

CPSC 340: Machine Learning and Data Mining

Non-Parametric Models

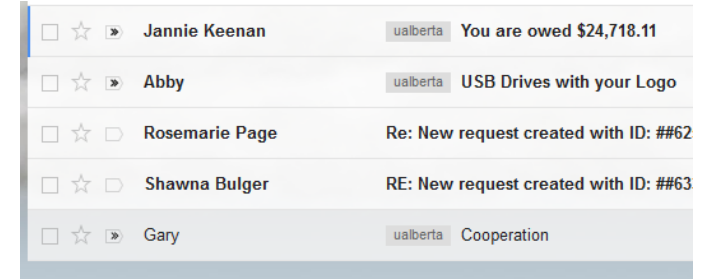
Fall 2019

Admin

- Course webpage:
 - <https://www.cs.ubc.ca/~schmidtm/Courses/340-F19/>
- **Assignment 1:**
 - 1 late day to hand in tonight, 2 for Wednesday.
- **Assignment 2** is out.
 - Due Friday of next week. It's long so start early.
- **Add/drop deadline** is tomorrow.
- **Auditing/Exchange:**
 - Bring your form to me after class.

Last Time: E-mail Spam Filtering

- Want to build a system that filters spam e-mails:
- We formulated as **supervised learning**:
 - $(y_i = 1)$ if e-mail 'i' is spam, $(y_i = 0)$ if e-mail is not spam.
 - $(x_{ij} = 1)$ if word/phrase 'j' is in e-mail 'i', $(x_{ij} = 0)$ if it is not.



\$	Hi	CPSC	340	Vicodin	Offer	...	Spam?
1	1	0	0	1	0	...	1
0	0	0	0	1	1	...	1
0	1	1	1	0	0	...	0
...

Last Time: Naïve Bayes

- We considered spam filtering methods based on **naïve Bayes**:

$$p(y_i = \text{"spam"} \mid x_i) = \frac{p(x_i \mid y_i = \text{"spam"}) p(y_i = \text{"spam"})}{p(x_i)}$$

- Makes **conditional independence** assumption to make learning practical:

$$p(\underbrace{\text{hello}=1, \text{vicodin}=0, \text{340}=1}_{\text{HARD}} \mid \text{spam}) \approx \underbrace{p(\text{hello}=1 \mid \text{spam})}_{\text{easy}} \underbrace{p(\text{vicodin}=0 \mid \text{spam})}_{\text{easy}} \underbrace{p(\text{340}=1 \mid \text{spam})}_{\text{easy}}$$

- Predict “spam” if $p(y_i = \text{"spam"} \mid x_i) > p(y_i = \text{"not spam"} \mid x_i)$.
 - We don’t need $p(x_i)$ to test this.

Naïve Bayes

- Naïve Bayes formally:

$$p(y_i | x_i) = \frac{p(x_i | y_i) p(y_i)}{p(x_i)} \quad (\text{first use Bayes rule})$$

$$\propto p(x_i | y_i) p(y_i) \quad (\text{"denominator doesn't matter"})$$

same for all y_i values

$$\approx \prod_{j=1}^d \left[p(x_{ij} | y_i) \right] p(y_i) \quad (\text{conditional independence assumption})$$

Only needs easy probabilities.

- Post-lecture slides: **how to train/test by hand** on a simple example.

Laplace Smoothing

- Our estimate of $p(\text{'lactase'} = 1 \mid \text{'spam'})$ is:

$$\frac{\# \text{spam messages with lactase}}{\# \text{spam messages}}$$

- But there is a problem if you have **no spam messages with lactase**:

- $p(\text{'lactase'} \mid \text{'spam'}) = 0$, so spam messages with lactase automatically get through.

- Common fix is **Laplace smoothing**:
$$\frac{(\# \text{spam messages with lactase}) + 1}{(\# \text{spam messages}) + 2}$$
 - **Add 1 to numerator**,
and 2 to denominator (for binary features).

- Acts like a “fake” spam example that has lactase,
and a “fake” spam example that doesn’t.

Laplace Smoothing

- Laplace smoothing:
$$\frac{(\text{\#spam messages with lactase}) + 1}{(\text{\#spam messages}) + 2}$$
 - Typically you **do this for all features**.
 - Helps against overfitting by biasing towards the uniform distribution.
- A common variation is to use a **real number β** rather than 1.
 - Add ' **βk** ' to **denominator** if feature has 'k' possible values (so it sums to 1).

$$p(x_{ij}=c | y_i=\text{class}) \approx \frac{(\text{number of examples in class with } x_{ij}=c) + \beta}{(\text{number of examples in class}) + \beta K}$$

This is a “**maximum a posteriori**” (MAP) estimate of the probability. We’ll discuss MAP and how to derive this formula later.

Decision Theory

- Are we **equally concerned about “spam” vs. “not spam”**?
- **True positives, false positives, false negatives, true negatives:**

Predict / True	True 'spam'	True 'not spam'
Predict 'spam'	True Positive	False Positive
Predict 'not spam'	False Negative	True Negative

- The costs mistakes might be different:
 - Letting a spam message through (false negative) is not a big deal.
 - Filtering a not spam (false positive) message will make users mad.

Decision Theory

- We can give a **cost** to each scenario, such as:

Predict / True	True 'spam'	True 'not spam'
Predict 'spam'	0	100
Predict 'not spam'	10	0

- Instead of most probable label, take \hat{y}_i **minimizing expected cost**:

$$\mathbb{E}[\text{cost}(\hat{y}_i, \tilde{y}_i)]$$

expectation of model with respect to \tilde{y}_i

cost of predicting \hat{y}_i if it's really \tilde{y}_i

- Even if “spam” has a higher probability, predicting “spam” might have a expected higher cost.

Decision Theory Example

Predict / True	True 'spam'	True 'not spam'
Predict 'spam'	0	100
Predict 'not spam'	10	0

- Consider a test example we have $p(\tilde{y}_i = \text{"spam"} \mid \tilde{x}_i) = 0.6$, then:

$$\begin{aligned} E[\text{cost}(\hat{y}_i = \text{"spam"}, \tilde{y}_i)] &= p(\tilde{y}_i = \text{"spam"} \mid \tilde{x}_i) \text{cost}(\hat{y}_i = \text{"spam"}, \tilde{y}_i = \text{"spam"}) \\ &\quad + p(\tilde{y}_i = \text{"not spam"} \mid \tilde{x}_i) \text{cost}(\hat{y}_i = \text{"spam"}, \tilde{y}_i = \text{"not spam"}) \\ &= (0.6)(0) + (0.4)(100) = 40 \end{aligned}$$

$$E[\text{cost}(\hat{y}_i = \text{"not spam"}, \tilde{y}_i)] = (0.6)(10) + (0.4)(0) = 6$$

- Even though "spam" is more likely, we should predict "not spam".

Naïve Bayes Training Phase

- Training a naïve Bayes model:

$$X = \begin{bmatrix} 0 & 1 \\ 1 & 1 \\ 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 1 \\ 1 & 0 \end{bmatrix}, \quad y = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

Naïve Bayes Training Phase

- Training a naïve Bayes model:

1. Set n_c to the number of times $(y_i = c)$.

$$X = \begin{bmatrix} 0 & 1 \\ 1 & 1 \\ 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 1 \\ 1 & 0 \end{bmatrix}, \quad y = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$n_1 = 6$



$n_0 = 4$



Naïve Bayes Training Phase

- Training a naïve Bayes model:

1. Set n_c to the number of times $(y_i = c)$.
2. Estimate $p(y_i = c)$ as $\frac{n_c}{n}$.

$$X = \begin{bmatrix} 0 & 1 \\ 1 & 1 \\ 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 1 \\ 1 & 0 \end{bmatrix}, \quad y = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$p(y_i = 1) = \frac{6}{10} \leftarrow n_1 = 6$

$p(y_i = 0) = \frac{4}{10} \leftarrow n_0 = 4$

Naïve Bayes Training Phase

- Training a naïve Bayes model:

1. Set n_c to the number of times $(y_i = c)$.
2. Estimate $p(y_i = c)$ as $\frac{n_c}{n}$.
3. Set n_{cjk} as the number of times $(y_i = c, x_{ij} = k)$

$p(y_i = 1) = \frac{6}{10} \leftarrow n_1 = 6$

$X = \begin{bmatrix} 0 & 1 & 1 \\ 1 & 1 & 1 \\ 0 & 0 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix}, y = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$

$n_{121} = 4$

$p(y_i = 0) = \frac{4}{10} \leftarrow n_0 = 4$

Naïve Bayes Training Phase

- Training a naïve Bayes model:

1. Set n_c to the number of times $(y_i = c)$.
2. Estimate $p(y_i = c)$ as $\frac{n_c}{n}$.
3. Set n_{cjk} as the number of times $(y_i = c, x_{ij} = k)$.
4. Estimate $p(x_{ij} = k, y_i = c)$ as $\frac{n_{cjk}}{n}$.

$p(y_i = 1) = \frac{6}{10} \leftarrow n_1 = 6$

$X = \begin{bmatrix} 0 & 1 & 1 \\ 1 & 1 & 1 \\ 0 & 0 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix}, y = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$

$n_{121} = 4$

$p(x_{12} = 1, y_i = 1) = \frac{4}{10}$

$p(y_i = 0) = \frac{4}{10} \leftarrow n_0 = 4$

Naïve Bayes Training Phase

- Training a naïve Bayes model:

1. Set n_c to the number of times $(y_i = c)$.
2. Estimate $p(y_i = c)$ as $\frac{n_c}{n}$.
3. Set n_{cjk} as the number of times $(y_i = c, x_{ij} = k)$.
4. Estimate $p(x_{ij} = k, y_i = c)$ as $\frac{n_{cjk}}{n}$.
5. Use that $p(x_{ij} = k | y_i = c) = \frac{p(x_{ij} = k, y_i = c)}{p(y_i = c)}$

$$= \frac{n_{cjk}/n}{n_c/n} = \frac{n_{cjk}}{n_c}$$

$$p(x_{12}=1 | y_i=1) = \frac{4}{6} = \frac{2}{3}$$

$$p(x_{12}=1, y_i=1) = \frac{4}{10}$$

$$p(y_i=0) = \frac{4}{10} \leftarrow n_0=4$$

$$p(y_i=1) = \frac{6}{10} \leftarrow n_1=6$$

$X =$

0	1	1
1	1	1
0	0	1
1	1	1
1	1	1
0	0	1
1	0	0
1	0	0
1	1	0
1	0	0

$y =$

1
1
1
1
1
1
0
0
0
0

$$n_{121}=4$$

$$n_0=4$$

Naïve Bayes Prediction Phase

- Prediction in a naïve Bayes model:

Given a test example \tilde{x}_i we set prediction \hat{y}_i to the 'c' maximizing $p(\tilde{x}_i | \tilde{y}_i = c)$

Under the naïve Bayes assumption we can maximize:

$$p(\tilde{y}_i = c | \tilde{x}_i) \propto \prod_{j=1}^d [p(\tilde{x}_{ij} | \tilde{y}_i = c)] p(\tilde{y}_i = c)$$

Naïve Bayes Prediction Phase

- Prediction in a naïve Bayes model:

Consider $\tilde{x}_i = [1 \ 1]$ in this data set \rightarrow

$$X = \begin{bmatrix} 0 & 1 \\ 1 & 1 \\ 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 1 \\ 1 & 0 \end{bmatrix}, \quad y = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

Naïve Bayes Prediction Phase

- Prediction in a naïve Bayes model:

Consider $\tilde{x}_i = [1 \ 1]$ in this data set \rightarrow

$$p(\tilde{y}_i = 0 | \tilde{x}_i) \propto \underbrace{p(\tilde{x}_{i1} = 1 | \tilde{y}_i = 0)}_{(1)} \underbrace{p(\tilde{x}_{i2} = 1 | \tilde{y}_i = 0)}_{(0.25)} \underbrace{p(\tilde{y}_i = 0)}_{(0.4)} = 0.1$$

$$X = \begin{bmatrix} 0 & 1 \\ 1 & 1 \\ 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 1 \\ 1 & 0 \end{bmatrix}, \quad y = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

Naïve Bayes Prediction Phase

- Prediction in a naïve Bayes model:

Consider $\tilde{x}_i = [1 \ 1]$ in this data set \rightarrow

$$p(\tilde{y}_i = 0 \mid \tilde{x}_i) \propto p(\tilde{x}_{i1} = 1 \mid \tilde{y}_i = 0) p(\tilde{x}_{i2} = 1 \mid \tilde{y}_i = 0) p(\tilde{y}_i = 0)$$

$$= (1) (0.25) (0.4) = 0.1$$

$$p(\tilde{y}_i = 1 \mid \tilde{x}_i) \propto p(\tilde{x}_{i1} = 1 \mid \tilde{y}_i = 1) p(\tilde{x}_{i2} = 1 \mid \tilde{y}_i = 1) p(\tilde{y}_i = 1)$$

$$= (0.5) (0.666...) (0.6) = 0.2$$

$$X = \begin{bmatrix} 0 & 1 \\ 1 & 1 \\ 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 1 \\ 1 & 0 \end{bmatrix}, \quad y = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

Naïve Bayes Prediction Phase

- Prediction in a naïve Bayes model:

Consider $\tilde{x}_i = [1 \ 1]$ in this data set \rightarrow

$$p(\tilde{y}_i = 0 \mid \tilde{x}_i) \propto p(\tilde{x}_{i1} = 1 \mid \tilde{y}_i = 0) p(\tilde{x}_{i2} = 1 \mid \tilde{y}_i = 0) p(\tilde{y}_i = 0) \\ = (1) (0.25) (0.4) = 0.1$$

$$X = \begin{bmatrix} 0 & 1 \\ 1 & 1 \\ 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 1 \\ 1 & 0 \end{bmatrix}, \quad y = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$p(\tilde{y}_i = 1 \mid \tilde{x}_i) \propto p(\tilde{x}_{i1} = 1 \mid \tilde{y}_i = 1) p(\tilde{x}_{i2} = 1 \mid \tilde{y}_i = 1) p(\tilde{y}_i = 1) \\ = (0.5) (0.666...) (0.6) = 0.2$$

Since $p(\tilde{y}_i = 1 \mid \tilde{x}_i)$ is bigger than $p(\tilde{y}_i = 0 \mid \tilde{x}_i)$, naïve Bayes predicts $\hat{y}_i = 1$

(Don't sum to 1 because we're ignoring $p(\tilde{x}_i)$)

“Proportional to” for Probabilities

- When we say “ $p(y) \propto \exp(-y^2)$ ” for a function ‘p’, we mean:

$$p(y) = \beta \exp(-y^2) \text{ for some constant } \beta.$$

- However, if ‘p’ is a probability then it must sum to 1.

– If $y \in \{1,2,3,4\}$ then $p(1) + p(2) + p(3) + p(4) = 1$

- Using this fact, we can find β :

$$\begin{aligned} & \beta \exp(-1^2) + \beta \exp(-2^2) + \beta \exp(-3^2) + \beta \exp(-4^2) = 1 \\ \Leftrightarrow & \beta [\exp(-1^2) + \exp(-2^2) + \exp(-3^2) + \exp(-4^2)] = 1 \\ \Leftrightarrow & \beta = \frac{1}{\exp(-1^2) + \exp(-2^2) + \exp(-3^2) + \exp(-4^2)} \end{aligned}$$

Probability of Paying Back a Loan and Ethics

- Article discussing predicting “whether someone will pay back a loan”:
 - <https://www.thecut.com/2017/05/what-the-words-you-use-in-a-loan-application-reveal.html>
- Words that **increase probability** of paying back the most:
 - *debt-free, lower interest rate, after-tax, minimum payment, graduate.*
- Words that **decrease probability** of paying back the most:
 - *God, promise, will pay, thank you, hospital.*
- Article also discusses an important issue: **are all these features ethical?**
 - Should you deny a loan because of religion or a family member in the hospital?
 - ICBC is limited in the features it is allowed to use for prediction.

Avoiding Underflow

- During the prediction, the **probability can underflow**:

$$p(y_i = c | x_i) \propto \prod_{j=1}^d [p(x_{ij} | y_i = c)] p(y_i = c)$$

→ All these are < 1 so the product gets very small!

- Standard fix is to (equivalently) maximize the logarithm of the probability:

Remember that $\log(ab) = \log(a) + \log(b)$ so $\log(\prod a_i) = \sum \log(a_i)$

Since \log is monotonic the 'c' maximizing $p(y_i = c | x_i)$ also maximizes $\log p(y_i = c | x_i)$,

so maximize $\log\left(\prod_{j=1}^d [p(x_{ij} | y_i = c)] p(y_i = c)\right) = \sum_{j=1}^d \log(p(x_{ij} | y_i = c)) + \log(p(y_i = c))$

Less-Naïve Bayes

- Given features $\{x_1, x_2, x_3, \dots, x_d\}$, naïve Bayes approximates $p(y|x)$ as:

$$\begin{aligned} p(y | x_1, x_2, \dots, x_d) &\propto p(y) p(x_1, x_2, \dots, x_d | y) \quad \text{product rule applied repeatedly} \\ &= p(y) p(x_1 | y) p(x_2 | x_1, y) p(x_3 | x_2, x_1, y) \dots p(x_d | x_1, x_2, \dots, x_{d-1}, y) \\ &\approx p(y) p(x_1 | y) p(x_2 | y) p(x_3 | y) \dots p(x_d | y) \quad (\text{naïve Bayes assumption}) \end{aligned}$$

- The assumption is very strong, and there are “less naïve” versions:
 - Assume independence of all variables except up to ‘k’ largest ‘j’ where $j < i$.
 - E.g., naïve Bayes has $k=0$ and with $k=2$ we would have:

$$\approx p(y) p(x_1 | y) p(x_2 | x_1, y) p(x_3 | x_2, x_1, y) p(x_4 | x_3, x_2, y) \dots p(x_d | x_{d-2}, x_{d-1}, y)$$

- Fewer independence assumptions so more flexible, but hard to estimate for large ‘k’.
 - Another practical variation is “tree-augmented” naïve Bayes.

Computing $p(x_i)$ under naïve Bayes

- **Generative models** don't need $p(x_i)$ to make decisions.
- However, it's **easy to calculate** under the naïve Bayes assumption:

$$p(x_i) = \sum_{c=1}^K p(x_i, y=c) \quad (\text{marginalization rule})$$

$$= \sum_{c=1}^K p(x_i | y=c) p(y=c) \quad (\text{product rule})$$

$$= \sum_{c=1}^K \left[\prod_{j=1}^d p(x_{ij} | y=c) \right] p(y=c) \quad (\text{naïve Bayes assumption})$$

These are the quantities
we compute during training

Gaussian Discriminant Analysis

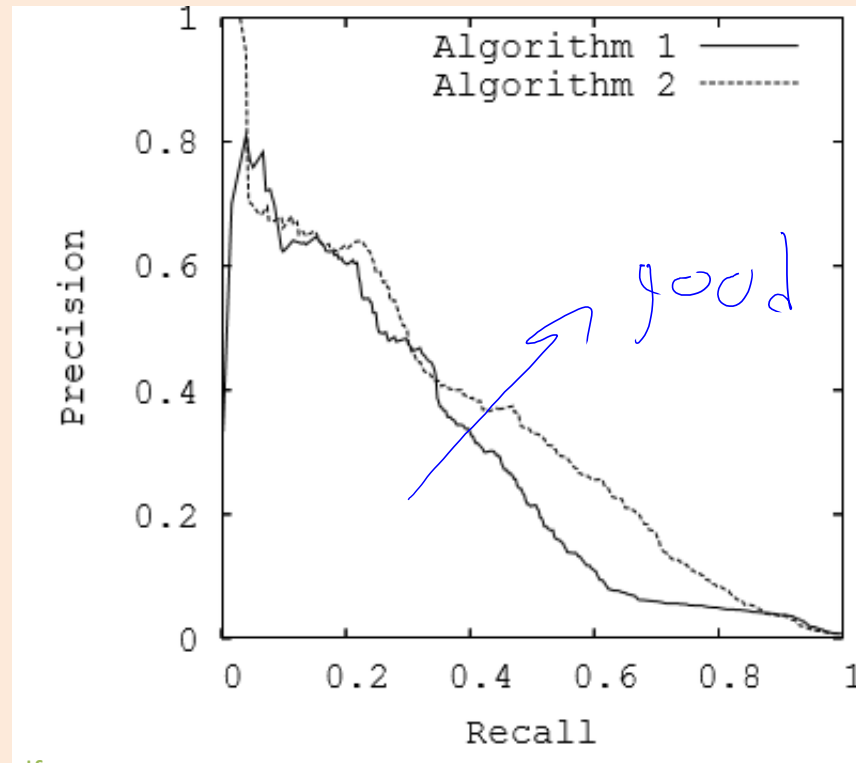
- Classifiers based on Bayes rule are called **generative classifier**:
 - They often work well when you have **tons of features**.
 - But they **need to know $p(x_i | y_i)$** , **probability of features given the class**.
 - How to “generate” features, based on the class label.
- To fit generative models, usually make BIG assumptions:
 - **Naïve Bayes** (NB) for discrete x_i :
 - Assume that each variables in x_i is independent of the others in x_i given y_i .
 - **Gaussian discriminant analysis** (GDA) for continuous x_i .
 - Assume that $p(x_i | y_i)$ follows a multivariate normal distribution.
 - If all classes have same covariance, it’s called “linear discriminant analysis”.

Other Performance Measures

- Classification error might be wrong measure:
 - Use weighted classification error if have different costs.
 - Might want to use things like Jaccard measure: $TP/(TP + FP + FN)$.
- Often, we report **precision** and **recall** (want both to be high):
 - Precision: “if I classify as spam, what is the probability it actually is spam?”
 - Precision = $TP/(TP + FP)$.
 - High precision means the filtered messages are likely to really be spam.
 - Recall: “if a message is spam, what is probability it is classified as spam?”
 - Recall = $TP/(TP + FN)$
 - High recall means that most spam messages are filtered.

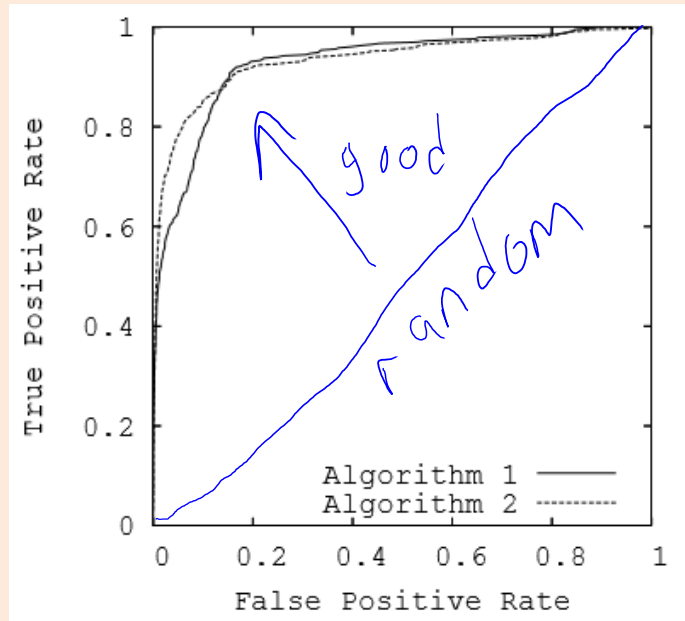
Precision-Recall Curve

- Consider the rule $p(y_i = \text{'spam'} \mid x_i) > t$, for threshold 't'.
- Precision-recall (PR) curve plots precision vs. recall as 't' varies.



ROC Curve

- Receiver operating characteristic (ROC) curve:
 - Plot true positive rate (recall) vs. false positive rate (FP/FP+TN).
(negative examples classified as positive)



- Diagonal is random, perfect classifier would be in upper left.
- Sometimes papers report area under curve (AUC).
 - Reflects performance for different possible thresholds on the probability.