train your model!

## What about testing?

training data
(labeled examples)

pre-process data

learn

model/ classifier

"better" training data

11

## What about testing?

test data

*classify*

model/
classifier

prediction

## What about testing?

test data

*pre-process data*

*classify*

model/
classifier

prediction

**How do we preprocess the test data?**

## Test data preprocessing

1. Throw out outlier examples
2. Remove noisy features
3. Pick "good" features
4. Normalize feature values
   1. center data
   2. scale data (either variance or absolute)
5. Normalize example length

**Which of these do we need to do on test data?
Any issues?**

## Test data preprocessing

1. ~~Throw out outlier examples~~
2. Remove irrelevant/noisy features        Remove/pick same features
3. Pick "good" features
4. Normalize feature values                Do these
   1. center data
   2. scale data (either variance or absolute)    Do this
5. Normalize example length

**Whatever you do on training, you have to do the
EXACT same on testing!**

## Normalizing test data

For each feature (over all examples):

Center:  adjust the values so that the mean of that feature is 0: subtract the mean from all values

Rescale/adjust feature values to avoid magnitude bias:
- Variance scaling: divide each value by the std dev
- Absolute scaling: divide each value by the largest value

What values do we use when normalizing testing data?

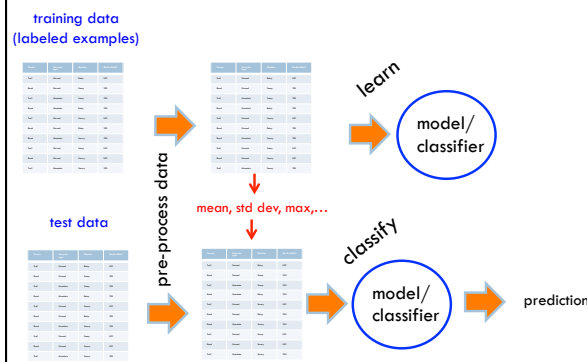## Normalizing test data

For each feature (over all examples):

Center:  adjust the values so that the mean of that feature is 0: subtract the mean from all values

Rescale/adjust feature values to avoid magnitude bias:
- Variance scaling: divide each value by the std dev
- Absolute scaling: divide each value by the largest value

Save these from training normalization!

## Normalizing test data

training data
(labeled examples)

learn

model/
classifier

pre-process data

mean, std dev, max,…

test data

classify

model/
classifier

prediction

## Features pre-processing summary

Many techniques for preprocessing data

Which will work well will depend on the data and the classifier

Try them out and evaluate how they affect performance on dev data

Make sure to do **exact same** pre-processing on train and test

1. Throw out outlier examples
2. Remove noisy features
3. Pick "good" features
4. Normalize feature values
   1. center data
   2. scale data (either variance or absolute)
5. Normalize example length