

# Machine Learning CS726 Fall '11

Angjoo Kanazawa

November 30, 2011

## 1 September 1st: Class one

### 1.1 Logistics

- Uname: cmisc726 Pass: generalize
- Programming projects **27%**: 3 total, after 48hrs 50% down, teams. (Python with numpy).
- Written hw **18%**: 13 of them, one per week out of 3 scales {0, 0.5, 1}. Individual, can't belate.
- Midterm **25%**:
- Final "practical" exam **25%**: Canned or your choice/teams, presentations during the final slot
- Class/Piazza participation **5%**

Time expectation

- 3 hrs in class
- 2 hrs reading
- 2 hrs on written asg
- 2 hrs on programming projects
- Don't ask questions that have already been answered

Do now:

- **Do hw00**: Due 6 Sept, next tuesday (Submit in .pdf only using *handin*)
- **Do the first reading**: pass protected, read the web page.
- **Sign up**: Subscribe to the Piazza group

## 1.2 Contents

### Classification

- When there's an ambiguity in the data, which function are you going to learn? = **Inductive Bias** (i.e. bird vs non-bird, fly vs land, background focus vs bg blurred)
- Which is the more reasonable distinction to make?
- The primary way different learning algorithms vary. There's no right answer. Does your inductive bias match your problem?
- 3 primary components: labels, features, model

### Ingredients for classification

1. Feature representation
  - not typically a focus of machine learning
  - Considered as “problem specific”, but very problematic if bad.
2. Training data: labeled examples
  - Often expensive
  - Sometimes available for free
3. Model
  - No single learning algorithm is always good (“no free lunch” theorem)
  - Requires some control over **generalization**

**Regression** is like classification, except that the *labels are real values*.  
**Structured Prediction**

## 2 September 6th: Class 2

Answer hw questions in few sentences.

### 2.1 Intro cont

- 95% of things in ML are binary classification. But many ML problems reduce to binary classification. We don't have to develop solutions for different problems
- **Reinforcement Learning** - unlike classification, regression, and unsupervised learning, RL does not receive examples, but it **experiences** by interacting with the world.
- RL always has *time* as a variable. i.e. chess, robot control. Comes with the "exploration vs exploitation" trade-off. Humans are really bad at exploitation.
- Defined by *state space*, the world, *actions*, and the *reward*.
- Is a natural extension to search etc we won't cover this too much.

Why we need some math:

- **Calc/Linear Algebra**: finding maxima/minima of functions, allows high-dimensional data analysis.
- **Probability**: the study of the outcome of repeated experiments (frequentist), the study of the plausibility of some event (bayesian)
- **Statistics**: the analysis and interpretation of data. Uses probability theory a lot. Stat is very similar to ML. But stat is about analysis/interpretation of data, but ML is about generalizing and using the analysis to make predictions.
  - A lady drinking tea..
  - model vs predictions
  - model fit vs generalization
  - explain the world vs predict the future

### 2.2 Decision Trees

- Asks questions at nodes, edges are possible options.
- Super efficient, the amount of time is the depth of the tree.
- Histogram representation: in the absence of extra information (no info, no question), make a frequency histogram of the end result. i.e. 62% play, 38% no play.
- Hope that you get about 62% right. Assumption: relative frequency of play/no play is indicative of the real distribution of play/no play.
- Try to ask one question at a time that's most beneficial.

- *example: is sunny?* =  $\{ \text{no (play-80\%, noplay-20\%)}, \text{yes (play-33\%, noplay-66\%)} \}$ . If no, I'll guess play and expect to get 80%, if yes, I'll guess no and expect about 2/3rds right.
- On 5/8th of the days it's not sunny, 3/8th of the days is sunny.
- $5/8 \times 80\% + 3/8 \times 66\% \approx 75\%$
- *example: is windy?* =  $\{ \text{yes (p-50\%, np-50\%)}, \text{no (p-75\%, np-25\%)} \}$ . Frequency of windy days is 1/2.
- So the training accuracy under this feature(question):  $1/2(75\%) + 1/2(50\%) = 5/8 = 62\%$
- So, do I want to ask "is windy" or "is sunny"? Is sunny!
- For each of the features, we're going to compute the training accuracy, get the feature that maximizes the accuracy. Then make that the first question on the decision tree.
- After the root of the tree is found, remove that feature, partition the data associated with the correct edge (not sunny days or sunny days) and repeat the step == The algorithm for decision tree:
  1. guess - get the most frequent answer
  2. base case: labels are unambiguous, no more features.
  3. else for all features, compute the score (training accuracy). Find the guy with max score, split the data, remove that guy, recurse on left and right
- A recursive function, so the depth of the recursion is the number of features, if the features are binary (other wise it's the number of possible feature values etc). So this will terminate.
- Eventually you'll ask every question, so we'll hit the second base case only if there is an inconsistency (two examples with exact same features but with different labels)
- **goal:** good future predictions given (label, example) pairs. A probability distribution,  $D$ , is the guy generating the examples of what the problem looks like. Is the God. Good future predictions means minimum test error.
- $P_D(\text{label}, \text{example})$ . If I have learned well, I should be correctly assign labels to examples that  $D$  thinks is likely.
- test error  $= \epsilon = E_{(x,y) \sim D}[f(x) \neq y]$ , where  $f(x)$  is the prediction,  $y$  is the true label.
- $x$  may be reasonable, or not, because it may not fit the type of data for the problem.
- We don't have  $D$ , but we have the training data  $(x, y)$ , examples drawn from  $D$ . Without  $D$ , minimizing  $\epsilon$  is too hard.

- **goal'**: good training predictions, i.e. minimize training error.
- $\hat{\epsilon} = 1/N \sum_{(x,y) \in \text{train}} [f(x) \neq y]$
- Can the decision tree algorithm minimize  $\hat{\epsilon}$ ? if there is no ambiguity in the training data, we can ask all the question (worst case), then we pinned down all examples (and there are no diff labels), and we can do this. If there is an ambiguity/inconsistency, it's not possible to always get the 0 training error.
- The best we can hope for is that whatever this minimum error is, it will achieve that. For this goal', the algorithm is optimum.
- Does solving goal' help us help solve goal? if the assumption that test comes from the same  $D$ , we can.
- Doesn't work when test data is small, and is not representative of the  $D$ .
- Say there are spurious correlations between features and the labels in the training data. (looks useful in training, but not in real test). The decision tree will pick these features, it will make  $\hat{\epsilon}$  low, but  $\epsilon$  will be high. This is **overfitting**.
- Solutions: constrain the height of the tree (you can only ask x questions and hope that it won't ask questions about this spurious feature), get more data.

## 3 September 8th Class 3

Project 1 is online now. Implement basic classifiers. Start early! Due 27th September.

### 3.1 HW1

1. Memorizing doesn't generalize, and overfits.
2. Look at 3.2
3. Goal: minimize  $\epsilon = \mathbf{E}_{(x,y)} [f(x) \neq y]$ . We can calculate  $\hat{\epsilon} = \frac{1}{N} \sum_n [f(x_n) \neq y_n]$ . If overfit, we do too well on the training data so  $\hat{\epsilon}$  will be tiny but not the  $\epsilon$ .
4. The algorithm should never look at the data because it'll overfit and we won't be able to say anything about the model's performance. We shouldn't look at it because we'll get a sense of what features will be important and fix the model accordingly.

### 3.2 Decision Tree

**Claim:** shallow trees are less prone to overfitting.

Example: two coins, A, B. Result of coin A will be the feature,  $x$ , that of coin B will be the label,  $y$  (they aren't correlated). Got  $(1, 1), (0, 0), (0, 0), (0, 1)$ . Looks like the feature is very useful, but not. The hope is that these accidental correlation will go down the tree, and the actual features that work go up the tree.

The depth is a *hyper parameter*. The goal of using hyperparameters is to reduce overfitting. If we treat it as a parameter, it'll just use max depth and lower  $\hat{\epsilon}$ . *Parameters* are what you estimate on the training data.

### 3.3 Geometry

We represent our inputs as vectors in high dimensional space.  $y = \pm 1$ ,  $x = (x_1, x_2, \dots, x_n) \in \mathbf{R}^D$

**K-NN:** Look at  $K$  nearest neighbors, label the new point as the majority/mean of the neighbors. If inf-NN, that's just taking the majority. This is underfitting. 1-NN would be overfitting. Choose odd  $K$  so that there won't be any ties.

**K-mediods Clustering.**

1. Pick  $K$  datapoints to be representatives
2.  $\forall n$ , put  $n$  in the closest cluster (closest  $K$  datapoints)
3. Choose a new prerepresentative that minimizes the average dist to all cluster members
4. Go back to 2

**K-means Clustering.**

1. Pick  $K$  datapoints randomly  $\mu_1, \dots, \mu_k$

2.  $\forall n, z_n = \arg \min_k |x_n - \mu_k|$ .
3.  $\forall k, \mu_k = \frac{1}{n_k} \sum_{n: z_n=k} [x_n]$  (Make a new mean prerepresentative who is an average of all cluster members.)
4. Go back to 2

Does  $K$ -means always converge? Measure the quality of the solution by the average distance from each point to its mean (score). (Convergence as in this measure of cluster quality won't change.)

Observation 1: Everytime step 2 or 3 is executed, the score goes down or stays the same.

Observation 2: If there are  $n$  members and  $k$  clusters, the possible clustering assignment is finite,  $k^n$ .

Does  $K$ -means always find the optimal clustering?

No! Not necessarily. The initial point assignment is random and it changes the result. The practical method is to do multiple initialization and pick the smallest score.

How many iteration does it take to converge?

In theory it takes a time that's exponential to  $k$ . In practice, it converges in like 4 iterations. Smooth analysis of why this converges quickly. Only certain initial points take a long time. If it's taking a loong time, if you perturb the initial points and move it away (like how simplex is really fast in real life).

$K = n$  ( $\hat{\epsilon} = 0$ ) or  $K = 1$  gives us no information about the data. As an implementation note, some cluster centers will disappear and never come back. So be able to handle that situation.

### 3.4 Curse of High Dimentionality

as we make the dimention higher, the distance between points get smaller. If every point is about the same distance from everything else, the  $k$ -nn are just random points. Picking the majority of those random points, is just another random guessing. As  $D$  increases, the distance between points concentrates to a point.

```
D = 100; hist(flatten(XtoD(rand(1000, D )) / sqrt(D)));
%XtoD takes the distance between points
```

Seems like  $K$ -nn shouldn't work, but because the  $D$  is uniform. It does work in real life, because in real life the actual data is correlated.

## 4 September 15th Class 4

### 4.1 HW2

2. Decision boundary for a one nearest neighbor classifiers on two datapoints is a linear plane perpendicular to both points.
3. Clustering as in given labeled data  $(x_n, y_n)$ ,  $\forall n \in N$  for each label run  $K$ -means on subset of data with this label, then run  $KNN$  using  $\mu_i$ s. A lot faster on  $KNN$  test time but more time on training, in a sense clustering throws away bunch of data. This could avoid overfitting. (Assuming that we're working on relevant feature set).

4. in 2D space, the dot product,  $\langle \bar{w}, \bar{x} \rangle$ , projects all the points,  $X$ , onto  $\bar{w}$  and makes it 1D.

If  $y = \text{sign}(w \cdot x)$ , is the decision rule, then the boundary  $B$  is  $\{x : w \cdot x = 0\}$ . The disadvantage is that it assumes the decision boundary lies on the origin. This won't work if all the points are positive. So now, modify  $y$  to  $y = \text{sign}(w \cdot x + b)$ , s.t.  $B = \{x : w \cdot x = -b\}$ . Bias is always  $-b$ , as  $b \rightarrow -\infty$ , nothing can be classified as positive.

Claim: using bias is the same as not using the bias but adding a new feature. i.e.  $x$  or  $\langle 1, x \rangle = \tilde{x}$ , where  $x$  uses  $w$  (to get  $w \cot x + b$ ) and  $\tilde{x}$  uses  $\tilde{w} = \langle bw \rangle$  so it's the same thing i.e.  $w \cot x + b = \tilde{w} \cdot \tilde{x} = \langle b, w \rangle \cdot \langle 1, x \rangle$ . So we can always write as if there's not bias.

### 4.2 Perceptron

General algorithm

```
for each (x,y)
  //if error, do an update
  if y(w\cdot x + b) \le 0: or if sign(y) == sign(w\cdot x + b):
    w = w + yx
    b = b + y
```

Usually people scale the update with some constant  $\gamma < 1$ , the learning rate, so that it won't be so influenced by  $x$ . (each single error);

*Question:* If you add  $\gamma$  does it change?? Doesn't really change anything, because after one update, the scale of  $w$  changes, but the direction is the same. With a bias, it'll shift it a lot of a little. After any number of updates, without using the learning rate, you'll get  $(w, b)$  in the end but with  $\gamma$ , all you'll get is  $(\gamma w, \gamma b)$ ., the decision boundary will have the same direction.

In the case of perceptrons, adding a  $\gamma$  makes no difference (assuming we have infinite floating point). In the case that  $w$  is a unit vector,  $b$  shifts the decision boundary  $b$  units, if  $w$  is not unit, it shifts  $b/\|w\|$  units. So if  $\|w\| \gg 1$ ,  $b$  won't really make that much difference.

More important point: the more you run, the less sensitive the algorithm will be for any update. Early on, the decision boundary moves all over the place, but as the number of epoches increases, the relative magnitude of the update is smaller than that of the beginning updates. I.e. even if it won't converge, it'll slow down in the end.



## 5 September 22nd Class 6

### 5.1 Linear Seperability

Perceptrons can't do XOR, because 2D XOR is not linearly seperable. But if you add a 3rd dimension, where  $-$ 's have 0 and  $+$ 's have 1. Then you can have a hyperplane that seperates them.

You can add all possible combinations linearly, this is what perceptrons do. But that's computationally expensive and makes the feature very high dimension. This is prone to overfitting, and you need more examples.

DT can combine few features but in a complicated way.

DT might fail on XOR because of all ties, features very slightly correlated with the label..

### 5.2 Evaluation/ROC

Given a ratio of dataset where 1% of them are positive and 99% is negative, the classifier will say all of them are negative and achieve 99% accuracy. Example for cancer detection, false negative is worse than false positive. So there are different costs associated with different classification. Perhaps my goal is to order the data to see what's interesting to look at.

In biology, people look at *Sensitivity vs Specificity*, we call it *Precision vs Recall*. Precision says of all the  $X$ s that you found, how many of them were actually  $X$ s? Recall asks: of all the  $X$ s that were out there, how many of them did you find? (should've found)

Choose the best threshold and get the harmonic mean, or their **f-measure**:

$$F = \frac{2PR}{P + R}$$

### 5.3 Cross-validation

Get bunch of data, divide it in  $n$  equal size. Say  $n = 10$ , train on  $n_2 \dots n_{10}$ , and test on  $n_1$ , but now, you can train on  $n_1 \dots n_9$ , then test on  $n_{10}$ . You can do this  $n$  times. Then you'll get performance over all of the data I have, not just a single slice.

Also you do this to tune your hyperparameters. For each hyperparameter setting, you get the score of each, you want the hyperparameter that does the best on average over all slices.

Say you run this on DT and get that depth = 17 is the best, so you have  $n$  trained classifier, get the one that does the best or, retrain a new classifier on all  $n - 1$  training slices. Usually learning a new classifier is a better idea, but it might be slightly underfit, because you're regularizing (with depth=17), but using more data.

4 choices: Number of slices  $n = 10$  if you're normal, you do  $n = 5$  if you're running out of time, you do  $n = 2$ , you'll make sure your two classifiers aren't trained on the same data points, one more is train on everything except 1, *leave-one-out*, great for *knn* but too expensive for others. Leave-one-out is the least succceptable to underfitting.

## 5.4 Statistical Significance

Try to answer the question, “given this new method, does it actually work?”. You should use the most common thing in your field. Two examples we cover is *Paired T-test* and *Jackknifing/Bootstrapping*.

Given DT and Perceptron’s performance, say perceptron always does better. “If I were to do a new experiment, with what probability would Perceptron beats DT”

This is what T-test does, fit gaussian to DT, fit gaussian to Perceptron, and if the overlap is tiny, then there’s small probability that DT will beat perceptron, this gives us a  $p$ -value.

Weakness is that there is a normal distribution assumption. LLN argument, it should be kind of, but... not really because accuracy lives between  $[0, 1]$ , normal lives between  $[-\infty, \infty]$ .

So some people say that you should do binomial tests, but ppl still do  $p$ -value.

**Jackknifing/Bootstrapping** (The word “Bootstrapping” is used a lot everywhere). This is suited where your evaluation method is holistic over the their data. Like  $F$ -measure.. Way around when T-test can’t really be applied. Doesn’t make normality assumption, it’s a *non-parametric* test.

If you have a big dataset, it’s easier to separate the two gaussian curves (because they’ll have more peaks), so you can get more examples to get more statistical significance.

## 6 September 27th Class 7

### 6.1 HW3

1. About non-linearly separable data and perceptron: Do the add one feature augmentation trick, why is this linearly separable now? For example,  $\mathbf{w}$   $D$  zeros and all indicators as corresponding  $y_n$  the label.

Does this affect generalization? The above construction will not work for test data at all. Generally it overfits, because it's like cheating after the first couple of passes when you start using the labels  $y_n$  as the components of your new augmented  $\mathbf{w}$

How long does it take to converge on this augmented data? Previously,  $\|\mathbf{x}\| \leq R$ , but now, we have  $\|\langle \mathbf{x}, 0, 0, \dots, 1, \dots \rangle\|$ ? Square it:

$$\begin{aligned}\|\langle \mathbf{x}, 0, 0, \dots, 1, \dots \rangle\|^2 &= \|\langle \mathbf{x}, \dots, 1, \dots \rangle \cdot \langle \mathbf{x}, \dots, 1, \dots \rangle\| \\ &= \mathbf{x} \cdot \mathbf{x} + 1 \\ &= \|\mathbf{x}\|^2 + 1\end{aligned}$$

So  $\hat{R} = \sqrt{R^2 + 1}$ , but now  $\gamma$  changes also. Using the  $\mathbf{w}$  above, the margin for this  $\mathbf{w}$  is  $\min_{(x,y) \in D} \frac{y}{\|\mathbf{w}\|} \mathbf{w} \cdot \mathbf{x} \|\mathbf{w}\| = \sqrt{N}$ , and  $\max \mathbf{w} \cdot \mathbf{x}$  is 1 so the new margin is at least  $\frac{1}{\sqrt{N}}$

2. Centering, Variance scaling vs DT, KNN, perceptron:

- **KNN** is not sensitive to centering but it is sensitive to variance scaling.
- **DT** is not sensitive to centering or scaling because the cut off is still the cut off.
- **Perceptron** is sensitive to centering and variance scaling, because if the feature value is huge, it'll move  $\mathbf{w}$  a lot to one direction, it'll take more time because  $R$  is huge.

### 6.2 Blue Box

- In the realm of *scale* of features, variance pruning might not work if done prior to scaling because it might be a good feature in the correct scale.
- *How do you get a confidence out of a DT or KNN?* For DT you can calculate the probability. Each node is a conditional probability. For KNN, we can calculate the average within  $k$  neighbors and get a percentage (with  $K = 6$ , 3 might be  $VE$ , 3 might be  $-VE$ , so 50% chance)

### 6.3 Multiclass Classification

**Reduction-based ML:** Is multiclass harder than binary? If I could do binary, then I'll use that to solve multiclass.

*Error-bounds:* If I can achieve  $\epsilon$  error in binary classification, we can have  $\leq f(\epsilon)$  error on multiclass classification.

The big assumption is that we can achieve  $\epsilon$  error.

**Definition:** Binary Classification.

- Assume some distribution  $\mathbf{D}$  of  $(x, y)$  pairs
- get sample  $D = (x_1, y_1), \dots, (x_n, y_n)$  drawn from  $\mathbf{D}$
- The goal is to compute  $f : X \rightarrow \{-1, +1\}$  s.t.  $\mathbf{E}_{(x,y) \sim D}[f(x) \neq y]$  is small.

**Weighted binary classification:** Weight the negative sample less and positive samples more, so that the classifier has to put more effort to classify positive examples correctly. Now the goal is different: it is to minimize (arbitrarily assuming that positive error costs more here):

$$\mathbf{E}_{(x,y) \sim D} \begin{cases} \alpha[f(x) \neq 1] & \text{if } y=1 \\ [f(x) \neq -1] & \text{if } y=-1 \end{cases}$$

For some  $\alpha > 0 \in \mathbf{R}$ . Above can be re-written as  $\mathbf{E}_{(x,y) \sim D}[\alpha^{y=1}[f(x) \neq y]]$   
Algorithm to use: **Subsampling**. We end up with a “balanced” data. If the cost is  $\alpha = 10$ , I want 10th as many positive examples. This is optimal for us. What we want

1. given  $\mathbf{D}^{WBC}$ , we want to reduce it to  $\mathbf{D}^{BC}$
2. learn a function  $f$  on the binary classification problem.
3. given a test point  $\hat{x}$ , use  $f$  to solve the weighted binary classification problem.

The algorithm

- 1a Retain all positive examples
- 1b For each negative examples, sample  $1/\alpha$  of them.
- 2  $\hat{y} = f(\hat{x})$  (i.e. 100 exmples, 98 -ve, 2ve,  $\alpha = 2$ , randomly select 49 negative examples.

We want  $\mathbf{E}_{(x,y) \sim D^{WBC}}[\alpha^{y=1}[f(x) \neq y]]$ , from weighted  $D$ .

$$\begin{aligned} \mathbf{E}_{(x,y) \sim D^{WBC}}[\alpha^{y=1}[f(x) \neq y]] &= \sum_{(x,y)} \mathbf{D}^{WBC}(x, y) [\alpha^{y=1}[f(x) \neq y]] \\ &= \sum_x \begin{cases} \mathbf{D}^{WBC}(x, +1) [\alpha[f(x) \neq 1]] \\ \mathbf{D}^{WBC}(x, -1) [f(x) \neq -1] \end{cases} \\ &= \sum_x \begin{cases} \mathbf{D}^{BC}(x, +1) [\alpha[f(x) \neq 1]] \\ \alpha \mathbf{D}^{BC}(x, -1) [f(x) \neq -1] \end{cases} \\ &= \alpha \sum_x \begin{cases} \mathbf{D}^{BC}(x, +1) [[f(x) \neq 1]] \\ \mathbf{D}^{BC}(x, -1) [f(x) \neq -1] \end{cases} \\ &= \alpha \epsilon \end{aligned}$$

Because we downselected negative examples, so to go from  $D^{WBC}$  to  $D^{BC}$ , we need to multiply it with  $\alpha$ .

Can we do better than  $\alpha\epsilon$  (on the weighted problem)? No, the error grows linearly with the size of  $\alpha$ . You could imagine to do slightly better because of the sampling more of +ve but something around there.

Sampling seems wasteful. SO instead, we can keep all negative examples and put  $\alpha$  times more positive examples (replacement of). This is **over-sampling**. Both gives us the same bounds, but  $\epsilon$  can be different.  $\epsilon$  is the error rate of the BC classifier. The big assumption is that we can achieve  $\epsilon$ , it's easier to achieve low  $\epsilon$  in oversampling than undersampling. *Even though the bound for both sampling is  $\alpha\epsilon$ ,  $\epsilon$  is smaller from oversampling, generally.*

## 6.4 Turning MC to BC

1. **OVA** (one-vs-all): with  $k$  classifiers  $f_1 \dots f_k$ , where  $f_1$  separates  $y = 1$  from  $y \neq 1$  and  $f_k$  separates  $y = k$  from  $y \neq k$ . Two bad things can happen: no one says yes, more than one says yes. Just pick at random for each, but these are the cases when the error occurs. Not great because of this. In practice, you pick the one with the highest confidence. Called brittle, because it takes only one error will shoot up the probability of error on MC. Any errors can cause grave harm, if each error is  $\epsilon$ , and you have  $k$  chances, it's pretty bad.
2. **AVA** (all-vs-all):  $\binom{k}{2}$  classifiers where  $f_{ij}$  separates class  $i$  from  $j$ .

## 7 September 29th Class 8

### 7.1 HW 4

1. Task for squared error regression: Given: input space  $\mathcal{X}$  and an unknown distribution  $\mathcal{D}$  over  $\mathcal{X} \times \mathbf{R}$ , Compute  $f : X \rightarrow \mathbf{R}$ , minimize  $\mathbf{E}_{(x,y) \sim \mathcal{D}}[f(x) \neq y]$
2. *To what types of error do these theorems apply?* All error in respect to data taken from the same distribution  $\mathcal{D}$ .
3. OVA vs AVA. If training for each classifier is  $O(N)$ , OVA:  $O(NK)$ , AVA:  $O(K^2N/K) = O(NK)$ , if training for each classifier is  $O(N^2)$ , OVA:  $O(N^2K)$ , AVA:  $O(K^2(N/K)^2) = O(N^2)$ , so AVA is better! so the worst it is to train, the better and better AVA gets.
4. A ranking preference function  $\omega$  that penalizes mispredictions linearly up to  $K$ :

$$\omega(x, y) = \begin{cases} \min\{k, |i - j|\} & \text{if } i \neq j \\ 0 & \text{else} \end{cases}$$

### 7.2 Continue on Multiclass classification

There's OVA, AVA (the default), and **Tree-based** classification. Idea is to divide and conquer, split all the classes in two, then keep on splitting. At the leaf we make decisions. So if we have  $K$  classes, we have  $K - 1$  nodes and so we train  $K - 1$  classifiers. This divides the data, so we beat OVA's  $O(K\epsilon^8)$ ,  $K$  factor.

As soon as we make an error, we're done. The chance of making an error is  $\epsilon^b$ , the number of times you can make an error is how many nodes you go down (pf union bound:  $P(A_1 \cup A_2 \cup \dots \cup A_n) \leq \sum P(A_i)$ , so at most the error caused by tree based is

$$\epsilon^{MC} \leq O(\epsilon^b \log(K))$$

If the number of features is not a power of 2, it might not be a full binary tree with height  $\log(K)$ , so more of the emphasis on  $\leq$ .

But here we can't specify which order/the layout of the tree. Certain layouts could make making a binary classifier easier. I.e. if data were clustered in  $1 : (0, 0), 2 : (0, 1), 3 : (1, 0), 4 : (1, 1)$ , shouldn't split classes  $\{1, 4\}$  vs  $\{2, 3\}$ .

### 7.3 Ranking

**Bipartite Ranking:** Some are labeled good and some are bad. The job is to rank the good, relevant ones before the bad ones. Given nodes  $a, b$ , the labels is 1 if  $a$  is good and comes before  $b$  and should,  $-1$  if  $b$  comes before  $a$ .

We use  $\omega(x, y)$  as the preference function just like binary classification's loss function. Where  $x$  is the true position,  $y$  is my ranking.

**Objective** Given an input space  $\mathcal{X}$ , and an unknown distribution  $\mathbf{D}$  over  $\mathcal{X}$ , compute: a function  $f : \mathcal{X} \rightarrow \Sigma_M$  minimizing

$$\mathbf{E}[\sum_{\mathbf{u} \neq \mathbf{v}} [\mathbf{sig}_{\mathbf{u}} < \mathbf{sig}_{\mathbf{v}}][\hat{\sigma}_{\mathbf{u}} < \hat{\sigma}_{\mathbf{v}}]\omega(\sigma_{\mathbf{u}}, \sigma_{\mathbf{v}})]$$

Where  $\hat{sig} = f(x)$ . Basically did we swap  $u$  and  $v$ , if so, pay  $\omega(u, v)$ .

$$f(x_{uv}) \begin{cases} +1 & \text{if } u < v \\ -1 & \text{if } v < u \end{cases}$$

$f(x_{uv})$  answers “is  $u <_f v$ ?” Now use this as a comparison function and use any sorting algorithm, like quicksort. So it’s just a quicksort with a minor tweak:

Consider  $f(x_{uv})$  as  $P(u < v)$ , so in quicksort, pick a pivot, is this other number greater than the pivot or not? We answer this by running  $f$ , and  $f$  will tell us if it’s less. So it’s like a randomized quicksort.

So the weighted problem looks like:, given  $(y = \{+1, -1\}, x, w)$ , (where  $w = \omega$ ,  $x = x_{uv}$ ) we want to minimize

$$\epsilon^w = \mathbf{E}_{(x,y,w)} \mathcal{D}[w(y \neq f(x))]$$

How would you solve this? Over sampling, add 2 more  $(x, y)$  if  $w = 2$ . Undersampling, we want to keep if  $w$  is high, so throw out data with probability  $1/w$ . So if  $w$  large, we most likely never throw out the data. And we have  $\epsilon^w \leq \epsilon^b[\mathbf{E}_w[w]]$  is the best we can hope for.

## 7.4 Collective Classification

Given a graph, where node is a feature, like a website, and if we were to label each node as  $+1$  not offensive, or  $-1$  offensive, the idea is that some information in the graph structure will tell us about the label. i.e. offensive websites tend to link to offensive websites.

So train a classifier with a bigger feature, where you include information about the node itself plus the number of positive neighbors and negative neighbors. Problem with this is that.....TODO!

The **Stacking algorithm**: Forget that we have a graph structure. First try to classify each page based on its individual features alone. So think:  $f_0 : \mathcal{X} \rightarrow \{+1, -1\}$ , and now I have guesses for all. Now we can train a second round classifier, that takes an input from the pages themselves and the neighbor information from the previous classifier (the new added info on the features are given by  $f_0$ ). So train  $f_1 : \mathcal{X} + f_0$ ’s output  $\rightarrow \{+1, -1\}$ . Now, the hope is that  $f_1$  will get everything correct. Now train  $f_2$  that takes in  $\mathcal{X}$  and  $f_0$  and  $f_1$ ’s output. etc.

Problem: it never gets the true label of the neighbors, but only the predicted data. This might be beneficial because at test time we never get to see the true label of the neighbors. But the issue is overfitting:  $f_0$  is trained and tested on the original data, so chances are it’ll do really good on the input data. So input to  $f_1$  will be unrealistically close to perfect. In the real world  $f_0$  will do much worse.

Concrete example: task for each 3D point, guess if it corresponds to a tree or not, don’t get info but angle and distance. It’s nearly impossible to classify the point just by using its features, you really need the entire graph structure. So at train time we do well, but ... not in test..

## 8 Oct 4th Class 8

### 8.1 Linear Classifiers

Linear decision boundary:  $\text{sign}(\mathbf{w} \cdot \mathbf{x} + b)$ . If  $b = 0$ , the plane passes through the origin. Computing the distance of this plane from any point  $\mathbf{z}$ , is  $\frac{\mathbf{w} \cdot \mathbf{z} + b}{\|\mathbf{w}\|}$ . So from the origin is  $\frac{b}{\|\mathbf{w}\|}$ . If margin  $\geq 0$ , loss is 1, if margin  $< 0$ , loss is 0.

Indicator function  $\infty(\cdot) : \begin{cases} 1 & \text{if argument is true} \\ 0 & \end{cases}$  Using 0-1 loss finding the decision boundary is NP-hard.

### 8.2 Regularization

We're concerned about test error, not training error. Just minimizing the training error will overfit. So Regularized Learning (Tikhonov ( $\lambda R(\mathbf{w})$ ), Ivanov (minimize the loss with constraints that  $R(\mathbf{w}) \leq \delta$ ), Morozov (minimize the regularizer with constraints on the loss function  $le\delta$ )) All three do the same thing, but most common is Tikhonov:

$$\mathbf{w} = \underset{\mathbf{w}, b}{\operatorname{argmin}} \frac{1}{N} \sum_n \infty(y_n(\mathbf{w} \cdot \mathbf{x}_n + b) < 0) + \lambda R(\mathbf{w})$$

This balances empirical loss and complexity of the classifier - prefers simpler one.  $R(\mathbf{w})$  is the regularizer for linear hyperplanes. We want both the loss and regularizer to be convex.

### 8.3 Convexity

**Definition:** A set  $S$  is convex if  $\forall x, y \in S, \alpha x + (1 - \alpha)y \in S$  where  $\alpha \in [0, 1]$ . The line segment joining  $x$  and  $y$  is contained in  $S$ .  $S = \{x : x^2 > 2\}$  is non-convex because  $x \geq \pm\sqrt{2}$ , inside it's empty  $S = \{x : \mathbf{x}^T \mathbf{x} < 1\}$ . Show in two points norm will always be  $\leq 1$  with  $\alpha$   $S = \{U : U^T U = 1\}$ , orthogonal matrices, is *not* convex. So eigenvalue problems are non-convex.

Operations that preserve convexity:

- Intersection
- Affine function  $f(x) = Ax + b$

**Definition:**  $f$  defined on a *convex* set is a *convex* if

$$f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y) \forall \alpha \in [0, 1]$$

The line is always above the function. Or equivalently,  $f$  is convex if its *epigraph* is a convex set:

$$E = \{(x, \mu) : \mu \in \mathbf{R}, f(x) \leq \mu\}$$

(region(set of points) that lies above the graph of the function  $f$ .)

Also,  $f$  convex  $\Rightarrow -f$  concave.

Checking for convexity: First check if the domain is convex

1. Use definition (chord lies above the function)



2. For differentiable  $f$ , function lies above all tangents i.e.

$$f(y) \geq f(x) + f'(x)(y - x)$$

3. For twice differentiable  $f$ , second derivative is non-negative. (These are positive definite matrices)

If the inequalities are strict, then  $f$  is *strictly convex*. Operations that preserve convexity:

- Positive scaling/addition
- Affine function composition:  $f$  convex  $\rightarrow f(Ax + b)$  convex
- Pointwise maximum:  $f, g$ , convex  $\rightarrow h(x) = \max(f(x), g(x))$  is convex
- $f$  convex,  $g$  convex non-decreasing  $\rightarrow h(x) = g(f(x))$  convex
- $f$  concave,  $g$  convex non-increasing  $\rightarrow h(x) = g(f(x))$  convex

exp good example of strictly convex, non-decreasing Examples:

- $e^x, x, x^2, x^4$  take the second derivati
- Any vector or matrix norm
- $h(x) = \max(|x|, x^2), x \in \mathbf{R}$
- $\exp(f(x))$ , convex, by the 4th property
- $\frac{1}{\log(x)}, x > 0$ .  $\log(x)$  is concave,  $1/x$  is convex non-increasing in  $\mathbf{R}^+$ , so by 5 this is convex.

### 8.3.1 Convex Loss Functions

Hinge loss, exponential loss, and logistic loss

- Hinge loss:  $L(y, \hat{y}) = \max\{0, 1 - y\hat{y}\}$
- Exponential loss:  $L(y, \hat{y}) = \exp(-y\hat{y})$
- Logistic Loss:  $L(y, \hat{y}) = \log_2(1 + \exp(-y\hat{y}))$

Even if margin is good, won't get 0 for exponential and logistic.

## 8.4 Weight Regularization

Why?

- We need our weights to be small i.e. we DON'T want  $\epsilon$  change in input ( $\hat{x} = (x + \epsilon)$ ) to make  $\hat{y}$  change a lot i.e.  $w\epsilon$  will have a lot of influence if  $\|w\| >> 1$ .
- Regularization is generalization, regularized learning is statistically stable and bounds generalization error. Stability here means that differences in  $\epsilon$  is bounded..
- We want regularizer to be convex for computational reasons

### 8.4.1 Norm based Regularizers

$l_p$  norm:

$$\|w\|_p = (\sum_i |w_i|^p)^{1/p}, p \in [1, \infty]$$

if  $p < 1$ , it won't satisfy triangle inequality and convexity. We use  $l_1, l_2$  usually. A unit ball in  $l_p$  is a set  $S = \{w : \|w\|_p \leq 1\}$ , so in  $l_2$  it will be a unit circle, with  $l_1$ , it will be a unit square rotated 90 degrees. As  $p \rightarrow 0$ , the sides approach the axis,  $l_0$  ball looks like a plus sign. As  $p \rightarrow \infty$  the circle approaches square. ( $l_\infty(w) = \max w$ ).

Properties:

- Solution often lies on the singularity (corners) of the ball when constrained enough. Why? (Why  $l_1$  gives us sparsity, because the corners/edges are 4 or constant)
- Rotational Invariance:  $l_2$  is invariant to rotations
- When to use  $l_1$ : less number of samples, lot of irrelevant features
- When to use  $l_2$ : otherwise, leads to good generalization.

Now we our final objective is:

$$L(w, b) = \frac{1}{N} \sum_n l(\hat{y}, y_n) + \frac{\lambda}{2} \|w\|_2^2$$

### 8.4.2 Gradient Descent

Gradient of a function

$$\nabla f(x) = (\frac{\partial f(x)}{\partial x_1}, \frac{\partial f(x)}{\partial x_2}, \dots, \frac{\partial f(x)}{\partial x_n})$$

**Directional derivative** Rate of change of function along a direction  $v$ .  $D_v(f) = \nabla f \cdot v$  The function increases the most along the direction of the gradient. Why? because  $\nabla f \cdot v = \cos(\theta) |\nabla f| |v|$ , when the angle is 0,  $\cos(0) = 1$  and gives us the biggest increase.

Geometric interpretation:

- **Level sets:** for function  $f(x)$ , the level set containing a point  $a$  is a set of points  $S = \{x : f(x) = f(a)\}$ . Function admits same value at all points in the level set  $S$
- Gradient at a point  $y$  is perpendicular to the level set containing  $y$ . i.e. a gradient at point  $(x, y)$  will be perpendicular to the level set going through  $(x, y)$ . (from definition, gradient is always normal to the tangent going through the point)
- i.e.  $f(x, y) = \sqrt{x^2 + y^2}$ , level set is a circle, for  $f(x, y, z) = \sqrt{x^2 + y^2 + z^2}$ , the level set is a sphere.

**Gradient Descent:**

- Move in the negative direction of gradient to minimize a function.

- Steepest gradient descent: To minimize a function  $f$ , start with some initial guess  $x^{(0)}$  with update rule:

$$x^{(t)} = x^{(t-1)} - \eta^{(t)} \nabla f(x^{(t-1)})$$

## 9 October 11th Class 10

### 9.1 Probabilistic Modeling

**Joint:**  $p(x, y)$ , **Conditional:**  $p(y|x)$  (some ppl think this is the discriminative).

#### 9.1.1 Joint Modeling Setup

Assumption  $\mathcal{D}(x, y)$ ,  $f(\hat{x}) = \operatorname{argmax}_y \mathcal{D}(\hat{x}, y)$  (pick the most likely thing). If there's unique  $y$  for every  $x$  with probability 1 there will be no error.

Training Data:  $(x_n, y_n)_n$   $\mathcal{D}$ , we want to estimate  $\hat{p}(x, y) \approx \mathcal{D}$ , where  $f(\hat{x}) = \operatorname{argmax}_y \hat{p}(\hat{x}, y)$

Typical Distributions:

- For labels with  $\{1, -1\}$  (coin flips), when  $p(h) = \beta$ ,  $p(t) = 1 - \beta$ , use *Bernoulli* distribution.
- For labels in the reals, people use *Gaussians*.

People approach joint modeling by factorizing  $p(x, y)$  into  $p(y)p(x|y)$  (chain rule). Backwards of what you think about classification. First decide if there's a bird or not, then draw the picture depending on if i decided if there's a bird or not.

$x$ s are usually feature vectors, so it's sometimes to come up with a good model over  $p(x|y)$ .

Example: label: play or not, features: wind, sunny, temp = many combinations. For each, given  $Y$ , assign probabilities to each combinations. where sum of all combinations is 1. But when the number of features grow, it's just a lot of combinations and there might be not enough data.

**Naive Bayes Assumption** Make smaller tables, condition on play for each feature, so  $p(W|P) = .6$ , then  $p(W|P) = .4$ , and if  $p(W|P) = .2$ , then  $p(W|P) = .8$  Using

$$p(x_1, \dots, x_n|y) = \prod_d p(x_d|x_1, \dots, x_{d-1}, y) \approx \prod_d p(x_d|y)$$

i.e. *Features are independent given the label*.

Note: In the real world it's clear that  $p(\text{temperature}|\text{outlook}) \neq p(\text{temp})$ . But naive bayes is not this bad, because  $p(\text{temp}|\text{outlook}, \text{play}) = p(\text{temp}|\text{play})$ , we are conditioning on the label. Not the best assumption but not totally wrong.

So we have,  $p(x, y) = p(y)p(x|y) \approx p(y) \prod_d p(x_d|y)$

Telling this fictional tale about where the data came from, the closer the tale is the better my model will be.

Probabilistic assumption = inductive bias.

### 9.2 Primary BUilding model in Naive Bayes

Short hand for if  $x = 1$ , probability is  $\theta$ , else  $1 - \theta$ :

$$\theta^{(x=1)}(1 - \theta)^{(x \neq 1)}$$

### 9.2.1 Parameter Estimation

: Goal is given data, estimate  $\mathcal{D}$ . Best way is **Maximum Likelihood**

- Pick parameters  $\theta$  that maximize  $P_\theta(D)$  where that's the likelihood of  $\theta$  (not likelihood of data, probability of data given  $\theta$  ( $P(D|\theta)$ )).

Goal: Maximize  $P((x_1, y_1), \dots, (x_n, y_n) | \theta)$  Make a first big assumption, that  $x$ s are i.i.d. Then we get  $\max_\theta \prod_n p(x_n, y_n | \theta)$ . The data isn't really iid because yesterday it was raining so it probably affects your samples.

Now using naive bayes, assume that we have

$$\begin{aligned} &\approx \max_\theta \prod_n p(y_n | \theta) \prod_d p(x_{n,d} | \theta, y_n) \\ &\approx \max_\theta \prod_n \theta_0^{(y_n=1)} (1 - \theta_0)^{(y_n \neq 1)} \prod_d \theta_{+1,d}^{(y_n=1, x_{n,d}=1)} (1 - \theta_{+1,d})^{(y_n=1, x_{n,d} \neq 1)} \theta_{-1,d}^{(y_n \neq 1, x_{n,d}=1)} (1 - \theta_{-1,d})^{(y_n \neq 1, x_{n,d} \neq 1)} \end{aligned}$$

( $\theta_0$  is how many days i play) Given Play=1, outlook=1 is  $\theta_{+1,out}$ , so play=1, outlook=0 is  $(1 - \theta)_{+1,-out}$ , and so forth for each feature and it's label combination. So  $\theta_{+1,d}^{(y_n=1, x_{n,d}=1)}$  is  $\theta_{+1,out}$ . (with D many features, we get 2D+1 coins and  $\theta_0$  is the guy who choses the label, then for each feature you'll get 2D possibilities)

*examples:* Given 50+ve, 150 -ve examples,  $\theta_0 = 50/(50+150) = 1/4 = p(y)$ . Take 50+ve ex and look at the outlook feature. Say 30 have outlook=1 and 20 have outlook=0. This is  $\theta_{+1,outlook} = 30/50 = 3/5$ . So training naive bayes is very easy, all you need to do is count. Very fast.

*Pf:* Just talking about labels so  $\theta_0 = \theta$  Super easy case, get bunch of data  $p(y_1, \dots, y_{200}) = \prod_n p(y_n) = \prod_n \theta^{(y_n=1)} (1 - \theta)^{(y_n \neq 1)}$  Now take the log so that it's easier to maximize. We can take the log because it's monotonically increasing. So:

$$\begin{aligned} &= \prod_n \theta^{(y_n=1)} (1 - \theta)^{(y_n \neq 1)} \\ &= \sum_n ([y_n = 1] \log(\theta) + [y_n \neq 1] \log(1 - \theta)) \\ \frac{\partial}{\partial \theta} &= \sum_n \left( \frac{y_n = 1}{\theta} + \frac{-[y_n \neq 1]}{(1 - \theta)} \right) = 0 \\ &= \frac{1}{\theta} \sum [y_n = 1] \qquad \qquad \qquad = \frac{1}{1 - \theta} \sum [y_n \neq 1] \\ (1 - \theta) \sum [y_n = 1] &= \theta \sum [y_n \neq 1] \\ \sum [y_n = 1] &= \theta (\sum [y_n \neq 1] + \sum [y_n = 1]) \\ \theta &= \frac{\sum_n [y_n = 1]}{N} \end{aligned}$$

This is what we wanted!!!

### 9.2.2 a

Objective:

$$\max_{\theta} \prod_k \theta_k^{[y_n=k]} \text{ subject to } \sum_k \theta_k = 1$$

Now rolling a die instead of the coin. Without the constraints,  $\theta_k$  will go to infinity. Take the log of that and get

$$g(\theta) = \sum_n \sum_k [y_n = k] \log(\theta_k) \text{ subject to } \sum_k \theta_k - 1 = 0$$

This is concave, with the restriction, it's still concave around the line that sums to one ( $\sum \theta_k = 1$ ), it's concave in the whole space and over my restriction, so if we maximize we'll get the the most.

Now construct the lagrangian  $\mathcal{L}(\theta, \lambda)$  For every constraint, add a  $\lambda$ , a lagrange multiplier. So here we only have one constraints and we get

$$\mathcal{L}(\theta, \lambda) = g(\theta) + \lambda(\sum_k \theta_k - 1)$$

Then,

$$\max_{\theta} \min_{\lambda} \mathcal{L}(\theta, \lambda) = \max_{\theta} g(\theta) \text{ subject to constraints}$$

(What ever i'm trying to do  $\lambda$  tries to do the opposite.) Now we can do  $g(\theta)$  without the constraints by solving  $\min_{\lambda} \max_{\theta} \mathcal{L}(\theta, \lambda)$ .

Idea: All i care is making  $g(\theta)$  big. But then there's  $\lambda$  who's trying to make  $\mathcal{L}$  small. Say I find  $\sum_n \theta_n = 2$ , and makes  $g(\theta)$  happy. But then  $\lambda$  will be set to  $-\infty$ , because  $\sum_n \theta_n - 1 = 1$ , so now  $\lambda$  has full control on making  $\mathcal{L}$  small. If  $\sum_n \theta_n = 0.5$ , then  $\sum_n \theta_n - 1 = -0.5$ , so  $\lambda$  will go to  $\infty$ . The only thing I can do is to actually satisfy my constraints  $\sum_n \theta_n - 1 = 0$ . i.e. if we get  $-\infty$ , that means we didn't satisfy the constraints, if we do get something other than that we know we've satisfied it.

## 10 Class 11 October 13th

HW

1.

$$\max_{\theta} \prod_d \theta_d^{x_d} \text{ s.t. } \sum_d \theta_d^2 = 1, (\theta_d \leq 0)$$

$\theta_d^2 \leq 0$  will always happen but technically need to include it. Since we're maximizing, we want to insure that the function is concave, and log is concave (second derivative should be non-positive). Make sure to check this! The overall goal is to  $\max_{\theta}$  and  $\lambda$  plays the other role. Take the log:

$$\begin{aligned} &= \max_{\theta} \sum_d x_d \log(\theta_d) \text{ s.t. } \sum_d \theta_d^2 = 1 \\ \mathcal{L}(\theta, \lambda) &= \sum_d x_d \log(\theta_d) + \lambda(\theta_d^2 - 1) \\ &= \max_{\theta} \min_{\lambda} \mathcal{L}(\theta, \lambda) \end{aligned}$$

The order of  $\max_{\theta} \min_{\lambda}$ : I pick  $\theta$  then  $\lambda$  picks, if I pick  $\theta$  that doesn't satisfy the constraint,  $\lambda$  will pick  $-\infty$ . If we had  $\min_{\lambda} \max_{\theta}$ ,  $\lambda$  picks  $\leq 0$ , then  $\theta$  will be  $\infty$ , if  $\lambda > 0$ , I will still pick bad  $\theta$ ... You need to choose  $\theta$  first,  $\min_{\lambda} \mathcal{L}(\theta, \lambda)$  is a function of  $\theta$ , to figure out what  $\lambda$  is, need to choose  $\theta$  first. That's why the order matters..? (unresolved) Whatever, pick  $\theta$  first! Now differentiate with respect to  $\theta_d$

$$\begin{aligned} \partial \mathcal{L} / \partial \theta_d &= x_d / \theta_d + 2\lambda \theta_d \\ -x_d / \theta_d &= 2\lambda \theta_d \\ 2\lambda \theta_d^2 &= -x_d \\ \theta_d &= \sqrt{\frac{-x_d}{2\lambda}} \end{aligned}$$

Now solve for  $\lambda$ . Substitute the constraint  $\sum_d \theta_d^2 = 1$

$$\begin{aligned} \sum_d \frac{x_d}{2\lambda} &= \frac{-1}{2\lambda} \sum_d x_d = 1 \\ \lambda &= -1/2 \sum_d x_d \end{aligned}$$

This means,  $\theta_d = \sqrt{\frac{x_d}{\sum_d x_d}}$ . Need to make sure that the positive solution is the right solution (unless there's a constraint that  $\theta_d \leq 0$ ).

### 10.1 Maximum Likelihood

Just to remind ourselves, we were trying to maximize the joint probability  $p(x, y)$  i.e.  $f(x) = \arg \max_y p(x, y)$ . The generative story is:

1. For each example  $n = 1 \dots N$ 
  - (a) choose  $y_n \sim \text{Ber}(\beta)$

- (b) for each  $d = 1 \dots D$ , choose  $x_{n,d} \sim N(\mu_{y_n,d}, \sigma_{y_n,d}^2)$  or  $Ber(\theta_{y_n,d})$  any reasonable distribution. For each class and each feature there is a distribution.
- (c)

This says if I know my model and the parameters,  $\mu$ s and  $\sigma$ , we can, so the idea is to reverse engineer these parameters. Example in matlab:

```
y = rand < 0.6
mu = zeros(2); mu(1,1) = 8; % where 1 is pos, 2 is neg, give random number
mu(1,2) = 11; mu(2, 1) = -4; mu(2,2) = 10;
mu(y+1, :) % ans = 8 11
X = repmat(mu(y+1, :), 100, 1) + randn(100,2))
X2 = repmat(mu(2, :), 100, 1) + randn(100,2))
plot(X(:,1), X(:,2), 'b.') hold on;
plot(X2(:,1), X2(:,2), 'r*')
mean(X) % solution/estimate, gaussian with these two mean
mean(X2) % => pretty close to the actual mean
```

Turns out that if we want the decision boundary, it just becomes the linear classification. **LLR**: Re-write as  $\log \frac{p(y=+1|x)}{p(y=-1|x)} = \mathbf{x} \cdot \mathbf{w} + b$ . The decision boundary is just linear. This is assuming bernoulli model on the features. If we did the gaussian case, it works out too. But not with any distribution. Point is that this is a linear model. The parameters of the model  $\theta_{\pm 1,d}$  you can compute by counting. (If gaussian model, find the  $\mu, \sigma$  of data for each feature and each label). Constrasting with perceptron or other iterative classifiers, this we get the closed form.

We have a linear model and closed form solution, but there's no guarantee about seperability or large margins.

## 10.2 Generative vs Conditional

Generative is trying to model  $p(x, y)$ , but all we care about is  $p(y|x)$  looks like. This is the conditional model. The big question in the generative was what probability distribution do we chose, then estimate the parameters. Same in conditional model, choose a distribution for  $p(y|x)$ , then estimate the parameters.

We like linear models, so maybe we want  $\mathbf{w} \cdot \mathbf{x} + b$ , wehre  $\mathbf{w}$  is the parameters of the model. Since this is in the  $\mathbf{R}$ , we need to squish it down to interval  $[0, 1]$

$$\text{sigmoid: } \sigma(z) = \frac{1}{1 + e^{-z}}$$

If  $wx + b \gg$ ,  $\sigma = 1$ , if  $wx + b = 0$ ,  $\sigma = 0.5$ , if  $xw + b \ll$ ,  $\sigma = 0$ . So we have

$$p(y = +1|x) = \sigma(\mathbf{w} \cdot \mathbf{x} + b)$$

$$p(y = -1|x) = 1 - \sigma(\mathbf{w} \cdot \mathbf{x} + b) = \sigma(-(\mathbf{w} \cdot \mathbf{x} + b))$$

Since  $1 - \sigma(z) = \sigma(-z)$  (symetric property on sig) So

$$p(y|x) = \sigma(y(\mathbf{w} \cdot \mathbf{x} + b))$$



We're saying the distribution over the label is a function of the margin, which makes sense. Any function that maps  $[-\infty, \infty]$  to  $[0, 1]$  is good.

What we really care about is  $\log(p(y|x)) = \log \sigma(y(wx+b))$  Since  $\log(\sigma(z)) = -\log(1 + e^{-z})$ ,

$$\log p(y|x) = -\log(1 + e^{-y(wx+b)})$$

The *logistic loss function*!!

So if you optimize a linear classifier under logistic loss function, it's the same as optimizing  $p(y|x)$ .

### 10.3 Discriminative Models:SVM

Generative and Conditional models are always probabilistic, for Discriminative there's no probability. We're explicitly optimizing the decision boundry for discriminative. Conditional models optimizes the probability of label given data. Generative models optimizes predicting both the label and the data.

If you care about probabilities, you should do generative, conditional models. When to prefer either or?

#### 10.3.1 Current State of Generative vs Conditional

1. Assume the generative model is "correct". (Tell the story about how the data got generated and it's actually true, always false but have it for now). Then,
  - (a) as  $N \rightarrow \infty$ ,  $\epsilon^{gen} \rightarrow \epsilon^{Bayes}$ ,  $\epsilon^{cond} \rightarrow \epsilon^{Bayes}$  Asymptotic
  - (b) finite  $N$ ,  $\epsilon^{gen} \rightarrow \epsilon^{Bayes}$  faster than  $\epsilon^{cond} \rightarrow \epsilon^{Bayes}$
2. Assume model not correct
  - (a) as  $N \rightarrow \infty$ ,  $\epsilon^{gen} > \epsilon^{cond}$ , Conditional is better because if the model is even slightly wrong, if we have infinite amount of data, it blows up.
  - (b)  $N$  finite,  $\epsilon^{gen} \rightarrow$  it's optimum faster than  $\epsilon^{cond} \rightarrow$  it's optimum. We'll never get to  $\epsilon^{Bayes}$ . It's no longer asymptote at the Bayes line, and they asymptote at different places, and conditonal's optimum might be better, but generative goes to its optimum faster.

Hall's heuristic is that if  $N > 2D$ , use conditional,  $N < D$  use generative. Generative adds more information, where data comes from, to the model. So when  $N$  is small, it's very good. But when  $N$  is large, it relies too much on your inductive bias about the data.

Conditional and Discriminative only makes sense for supervised.

## 11 Class 11 October 18th 2011

### 11.1 Lagrangian

*Goal:*  $P(y|x)$ , given examples  $(x_i, y_i), i = 1, \dots, N$  One thing to do is want my feature's distribution to match the data. My model has some  $E_{(x,y)\tilde{p}}[y = 1 \wedge \text{sunny}] = \frac{1}{N} \sum_{(x_n, y_n)} [y = 1 \wedge \text{sunny}]$  i.e. want  $E_P[f_i(x, y)]$  to match the empirical data. So:  $E_P[f_i(x, y)] = E_D[f_i(x, y)]$   $P$  what i'm going to learn,  $D$  is the training data.

if we store the training data and say  $p(y|x) = \begin{cases} 1/N & \text{if } (x, y) \in D \\ 0 & \text{otherwise} \end{cases}$  it will

work, but this is overfitting.

Extreme case of underfitting is uniform distribution. Want something in between, so want to match the data, but we also want  $p(y|x)$  to be uniform as possible.

**Entropy:**  $H(q) = -\sum_z q(z) \log(q(z))$

So we want to maximize of the entropy of  $p$ . so

$$\max_{p(y|x)} - \sum_y p(y|x) \log(p(y|x)) \text{ s.t. } E_P[f_i(x, y)] = E_D[f_i(x, y)] \forall i$$

This is maximizing over a function, which is scary. But from "Calculus of variations" that study functions over functions, the magic is that everything we know holds.

Now remove the - to change max to min, and add lagrange multipliers and get:

$$\min_{p(y|x)} \max_{\lambda} \sum_{x \in D} \sum_y p(y|x) \log p(y|x) + \sum_i \lambda_i [E_P(f_i) - E_D(f_i)]$$

$E_P[f_i]$  is  $\sum_{x \in D} \sum_y p(y|x) f_i(x, y)$ ,  $E_D[f_i]$  is some number.

Now take the derivative with respect to  $p(y|x)$ :

$$\begin{aligned} \frac{\partial}{\partial p(y|x)} &= \sum_{x \in D} \sum_y [\log p(y|x) + 1] + \sum_i \lambda_i \sum_{x \in D} \sum_y f_i(x, y) = 0 \\ \sum_{x \in D} \sum_y [\log p(y|x) + 1] &= - \sum_i \lambda_i \sum_{x \in D} \sum_y f_i(x, y) \\ \sum_{x \in D} \sum_y [\log p(y|x) + 1] &= - \sum_{x \in D} \sum_y \sum_i \lambda_i f_i(x, y) \end{aligned}$$

Then we need to show  $\log p(y|x) = -\sum_i \lambda_i f_i(x, y) - 1 \forall x, y$  Now take the exponential:

$$p(y|x) = \exp[-\sum_i \lambda_i f_i(x, y) - 1] = \exp[-\sum_i \lambda_i f_i(x, y)] / e$$

This is like  $\sigma(w \cdot x) = 1 / (1 + \exp(-w \cdot x))$ , the sigmoid linear regression! it's linear model again, where  $\lambda$ s are the  $w$ , even though we didn't start off with a linearity expectation.

## 11.2 Neural Networks

Dealing with non-linearity..

Each of the feature of  $x$  will be an input, and the goal is to get  $y \in \mathbf{R}$  from this. A linear model has weight on each edge from the input to the last node, where the value of the last node is the sum of  $w_i x_i$ . This can't solve XOR, so we introduce *hidden units*.

Connect each input to each hidden units, a fully connected graph, where each edge will have a weight  $w_{d,i}$ , where it goes from input  $x_d$  to hidden node  $i$ , and hidden nodes will have weights  $u_i$  on edges that connect to the output unit.  $y = \sum_i u_i (\sum_d w_{d,i} x_d)$ , but this is still linear! To introduce non-linearity, we'll introduce a link function,  $f$ , usually sigmoid,  $f(z) = \tanh(z)$

Let's do XOR: OR:  $y = x_1 \vee x_2$ , where  $(1, 1), (-1, 1), (1, -1) \rightarrow y = 1$ , and  $y = -1$  for  $(-1, -1)$ . Here you want the bias to be positive so that the only way you can get -1 is if both are negative. Let  $b = 1000$ , then if  $w_1 = w_2 = 1000$ , then we'll get the OR behavior. Construct AND, then using OR and AND (with a bias), we can do XOR.

Three parts involved training, architecture selection, and selection of link functions.

### 11.2.1 Forward and Back propagation

Forward: given  $\mathbf{x}$ , get  $y$ . Backward:  $\partial L / \partial u$  is how much the error on my final prediction changes as  $u$  changes and so forth.

### 11.2.2 Architecture

*Terminology:*

- no hidden units  $\rightarrow$  one layer network, the linear model, no xor
- 1 set of hidden units  $\rightarrow$  2 layer network, the universal function approximator (essentially infinite representation power), but exponentially many hidden units.
- "deep" ( $\geq 3$ ), also a universal function approximator, the hope is fewer hidden units.

## 12 Class 12 October 20th 2011

### 12.1 more on backpropagation

First calculate the error: the derivative of squared loss is the difference at the output:  $e = y - \hat{y}$ . Then calculate the gradient for the last edge  $-eh_i$ . Gradient descent subtracts, so we add all the gradients, then subtract it from the weight later.

Then compute the errors as  $h_i$  nodes. Need to look at all edges connected to  $h_u$ .

$$e_u \leftarrow e_u + e_v w_{u,v} (1 - \tanh^2(a_u))$$

For square loss, we're trying to minimize  $L = \frac{1}{2}(y - g(x))^2$ , where  $g(x)$  is the output of ANN. We want  $\frac{\partial L}{\partial w_{u,y}} = (y - g(x)) - g'(x) = -e_y \frac{\partial g(x)}{\partial w_{u,y}}$  here  $g(x) = w_{u,y} \cdot f(a_u)$ , so  $\frac{\partial g(x)}{\partial w_{u,y}} = f(a_u)$

Now, we need  $\frac{\partial L}{\partial w_{v,u}} = \frac{\partial L}{\partial g(x)} \frac{\partial g(x)}{\partial w_{v,u}}$

$$\begin{aligned} \frac{\partial g(x)}{\partial w_{v,u}} &= w_{u,y} \frac{\partial f}{\partial w_{v,u}} f\left(\sum_{z \rightarrow u} w_{z,u}\right) \\ &= w_{u,y} f'(a_u) \frac{\partial}{\partial w_{v,u}} w_{z,u} \cdot h_z \\ &\dots \\ &= \left( \prod_{c \rightarrow d \text{ is in path to } y} w_{c,d} f'(a_c) \right) \cdot \mathbf{h}_{\text{last thing in path}} \end{aligned}$$

So for the hw, using logistic,,

$$\begin{aligned} \frac{\partial}{\partial w_j} \log(1 + e^{-yg}) &= \frac{1}{1 + e^{-yg}} \frac{\partial}{\partial w_j} e^{-yg} \\ &= \sigma(yg) e^{-yg} \frac{\partial}{\partial w_j} (-yg) \\ &= \sigma(yg) e^{-yg} (-y) \frac{\partial}{\partial w_j} g \end{aligned}$$

and now the first error  $e_y$  is not  $y - g$  but  $\sigma(yg)e^{-yg}(-y)$ .

### 12.2 Architecture Selection

In a 2 layer network, there are 3 layers, fully connected.  $D$ -dimensional inputs,  $X$  number of input, than 1 output. Question is how many hidden units?

If there are only 1 hidden units, underfitting, because it's linear.

If you go from  $D/2$  to  $D/2 + 1$  units, do I get more expressive power? Yes! It's like adding another bit in binary. (another combination/expression)

$W = \mathbf{R}^{D \times K}$ , if  $D = 1000$ , you don't want  $K$  to be more than 100 because it's just a lot of computation and requires a lot of examples (2 times the number of parameters  $DK$ ).  $K > 100$  is rare.

$N$ -layer network,  $D$  inputs, usually number of hidden nodes decrease  $D \geq |K_1| \geq |K_2| \geq \dots \geq 1 = |K_N|$ . Very slow to train.

## 13 Class 13 Oct 25th

History

- Tikhonov regularization minimize loss (linear) + regularizer
- Kernelization of linear models (how can i take a linear model and make it non-linear for free)

Neither works really well, but add them together was SVM which worked well. SVM is just another kernelized model.

### 13.1 Kernel Methods

Linear model  $w \cdot x$  want to map this to higher dimension  $w_{big} \cdot \phi(x)$ .

In 1D linearly seperable data, a linear seperator is just a dot, a threshold. If it's not linearly seperable, map  $x$  to  $(x, x^2)$ . Now we have 2D data. Under this feature mapping,  $\phi(x) = (x, x^2)$ , this is linearly seperable. In 1D, this seperation is not linear.

Problem is that with P1 code, this will take forever.  $\phi(\mathbf{x}) = (x_1, x_2, \dots, x_d, x_1^2, x_1x_2, \dots, x_1x_d, \dots, x_dx_1, \dots)$ . Roughly  $d^2$  more points. (+ + - + + - requires cubic map, + + - + + - requires quadratic map)

For all of these mappings,  $\phi(x) \cdot \phi(z)$  can be computed sublinearly in the expanded mapping. A purely computational trick. With quadratic feature mapping,

$$\begin{aligned}\phi(x) \cdot \phi(z) &= (1 + x \cdot z)^2 \\ &= (1 + x \cdot z)(1 + x \cdot z) \\ &= 1 + x \cdot z + (x \cdot z)^2\end{aligned}$$

gives you same thing with some constants.

1D case:  $(1+xz)(1+xz) = 1+2xz+x^2z^2$ . Compare this to  $(x, x^2) \cdot (z, z^2) = xz + x^2z^2$ . We have our linear term, quadratic, and a constant

**Why does this constant not matter?** 1 will just affect the bias. the 2 is just the weight on  $xz$ , so the representation is not exactly the same. But this will affect my regularizer, my weights won't be small anymore, may not have the best margin anymore. But we're happy with these constants.

Problem is, all of my algorithm depends on  $w_{big} \cdot \phi(x)$ . But  $w$  is not in  $\phi$ .  $w$  is in the big space, also  $w \neq \phi(\text{something})$ , most likely.

How can we take our current algorithms so that it only depends on  $\phi$  terms.

### 13.2 Kernelized Perceptron

```
for each  $(x, y)$  do
   $a \leftarrow w \cdot \phi(x)$ 
  if  $y_na \leq 0$  then
     $w \leftarrow w + y_n\phi(x_n)$ 
  end if
end for
return  $w$ 
```

**Representer Theorem:** (Any algorithm you have should have a corresponding representer theorem) I can write my vector, somehow, so that it only depends on the feature vector of the training points  $\phi(x)$

Want to say that optimal weight vector also lies in  $\phi$  space. If  $w$  points in the direction of data, but if  $w$  moves to the nullspace of data, there's a component of  $w$  that's entirely ignored, which won't help my margin (i'd paying a margin penalty and getting nothing in return). So if  $w$  points in the direction of data, we'll have large margin.

*From representer theorem for perceptron:* At any point we can write  $w$  as

$$\sum_m \alpha_m \phi(x_m)$$

for some  $\alpha_m \in \mathbf{R}$

Any time I do an update,  $w$  changes by a constant times feature vector of a single point.

Now, our objective is to learn  $\alpha$ .

How many components does  $w$  have? Say  $\phi(x)$  is quadratic, then  $w$  have  $d^2$  components. How many  $\alpha$ s?  $N$ , one for each data. So we win  $N$  is smaller than  $D^2$ , otherwise don't do it!!

Now:

```

 $\alpha \leftarrow 0$ 
for each  $(x, y)$  do
   $a \leftarrow \sum_m \alpha_m \phi(x_m) \cdot \phi(x)$ 
  if  $y_n a \leq 0$  then
     $\alpha_n \leftarrow \alpha_n + y_n$ 
  end if
end for
return  $\alpha$ 

```

It's not always good to go to higher dimension, because if your data was already separable in lower dimension, the result might be that your margin get smaller. Usually the time it takes to train is proportional to the size of the margin, so it'll take longer and make the margin smaller. So in practice ppl only do cubic and quadratic.

### 13.3 Why is it called a kernel?

We've re-written the algorithm so that it'll depend on  $\phi(x) \cdot \phi(z) = (1 + x \cdot z)^d$ . Write this as  $K(x, z) = (1 + x \cdot z)^d$ , the polynomial kernel. Now make up a crazy kernel function, and does it maintain the property we care about..? (like converge for sure if it exists, converge in the bounded time  $(1/\gamma^2)$ ) and this is NO, not any function works.

$K(x, z)$  is okay if  $\exists \phi$  s.t.  $K(x, z) = \phi(x) \cdot \phi(z)$ . Functions that are okay to use are the ones that actually do correspond to the dot product after some mapping. This space where dot product is defined is *Hilbert space*.

How do you do this? You provide the mapping, or come up with equivalent property. In this case,  $K(x, z)$  is okay if  $K$  is *symmetric positive semi-definite*.

**Definition:**  $a \in \mathbf{R}$  is non-negative if  $\forall x \in \mathbf{R}, xax \leq 0$ . i.e.  $ax^2 \leq 0$

**Generalization:**  $A \in \mathbf{R}^{D \times D}$  is positive semi-definite if  $\forall x \in \mathbf{R}^D, x^T A x \leq 0$ . ( $x^T A x$  is  $(1 \times D)(D \times D)(D \times 1) \in \mathbf{R}$ )

**Symmetry:** Say  $x_1, x_2, \dots, x_n$ , we have  $N$  by  $N$  matrix  $A$ , where position  $K(x_3, x_2) = K(x_2, x_3)$ , called the Gram matrix.

Equivalent of saying  $K$  is a symmetric positive semi-definite function is,  $\forall$  datasets  $D$ , the matrix  $A$  defined by  $K$  on  $D$  is symmetric positive semi-definite.

**the real definition:**  $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbf{R}$  is symmetric positive semi-definite if  $\forall f$ ,

$$\int \int f(x)K(x, z)f(z)dx dz \geq 0$$

This definition is useful because it gives us a way to combine multiple kernels to get new kernels.

Example: Suppose  $K$  is a kernel (symmetric, positive semi-definite). Then  $L(x, z) = aK(x, z) \forall a \in \mathbf{R} > 0$  is also a valid kernel. If  $K$  is symmetric,  $L$  is also symmetric.  $a$  is just a constant, so we're good.

### 13.3.1 This is called the kernel method because:

$\phi : \mathcal{X} \rightarrow \mathcal{H}$ ,  $\mathcal{H}$ , hilbert space =  $\{f : \mathbf{R}^D \rightarrow \mathbf{R} \mid \int f(x)^2 dx < \infty\}$  Points in this space are functions. aka  $\mathcal{L}_2$  The natural definition of dot product in this space is  $f \cdot g = \int f(x)g(x)dx$ .

So our kernel  $K$  now takes 2 points  $x, z$ ,  $K(x, z) = \phi(x) \cdot \phi(z)$ .. next time!

**14 Class 14 on paper..**



## 15 Class 15 November 1st 2011

Take-home midterm handed out this thursday. Up to end of kernel methods.

### 15.1 Ensemble Methods

1. Independent ensembles: build multiple classifiers independently. Bagging, random forests
2. Dependent ensembles: typically greedy, build one classifier, and using that as a basis build a second one, i.e. boosting

Notion of independent/dependent has to do with training time, not test time.

Pro: for independent ensembles: can train them in parallel. Con: learning a bunch of stuff, and hoping that it does well together, where dependent ensembles actually makes a weak learner stronger. Dependent usually requires smaller number of ensembles by its construction.

#### 15.1.1 Bagging

Dataset  $\tilde{\mathcal{D}} : \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ , a finite distribution over these examples.  $\tilde{\mathcal{D}}$  will put probability in each of these examples i.e.

$$\tilde{\mathcal{D}} = \sum_{n=1}^N \frac{1}{N} \delta_{x_n}$$

We have bunch of examples on  $\mathcal{X} \times \mathcal{Y}$ , each of  $\delta_{x_n, y_n}$  is the probability point mass centered at  $(x_n, y_n)$ , that has probability  $\frac{1}{N}$  but 0 everywhere else. Our empirical probability distribution over the training data.

It's easy to sample from this distribution, just pick.

**Bootstrap resampling:** Assume  $\mathcal{D}$ , the true distribution, where these examples come from.  $D$  is our training data, drawn from  $\mathcal{D}$ . We would like to draw from  $\mathcal{D}$  a lot, but that's expensive (to get labeled training data). Using  $\mathcal{D}$  and  $D$ , define  $\tilde{\mathcal{D}}$ , and draw  $D_1, D_2, \dots, D_M \sim \tilde{\mathcal{D}}$ . Magic: each  $D_M$  is a draw from  $\mathcal{D}$ .

How to draw  $D_m \sim \tilde{\mathcal{D}}$ :

```
for  $n = 1, \dots, N$  do  
  choose  $I$  uniformly at random from  $[1, \dots, N]$   
  add  $(x_i, y_i)$  to  $D_m$   
end for
```

We're sampling with replacement so it's not the exact same set. Ensemble is only interesting if  $D_i$ 's are all different. Look at  $D_i$  and  $(x_1, y_1)$ . Probability of picking  $(x_1, y_1)$  in  $D_i$  is  $\frac{1}{N}$ , so the probability of NOT picking that is  $(1 - \frac{1}{N})$ . The probability that  $D_5$  does not contain  $(x_1, y_1)$  at ALL is

$$(1 - \frac{1}{N})^N \approx \frac{1}{e} \approx 60\%$$

In general, only a 60% chance that the data will be in the dataset. So  $D_i$ s will be pretty diverse (not so much overlap).

Now train classifiers on these  $D_i$ s and get  $f_1, \dots, f_n$ , predict and get  $\hat{y}_1, \dots, \hat{y}_n$ , then we take the majority vote as our answer.

**When is bagging a good idea?** If you're worried about overfitting, it's a good idea to use bagging. If each  $f_i$  is very different from other classifiers, and each one is massively overfitting, so predictions are very different,  $\hat{y}_i$ 's have very high variance, so the final majority vote has low variance. Bagging and regularization like two sides on a same coin. So if you have a classifiers that's hard to regularize (like perceptron), you can do bagging.

### 15.1.2 Boosting

What kind of classifiers are good? If it gets low test error on *most* datasets → **strong learning**. We could ask for 0 error, which is impossible. So best I can hope for is low error. But if I ask for all datasets, that's also impossible. Formally, we want

$$P(\text{test error} < \epsilon) \geq 1 - \delta$$

If  $\epsilon = 0.1, \delta = 0.01$ , we want  $P(\text{test error} < 10\%) \geq 99\%$ . This is called **PAC-learning** (kind of dead, unless you're in theory). Strong learning is hard.

**Weak learning**: low error ( $< 50\%$ ) on at least  $\frac{2}{3}$  or 51% of the dataset. Something that does just slightly better than chance. Can we convert this into a strong learner? yes by boosting (a dependent ensemble).

The first boosting algorithm in polynomial time was *AdaBoost*, the first practical algorithm.

### 15.1.3 AdaBoost

$D$ , our dataset, an empirical distribution over the training data  $\mathcal{X} \times \mathcal{Y}$ . Say  $N = 5$

Step 1 study each question (data points) the same amount

Step 2 evaluate yourself. (must at least more than 50% of the weighted importance). Say we get 1,3,5 correct

Step 3 Make 4 and 2 (the errors) more important

Step 4 go to step 1

Train:

```

 $d_n = \frac{1}{N}$ 
for  $k = 1, \dots, K$  do
  train  $f^{(k)}$  on  $(D, d)$ 
   $\epsilon^{(k)} \leftarrow$  weighted error of  $f^{(k)}$  on  $d < 0.5$  (by assumption of weak learner)
  i.e  $\sum_n [f^{(k)}(x_n) \neq y_n]$ 
   $\alpha^{(k)} = \log(\frac{1-\epsilon^{(k)}}{\epsilon^{(k)}}) > 0$  (if the weak learner is perfect, this will be  $\infty$ )
  update weights  $d_n \leftarrow d_n \exp(-\alpha^{(k)} y_n f^{(k)}(x_n))$  ( $y_n f^{(k)}(x_n)$  is +1 if correct, -1 if wrong). So if we get correct, we'll decrease the weight, if we don't we'll increase the weight. Also because exponential, if we routinely get certain examples correct it'll go down very quickly.
  renormalize so  $\sum_n d_n = 1$ 
end for

```

Test time, predict

$$\mathbf{1}[\sum_{k=1}^K \alpha^{(n)} f^{(k)}(x)]$$

So if we get 0 error with  $f^{(k)}$ , that  $\alpha^{(k)} = \infty$ , so all the decision is made by  $f^{(k)}$  at test time, which is fine because that's the best. (Question of overfitting aside).

Maybe the weak learner just memorizes couple (to satisfy the requirement) and maybe this is overfitting. But boosting doesn't overfit as much as you might expect. We don't understand why really. Works really well!

Con: we can't parallelize training time. But we can at test time.

## 16 Class 16 Nov 3rd 2011

### 16.1 Stochastic GD

$$\begin{aligned}\mathcal{F}(w) &= \sum_n l(y_n, f(x_n)) + \lambda R(f) \\ &= \text{sum}_n[l(y_n, f(x_n)) + \frac{\lambda}{n} R(f)]\end{aligned}$$

Then  $\nabla F \text{sum}_n[\nabla_f l + \frac{\lambda}{N} \nabla_f R]$

in GD: the update is  $f \leftarrow f - \eta \sum_n [\nabla_f l_n + \frac{\lambda}{N} \nabla_f R]$

in SGD: for  $n = 1, \dots, N$ ,  $f^k \leftarrow f^{k-1} - \eta [\nabla_f l_n|_{f^k} + \frac{\lambda}{N} \nabla_f R|_{f^k}]$

So hw, with hinge loss:  $l(y_n, w \cdot x_n) = \begin{cases} 0 & \text{if } y_n w \cdot x_n \leq 1 \\ 1 - y_n w \cdot x_n & \text{otherwise} \end{cases}$  and

$R(w) = \frac{1}{2} \|w\|^2$ . So

$$\nabla_w l(y_n, w \cdot x_n) = \begin{cases} 0 & \text{if } y_n w \cdot x_n \leq 1 \\ y_n x_n & \text{otherwise} \end{cases}$$

and  $\nabla_w R = w$ .

so

```
for  $n = 1, \dots, N$  do
  if  $y_n w^k \cdot x_n < 1$  then
     $w^{k+1} \leftarrow w^k + \eta y_n x_n$ 
  end if
   $w^{k+1} \leftarrow w^{k+1} - \eta \frac{\lambda}{N} w^k$ 
   $k \leftarrow k + 1$ 
end for
```

This is just like perceptron! to make it into a perceptron, set  $\lambda = 0$  (no regularizer),  $\eta = 1$  (no learning rate), and shift hinge loss to  $-y\hat{y}$  rather than  $1 - y\hat{y}$ . So in othewords, perceptron is doing stochastic GD with shifted hinge loss, no regularizer, and with no learning rate.

$\eta \frac{\lambda}{N} w^k$  is called the *weight decay*. Most of your effort is used here... In practice we really need to do this. very important. makes a huge difference. Otherwise the algorithm is impractical

Hinging at 1 forces you to have a margin of 1, but hinging at 0 just forces you to have a decision boundary. This is why perceptron does not give us margin. It's not regularized so perceptron might over fit, and since the learning rate is constant, the last example it process could potentially have a large effect, and it's bad when it's an arbitrary data. i.e. particularly sensitive to ordering of data. So we want the learning rate to decrease over time.

Why is it called "stochastic"?

In normal optimization, we say  $\min_x F(x)$ , assuming we can compute  $\nabla F|_x$  (at any  $x$ ). In stochastic optimization, we're trying to  $\min_x \mathbb{E}[F(x)]$ , and we can't compute this expectation nor the  $\nabla \mathbb{E}[F]$  (randomness is the choice of example you're looking at at a single point).

## 16.2 Convergence

We're trying to  $\min_x f(x)$ , say the true minimum is  $x^*$ ,  $f(x^*) = f^*$ . The function we learn at iteration  $k$ :

$$f_{best}^{(k)} \leq f^* + \frac{\|x^{(0)} - x^*\|^2 + G^2 \sum_k \eta_k^2}{2 \sum_k \eta_k}$$

$G$  is something that says that  $f$  isn't too steep. It measures the magnitude of  $f$ . Ideally, we want the right hand side (of the rhs, the fraction) to go to 0. To do that, the denominator needs to go to  $\infty$ . So we want  $\sum_k \eta_k \rightarrow k$  as  $k \rightarrow \infty$ . But if  $\eta_k$  diverge too quickly, the numerator will also go to  $\infty$ . So we also want  $\sum_k \eta_k^2 < \infty$ . Choose  $\eta_k = \frac{a}{b+k}$ , for some constant  $a, b$ . In practice, everyone uses  $\eta_k = \frac{1}{\sqrt{k}}$ .

## 16.3 Feature Hashing

Say we have  $x \in \mathbf{R}^{10^9}$ . So big. So we want have a hash function  $h : \mathbf{R}^{10^9} \rightarrow \mathbf{R}^{11}$ . to make  $x$  into  $\phi(x)$ . like  $h(x) = az + b\%11$ . When you compute  $\phi(x) \cdot \phi(z)$ , you get  $x \cdot z +$  quadratic "stuff". The difference between  $\phi(x) \cdot \phi(z)$  and  $x \cdot z$  is precisely the quadratic terms. it's not that bad! We need pair wise independence. (Stemming in NLP does not satisfy)

## 16.4 Final Project

- Arbitrary sized groups (more the ppl higher expectation)
- Research project is fine, double countint is fine.

You should demonstrate that you learned something in this class. Anything is fair game as long as you do this. Just turn-in write up. (he won't look at the code) In hw format, 4pg (not more than 7 or 8). Aim for 4. Need a short presentation (5min-ish, 3 slides). (depends on how many groups total)

## 17 Class 17 November 8th

Went over midterm..

## 17.1 Unsupervised Learning: K-means

Does it converge? Yes. How fast? In theory, very bad (long), in practice, very fast. How good is the solution? We don't know.

The K-means problem:

$$\min_{\mu} \sum_n \min_k \|x_n - \mu_k\|^2$$

It's highly convex, easy to get stuck. it uses a squared distance. The optimum is the solution to this problem. This is what we want. Floyd's algorithm. Furthest first heuristic (pick one, then choose the furthest one) to do the first initialization. If we have outliers, the first heuristic will likely to choose them as the initial points which could lead to weird behavior, but in general it works better than random selection.

5 years ago, **K-means++** algorithm came up: The furthest first heuristic is roughly deterministic (other than the first step).  $\mu_1$  is a random point, just like before, but to select the second one, compute the distance between  $\mu_1$  and every point, get smallest  $d_i$ 's, and pick the point with probability proportional to  $d_i^2$ . Then do Lloyd's algorithm  $\Rightarrow$  value  $\phi$

We get the guarantee that what we find is "not that bad" than the optimum we want:  $\phi \leq 8(2 + \log k)\phi^{opt}$  with high probability. This bound comes from the initialization, not about Lloyd's algorithm, so it's not very satisfying. We don't know how much better it makes Lloyd.

**How to choose K:** No one really knows how. Not a question we can formulate well as a theory question. By nature so dependant on the data. We can regularize K. AIC and BIC criteria. BIC is the bayes information criteria, Akaike information criteria. On data normalized s.t.  $\|x\| = 1$ . AIC: choose  $K$  to minimize objective value +  $2KD$  BIC: choose  $K$  to minimize objective value +  $K \log(D)$  (likes more clusters. K is punished less).

## 18 Class 21? November 22nd 2011

### 18.1 Semi-supervised Learning

Both labeled and unlabeled data at training time. The unlabeled data could be our test data, if we do this this is called *transductive learning*. What's out there:

1. Generative models for S.S. learning.
2. Self-training: Learn  $f$  on  $L$ , labeled data, run  $f$  on  $U$ , unlabeled data, make it into  $\hat{U}$  (entire thing or only the high confidence ones), re-learn  $f$  on  $L \cup \hat{U}$ .

Can be quite effective. Why? say there are two features that are very effective, but one dominates, so we learn  $f$  mostly on that high frequency data, so when predict on  $U$ , the good feature that wasn't so apparent will come out more (more number, higher correlation), so when we re-learn we get to learn on two good features. May not work in many cases.

3. Co-training: attempt to take something that look like self-training and prove that it works. Someone gives me data split into two disjoint, relatively independent sets. (i.e. webpages split into content  $\bar{x}_n^1$  and link text to the page  $\bar{x}_n^2$ , with labels  $y_n$ ). Where  $U = (\bar{x}_m^1, \bar{x}_n^2)$ . Learn  $f_1, f_2$  on each data independently, then classify on  $U$ , if  $f_1$  and  $f_2$  agree, put the “labeled” data back in training.

Tries to mitigate against mistakes that  $f$  may make.

```

Learn  $f^1$  on  $L^1$  and  $f^2$  on  $L^2$ 
run  $f^1$  and  $f^2$  on  $U$ 
add unlabeled  $\bar{x}_m$ 's to  $L$  only if  $f^1(x_m^1) = f^2(x_m^2)$ 
relearn  $f^1, f^2$  on  $L \cup \{\text{new examples}\}$ 

```

4. Cluster assumptions in discriminative models: decision boundary goes through low-density regions. Say we start with SVM:  $\min_w \frac{1}{2} \|w\|^2 + C \sum_n \xi_n$ , subject to  $y_n(w \cdot x_n + b) \geq 1 - \xi_n \forall n \in L$ . How to use the unlabeled data? regardless of the label, I also want a large margin. i.e. add a constraint:  $|w \cdot x_m + b| \geq 1 - \eta_m \forall m \in U$ , and the entire objective becomes

$$\min_w \frac{1}{2} \|w\|^2 + C \sum_n \xi_n + c_2 \sum_m \eta_m$$

But the new constraint is not convex, so it makes life hard.

Trade off between labeled and unlabeled data. If unlabeled data dominates, SVM may focus on the unlabeled data and ignore the unlabeled data. Very tricky! Standard thing to do is to weigh them equally by using the same  $C$  and divide both by  $N$  (down weight unlabeled data).

5. Label propagation Two half moons dataset. Svms won't work. It's manifold, so we need to construct a nearest neighbor graph across the data (graph Laplacian). Basic idea: heat one end, put an ice to another end, all warm points are positive, cold points are negative etc.

## 18.2 EM for SS-learning

On binary features (binomial)  $L = \{(x = \langle x_1, \dots, x_d \rangle, y)\}$ ,  $U = \{(x = \langle x_1, \dots, x_d \rangle)\}$ . *Generative story:*

```

for  $n = 1, \dots, N$  do
  choose  $y_n \approx \text{Bin}(p)$ 
  for  $d = 1, \dots, D$  do
    choose  $x_{n,d} \approx \text{Bin}(\theta_{y_n,d})$ 
  end for
  return  $(x_n, y_n)$ 
end for
for  $m = 1, \dots, M$  do
  choose  $y_m \approx \text{Bin}(p)$ 
  for  $d = 1, \dots, D$  do
    choose  $x_{m,d} \approx \text{Bin}(\theta_{y_m,d})$ 
  end for
  return  $(x_m)$ 
end for

```

Do MLE:  $\max_{\rho, \theta} P(L, U | \rho, \theta)$

$$\begin{aligned} &= \max_{\rho, \theta} \left[ \prod_n P(x_n, y_n | \rho, \theta) \right] \left[ \prod_m P(x_m | \rho, \theta) \right] \\ &= \max_{\rho, \theta} \left[ \prod_n P(y_n | \rho, \theta) \prod_d P(x_{n,d} | y, \theta) \right] \left[ \prod_m P(x_m | \rho, \theta) \right] \text{ naive bayes assumption} \end{aligned}$$

re-write  $P(x_m | \rho, \theta)$  as  $\sum_{y_m} P(x_m, y_m | \rho, \theta)$ , and break it down as  $\sum_{y_m} P(y_m | \rho, \theta) P(x_m | \rho, \theta)$ .  
Now:

$$= \max_{\rho, \theta} \left[ \prod_n \rho^{y_n=1} (1-\rho)^{y_n=-1} \prod_d \theta_{y_n,d}^{x_{n,d}=1} (1-\theta_{y_n,d})^{x_{n,d}=0} \right] \left[ \prod_m \sum_{y_m} \rho^{y_m=1} (1-\rho)^{y_m=-1} \prod_d \theta_{y_m,d}^{x_{m,d}=1} (1-\theta_{y_m,d})^{x_{m,d}=0} \right]$$

If we were doing only labeled, then the updates are :  $\rho = \frac{\# \text{ of positives}}{N}$ ,  
 $\theta_{1,5} = \frac{\# \text{ of positives with } x_5 = 1}{\# \text{ of positives}}$ .

For unlabeled, it's incomplete because we don't have  $y_m$ , so we want to sum over, but that sum will interfere with the log! so use jensen's inequality:  $\log \sum_{y_m} \dots \Rightarrow \sum_{y_m} q(y_m) \log \dots$ , this is the EM trick. Based on current parameters, guess  $q(y_m) = P(y_m | x_m, \rho, \theta)$ .  $q$  is what do I think the label is supposed to be.

$$\begin{aligned} q(y_m) &= P(y_m | x_m, \rho, \theta) \\ &= \frac{1}{\text{stuff}} P(y_m | \rho, \theta) p(x_m | y_m, \rho, \theta) \text{ bayes rules} \end{aligned}$$

Now we take the log of MLE, then plug in  $q$ :

$$\begin{aligned} &\geq \max_{\theta, \rho} \sum_n [(y_n = +1) \log \rho + (y_n = -1) \log(1-\rho) + \sum_d [(x_{n,d} = 1) \log \theta_{y_n,d} + (x_{n,d} = 0) \log(1-\theta_{y_n,d})]] + \sum_m \\ &= \sum_{i=1}^{N+M} \sum_{y_i} q(y_i) = \begin{cases} 1 & \text{if } i \leq N, \text{ and } y_n = y_i \\ 0 & \text{if } i \leq N, \text{ and } y_n \neq y_i \text{ (} y_i = +1) \log \rho + (y_i = -1) \log(1-\rho) + \sum_d [(x_{i,d} = 1) \log \theta_{y_i,d} + (x_{i,d} = 0) \log(1-\theta_{y_i,d})] \\ p(y_i | x_i, \rho, \theta) & \text{if } i > N \end{cases} \end{aligned}$$

Let's take  $\frac{\partial}{\partial \rho}$  ( $\rho$  and  $\theta$  in  $q$  are constants here!)

$$\begin{aligned} &= \sum_i \sum_{y_i} q(y_i) \left[ \frac{y_i = +1}{\rho} - \frac{y_i = -1}{1-\rho} \right] \\ \sum_i \sum_{y_i} q(y_i) (y_i = +1) (1-\rho) &= \sum_i \sum_{y_i} q(y_i) (y_i = -1) (\rho) \\ \sum_i q(y_i = +1) - q(y_i = +1) \rho &= \rho \sum_i q(y_i = -1) \text{ distribute} \\ \sum_i q(y_i = +1) &= \rho (\sum_i q(y_i = -1) + q(y_i = +1)) \\ \sum_i q(y_i = +1) &= \rho \times (N + M) \\ \rho &= \frac{\sum_i q(y_i = +1)}{N + M} \end{aligned}$$

So the final update for

$$\rho = \frac{\sum_i q(y_i = +1)}{N + M}$$

## 19 Class 22? November 29th 2011

### 19.1 HMM

Reasoning over time: of a sequence of observations. Involves time into our models.

#### 19.1.1 Markov Models

- a chain-structured Bayes Net
- each node identically distributed
- We have  $P(X_1)$ , then we assume  $P(X_j|X_{j-1})$  uniformly
- Parameters: **transition probabilities**
- Conditional independence assumption: first order Markov property, that each time step only depends on the previous
- example:  $P(X_2 = \text{sun}) = P(X_2 = \text{sun}|X_1 = \text{sun})P(X_1 = \text{sun}) + P(X_2 = \text{sun}|X_1 = \text{rain})P(X_1 = \text{rain})$
- generalized  $P(X_t = \text{sun}) = \sum_{x_1 \dots x_{t-1}} P(x_1, \dots, x_{t-1}, \text{sun})$  (Stupid way, because there are  $t-1$  days, sum over all values of 1 to  $t-1$ .  $2^{t-1}$  values in the sum!)
- Nice thing is that we can apply the same chain rule and move sums around, do dynamic programming. If we look at one configuration  $P(X_1 = \text{sun})$ , when second day is rain or second day is sun, the rest is the same.
- **Mini-forward algorithm** draw in lattice. Say we know  $P(x_1)$  is known, then we also know that

$$P(x_t) = \sum_{x_{t-1}} P(x_t|x_{t-1})P(x_{t-1})$$

the forward simulation.

- So in the example,  $P(X_t = \text{sun}|X_{t-1} = \text{sun}) = 0.9$ ,  $P(X_t = \text{sun}|X_{t-1} = \text{rain}) = 0.1$ , so  $P(X_3 = \text{sun}) = P(X_2|X_1 = \text{sun})P(X_1 = \text{sun}) + P(X_2|X_1 = \text{rain})P(X_1 = \text{rain}) = 0.9(0.9) + 0.1(0.1) = 0.82$
- If we simulate the chain long enough, there's a stationary distribution exists. The states like to stay constant..

For us we don't care that stationary distribution exists.



### 19.1.2 HMM

Markov chains not so useful for robots/agents. Wherever I come from, the next distribution is always the same. With the stationary distribution, the robot will have no idea where he is. Completely independent of my initial configuration. This happens because the robot is not paying attention to his environments. You'd like it to be able to look around and see what's going on and adjust it's beliefs based on what it sees.

So HMM is defined by:

- initial distribution:  $P(X_1)$
- Transitions:  $P(X_T|X_{T-1})$
- Emissions:  $P(E|X)$

Conditional Independence assumptions:

- Future depends on past only by the time step before i.e. the present (1st order markov assumption)
- Current observation independent of all else given current state.

Notation:  $X$  is all possible configurations,  $x_t$  is a single value we've observed. Robot example:

- Assume we have current belief  $P(X|evidencetodate)$ , then after one time step,  $P(X_{t+1}|e_{1:t}) = \sum_{x_t} P(X_{t+1}|x_t)P(x_t|e_{1:t})$
- After a step, current belief:  $B'(X_{t+1}) = P(X_{t+1}|e_{1:t})$ , before observing the new world. Then

$$P(X_{t+1}|e_{1:t+1}) \propto P(e_{t+1}|X_{t+1})P(X_{t+1}|e_{1:t})$$

aka:

$$B(X_{t+1}) \propto P(e|X)B'(X_{t+1})$$

Does what I observe agree with what I think Where I am?

- Idea: beliefs reweighted by likelihood of evidence. We need to renormalize

The forward algorithm:

$$\begin{aligned} P(X_t|e_{1:t}) &\propto P(x_t, e_{1:t}) \\ &= \sum_{x_{t-1}} P(x_{t-1}, x_t, e_{1:t}) \\ &= \sum_{x_{t-1}} P(x_{t-1}, e_{1:t-1})P(x_t|x_{t-1})P(e_t|x_t) \text{ chain rule} \\ &= P(e_t|x_t) \sum_{x_{t-1}} P(x_t|x_{t-1})P(x_{t-1}, e_{1:t-1}) \text{ reorganized} \end{aligned}$$

Main thing is to compute my belief at time  $t$ , I need to only sum over values at  $t-1$ . (quadratic in the number of states, compute x many state beliefs, then also compute what state was yesterday = x options) Not good when the number of states  $\mathcal{L}$  or states are continuous.

Alternative: **Particle Filtering**

- Discrete approximation of forward algorithm
- Track samples of  $X$ , not all values. Time per step now is linear in the number of samples.
- For a grid of states, use  $n$  balls (num samples). Number of balls in each bin is the approximation to the probability of picking that bin (grid).
- Each particle is moved by sampling its next position from the transition model

$$x' = \text{sample}(P(X'|x))$$

. He's at  $x$ , where does he think he's going to be next time step?

- Instead of evaluating the probability, pick at random weighted by the probability. Moving the particles according to the transitions. Each particles move independently, so linear in the number of particles/samples
- Observation: right after the move we have  $B'(X)$ . We need to get these particles to account for what they actually observe. Like likelihood weighting, so we downweight our samples based on the evidence - reweight.
- if a particle can explain the observation well, we gets more, less otherwise.  
 $w(x) = P(e|x)$

$$B(X) \propto P(e|X)B'(X)$$

- Probabilities don't sum to one! sum to  $P(e)$ , the approximation
- This is repetition of move-resize, in the end one particle will dominate, that happens to get moved in the right way. To alleviate, we can **resampling**:
  - Think of moved particles as distribution of states. Then resample based on this distribution, all with weight 1. Now probability sums to one (equivalent to renormalizing).
  - repeat: Move, resize, resample.

## 20 Most-likely Question

Question: most likely sequence ending in  $x$  at  $t$ ? We don't have the observation (more like MM setting, not HMM). Slow answer: enumerate and score.

$$\max_{x_1, \dots, x_{t-1}} P(x_1, \dots, x_{t-1}, \text{sun})$$

$2^{t-1}$  state sequences. Previously, we were summing over all of this, not we max over them because we want the most likely state.

## 20.1 Mini-viterbi algorithm

(mini because it's for MM) What's the probability of most likely path that goes to this state, and in that sequence where did I come from? define:

$$m_t[x] = \max_{x_{1:t-1}} P(x_{1:t-1}, x_t)$$

$$a_t[x] = \operatorname{argmax}_{x_{1:t-1}} P(x_{1:t-1}, x)$$

$m$ 's are the probability,  $a$ 's are the back pointers.

Recursively,

$$\begin{aligned} m_t[x] &= \max_{x_{1:t-1}} P(x_{1:t-1}, x) \\ &= \max_{x_{1:t-1}} P(x_{1:t-1})P(x_t|x_{t-1}) \\ &= \max_{x_{t-1}} P(x|x_{t-1}) \max_{x_{1:t-2}} \text{ if the most likely sequence goes through } x_{t-1}, \text{ ask how we got to } x_{t-1} \\ &= \max_{x_{t-1}} P(x|x_{t-1})m_{t-1}[x_{t-1}] \end{aligned}$$

Where  $m_1[x] = P(x_1)$ .

## 20.2 Viterbi Algorithm

Q:What is the most likely state sequence given the observations?. Same thing but we need to factor in the observation probabilities. Now

$$m_t[x] = \max_{x_{1:t-1}} P(x_{1:t-1}, x_t, e_{1:t})$$

Same story:

$$\begin{aligned} m_t[x] &= \max_{x_{1:t-1}} P(x_{1:t-1}, x_t, e_{1:t}) \\ &= \max_{x_{1:t-1}} P(x_{1:t-1}, e_{1:t-1})P(x_t|x_{t-1})P(e_t|x_t) \\ &= P(e_t|x_t) \max_{t-1} P(x_t|x_{t-1}) \max_{x_{1:t-2}} P(x_{1:t-1}, e_{1:t-1}) \\ &= P(e_t|x_t) \max_{t-1} P(x_t|x_{t-1})m_{t-1}[x_{t-1}] \end{aligned}$$