

Artificial Intelligence

CE-417, Group 2

Computer Eng. Department

Sharif University of Technology

Spring 2020

By Mohammad Hossein Rohban, Ph.D.

Courtesy: Most slides are adopted from CSE-573 (Washington U.), original
slides for the textbook, and CS-188 (UC. Berkeley).

Machine Learning

Key Concepts

Machine Learning

- Up until now: how use a **model** to make optimal decisions.
 - **Machine learning:** how to acquire a model from data / experience
 - Learning parameters (e.g. Probabilities)
 - Learning structure (e.g. BN graphs)
 - Learning hidden concepts (e.g. Clustering)
 - Today: model-based **classification** with **naive Bayes**
- DATA



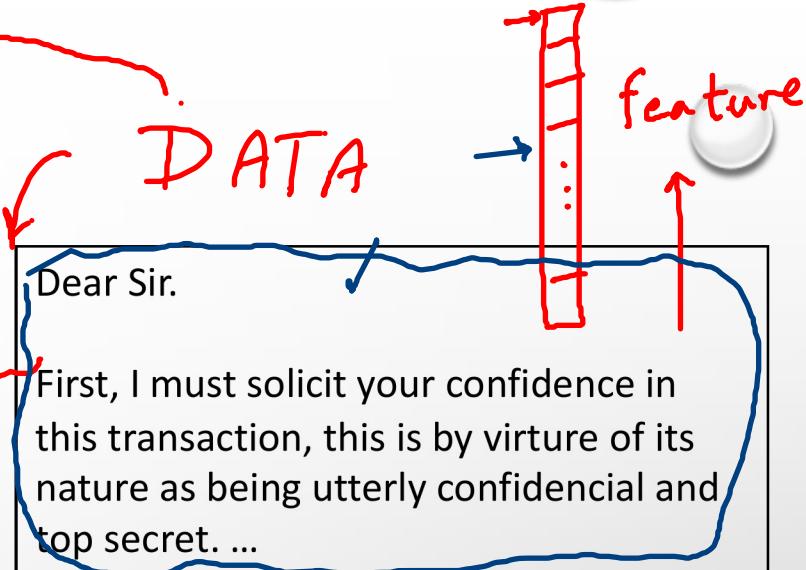
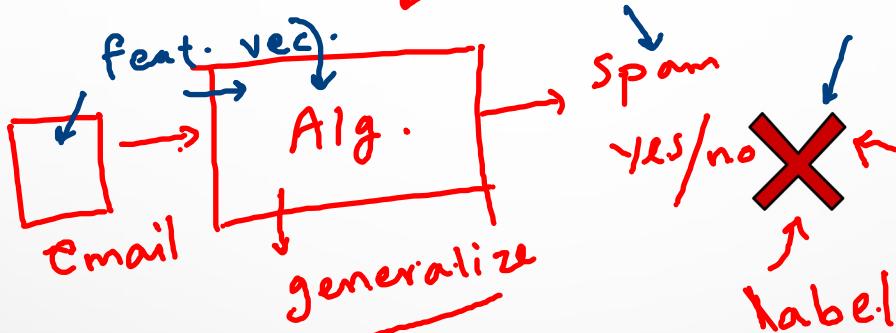
- Input: an email
- Output: spam/ham

- Setup:
 - Get a large collection of example emails, each labeled "spam" or "ham"
 - Note: someone has to hand label all this data!
 - Want to learn to predict labels of new, future emails

- Features: the attributes used to make the ham / spam decision

- Words: FREE!
- Text patterns: \$dd, CAPS
- Non-text: SenderInContacts
- ...

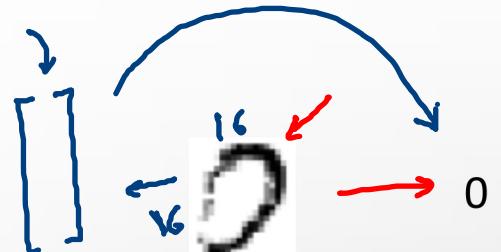
Example: Spam Filter



Example: Digit Recognition

- Input: images / pixel grids
- Output: a digit 0-9

- Setup:
 - Get a large collection of example images, each labeled with a digit
 - Note: someone has to hand label all this data!
 - Want to learn to predict labels of new, future digit images
- Features: the attributes used to make the digit decision
 - Pixels: (6,8)=ON
 - Shape patterns: NumComponents, AspectRatio, NumLoops
 - ...

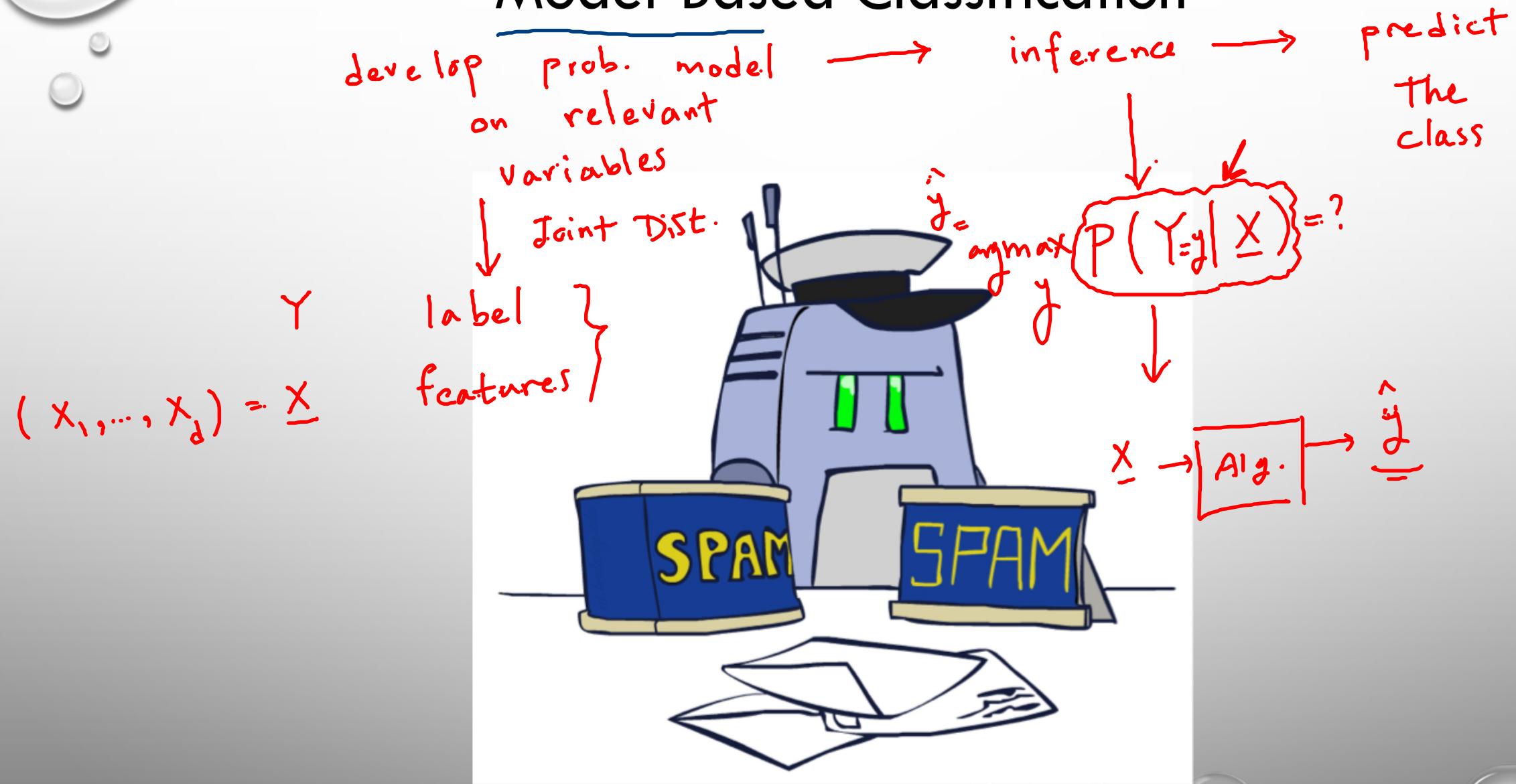


Other Classification Tasks

- Classification: given inputs x , predict labels (classes) y
- Examples:
 - Spam detection (input: document,
Classes: spam / ham)
 - OCR (input: images, classes: characters)
 - Medical diagnosis (input: symptoms,
Classes: diseases)
 - Automatic essay grading (input: document,
Classes: grades)
 - Fraud detection (input: account activity,
Classes: fraud / no fraud)
 - Customer service email routing
 - ... Many more
- Classification is an important commercial technology!



Model-Based Classification



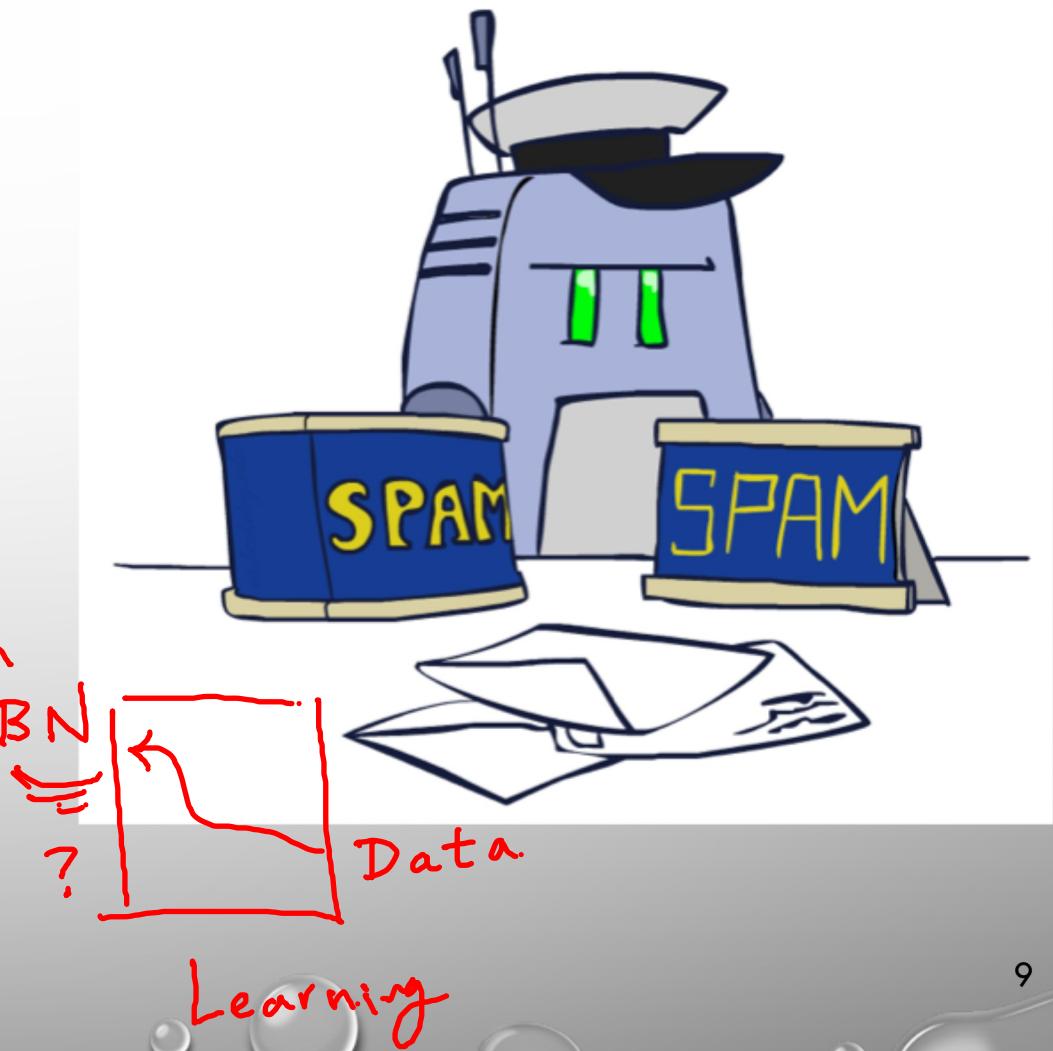
Model-Based Classification

- Model-based approach
 - Build a model (e.g. Bayes' net) where both the label and features are random variables
 - Instantiate any observed features
 - Query for the distribution of the label conditioned on the features

$$\text{pred} \leftarrow P(y | \underline{x}) \leftarrow P(\underline{x}, y)$$

- Challenges

- What structure should the BN have?
- How should we learn its parameters?



Naïve Bayes for Digits

- Naïve Bayes: assume all features are independent effects of the label

- Simple digit recognition version:

- One feature (variable) f_{ij} for each grid position $\langle i,j \rangle$
- Feature values are on / off, based on whether intensity is more or less than 0.5 in underlying image
- Each input maps to a feature vector, e.g.

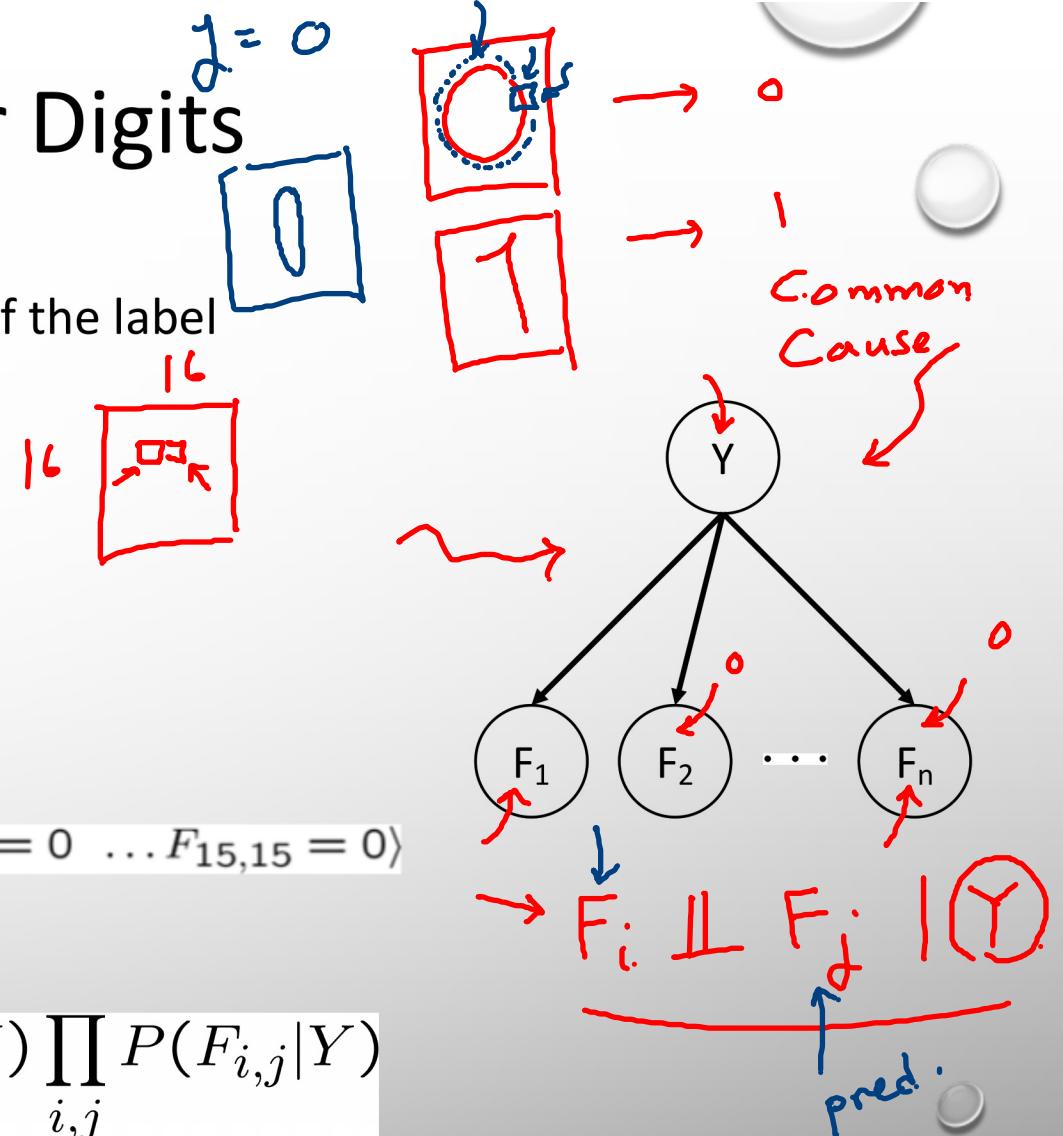


$\rightarrow \langle F_{0,0} = 0 \ F_{0,1} = 0 \ F_{0,2} = 1 \ F_{0,3} = 1 \ F_{0,4} = 0 \ \dots \ F_{15,15} = 0 \rangle$

- Here: lots of features, each is binary valued

- Naïve Bayes model:
$$P(Y|F_{0,0} \dots F_{15,15}) \propto P(Y) \prod_{i,j} P(F_{i,j}|Y)$$

- What do we need to learn?



General Naïve Bayes

- A general **naive Bayes** model:

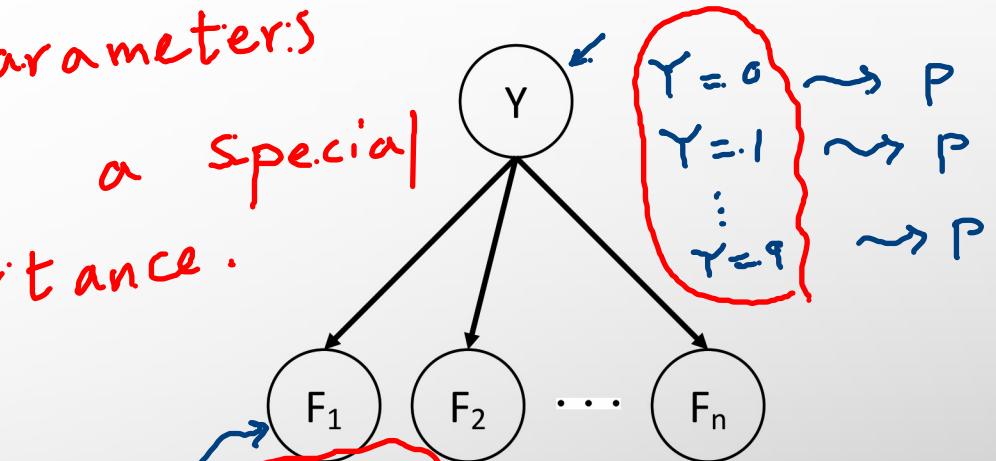
$$P(Y, F_1 \dots F_n) = P(Y) \prod_i P(F_i|Y)$$

$|Y| \times |F|^n$ values

$|Y|$ parameters

$n \times |F| \times |Y|$ parameters

Reduction of free parameters is of a special importance.



$$P(F_i|Y)$$

→ generalization

- We only have to specify how each feature depends on the class
- Total number of parameters is **linear in n**
- Model is very simplistic, but often works anyway

need $\sim \#$ of data points.

Inference for Naïve Bayes

- Goal: compute posterior distribution over label variable Y
 - Step 1: get joint probability of label and evidence for each label

$$P(Y, f_1 \dots f_n) = \begin{bmatrix} P(y_1, f_1 \dots f_n) \\ P(y_2, f_1 \dots f_n) \\ \vdots \\ P(y_k, f_1 \dots f_n) \end{bmatrix} \xrightarrow{\text{BN}} \frac{\begin{bmatrix} P(y_1) \prod_i P(f_i|y_1) \\ P(y_2) \prod_i P(f_i|y_2) \\ \vdots \\ P(y_k) \prod_i P(f_i|y_k) \end{bmatrix}}{P(f_1 \dots f_n)}$$

- Step 2: sum to get probability of evidence

$$\hat{y} = \underset{y}{\operatorname{argmax}}$$

- Step 3: normalize by dividing step 1 by step 2

$$P(Y_{=\hat{y}} | f_1 \dots f_n)$$

General Naïve Bayes

- What do we need in order to use naïve Bayes?



- Inference method (we just saw this part)

- Start with a bunch of probabilities: $P(Y)$ and the $P(F_i | Y)$ tables
- Use standard inference to compute $P(Y | F_1 \dots F_n)$
- Nothing new here

CPT ?

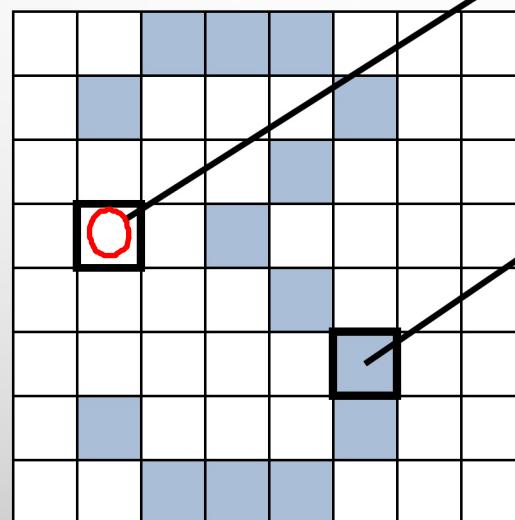
- Estimates of local conditional probability tables

- $P(Y)$, the prior over labels
- $P(F_i | Y)$ for each feature (evidence variable)
- These probabilities are collectively called the *parameters* of the model and denoted by θ
- Up until now, we assumed these appeared by magic, but...
- ...They typically come from training data counts: we'll look at this soon

Example: Conditional Probabilities

$P(Y)$

| | |
|---|-----|
| 1 | 0.1 |
| 2 | 0.1 |
| 3 | 0.1 |
| 4 | 0.1 |
| 5 | 0.1 |
| 6 | 0.1 |
| 7 | 0.1 |
| 8 | 0.1 |
| 9 | 0.1 |
| 0 | 0.1 |



$P(F_{3,1} = \text{on}|Y) \quad P(F_{5,5} = \text{on}|Y)$

| | |
|---|------|
| 1 | 0.01 |
| 2 | 0.05 |
| 3 | 0.05 |
| 4 | 0.30 |
| 5 | 0.80 |
| 6 | 0.90 |
| 7 | 0.05 |
| 8 | 0.60 |
| 9 | 0.50 |
| 0 | 0.80 |

| | |
|---|------|
| 1 | 0.05 |
| 2 | 0.01 |
| 3 | 0.90 |
| 4 | 0.80 |
| 5 | 0.90 |
| 6 | 0.90 |
| 7 | 0.25 |
| 8 | 0.85 |
| 9 | 0.60 |
| 0 | 0.80 |

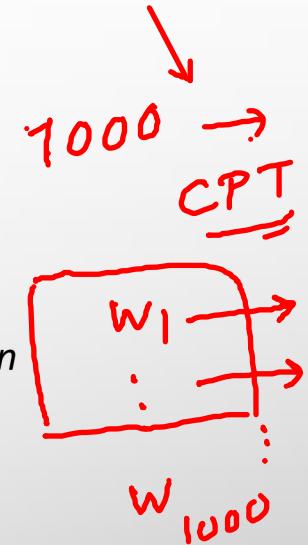
Naïve Bayes for Text

- Bag-of-words naïve Bayes:

- Features: W_i is the word at position i
- As before: predict label conditioned on feature variables (spam vs. Ham)
- As before: assume features are conditionally independent given label
- New: each W_i is identically distributed

$$P(Y, W_1 \dots W_n) = P(Y) \prod_i P(W_i|Y)$$

Word at position
 i , not i^{th} word in
the dictionary!



- Generative model:

- “Tied” distributions and bag-of-words

- Usually, each variable gets its own conditional probability distribution $P(F|Y)$
- In a bag-of-words model
 - Each position is identically distributed
 - All positions share the same conditional probs. $P(W|Y)$
 - Why make this assumption?
- Called “bag-of-words” because model is insensitive to word order or reordering

Example: Spam Filtering

- Model: $P(Y, W_1 \dots W_n) = P(Y) \prod_i P(W_i|Y)$
- What are the parameters?

$P(Y)$ ↗

| |
|------------|
| ham : 0.66 |
| spam: 0.33 |

$P(W|\text{spam})$ ↗

| |
|--------------|
| the : 0.0156 |
| to : 0.0153 |
| and : 0.0115 |
| of : 0.0095 |
| you : 0.0093 |
| a : 0.0086 |
| with: 0.0080 |
| from: 0.0075 |
| ... : |

$P(W|\text{ham})$ ↗

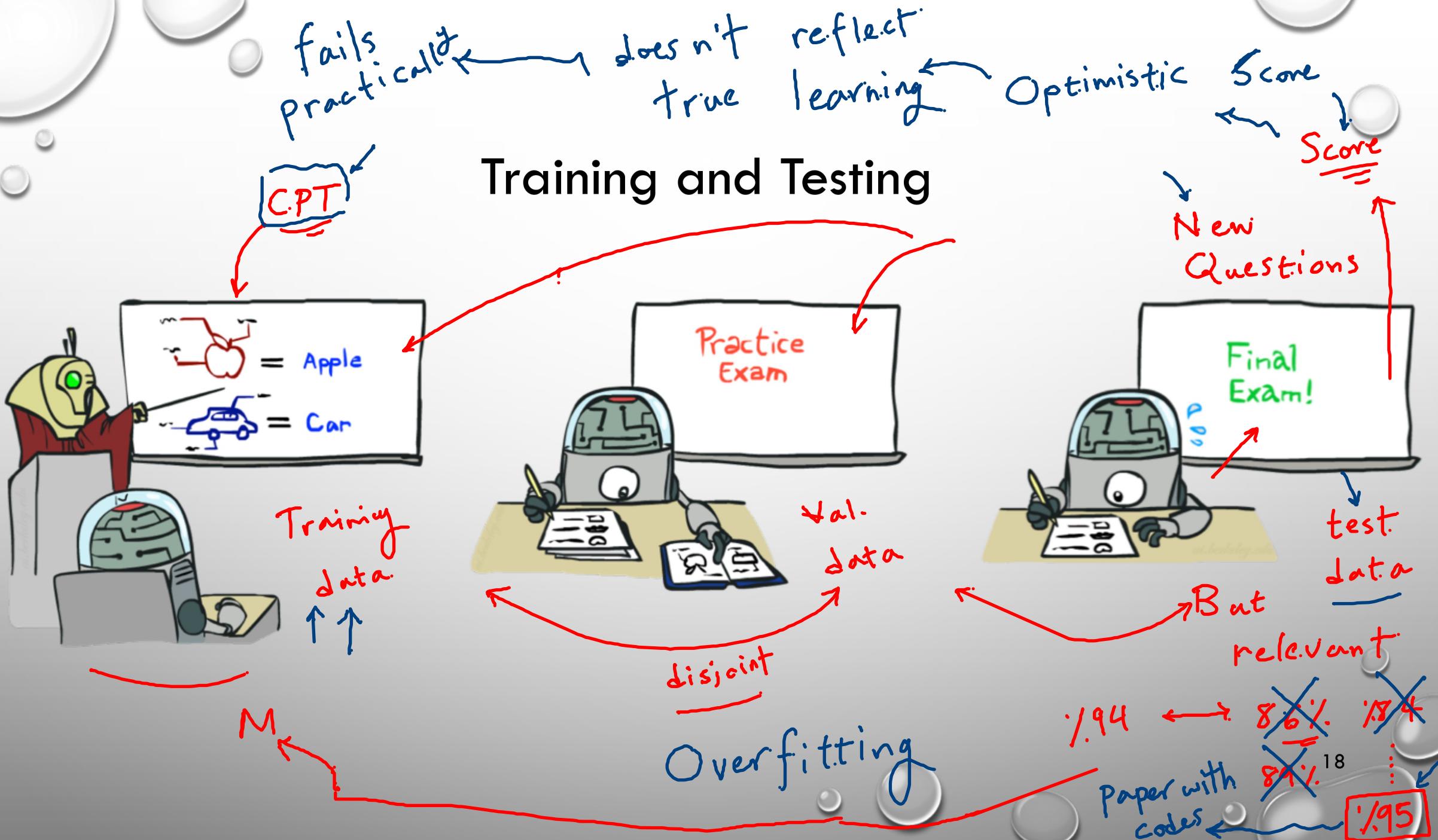
| |
|--------------|
| the : 0.0210 |
| to : 0.0133 |
| of. : 0.0119 |
| 2002: 0.0110 |
| with: 0.0108 |
| from: 0.0107 |
| and : 0.0105 |
| a : 0.0100 |
| ... : |

- Where do these tables come from?
↗

Spam Example

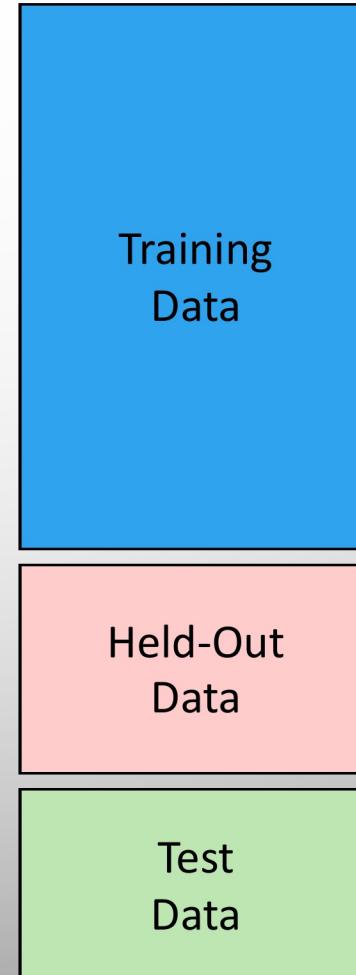
| Word | P(w spam) | P(w ham) | Tot Spam | Tot Ham |
|---------|-----------|----------|----------|---------|
| (prior) | 0.33333 | 0.66666 | -1.1 | -0.4 |

$$P(\text{spam} \mid w) = 98.9$$

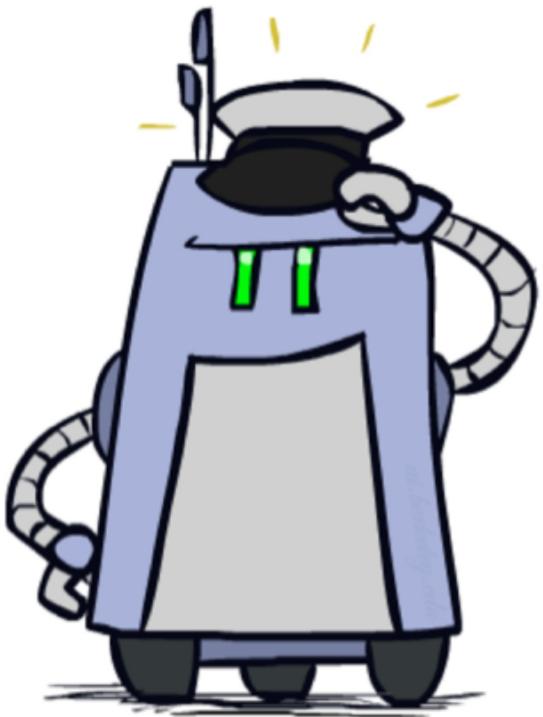
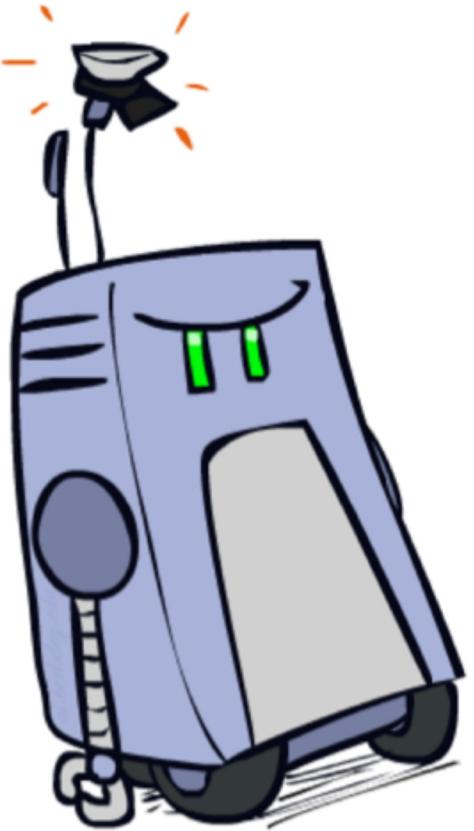


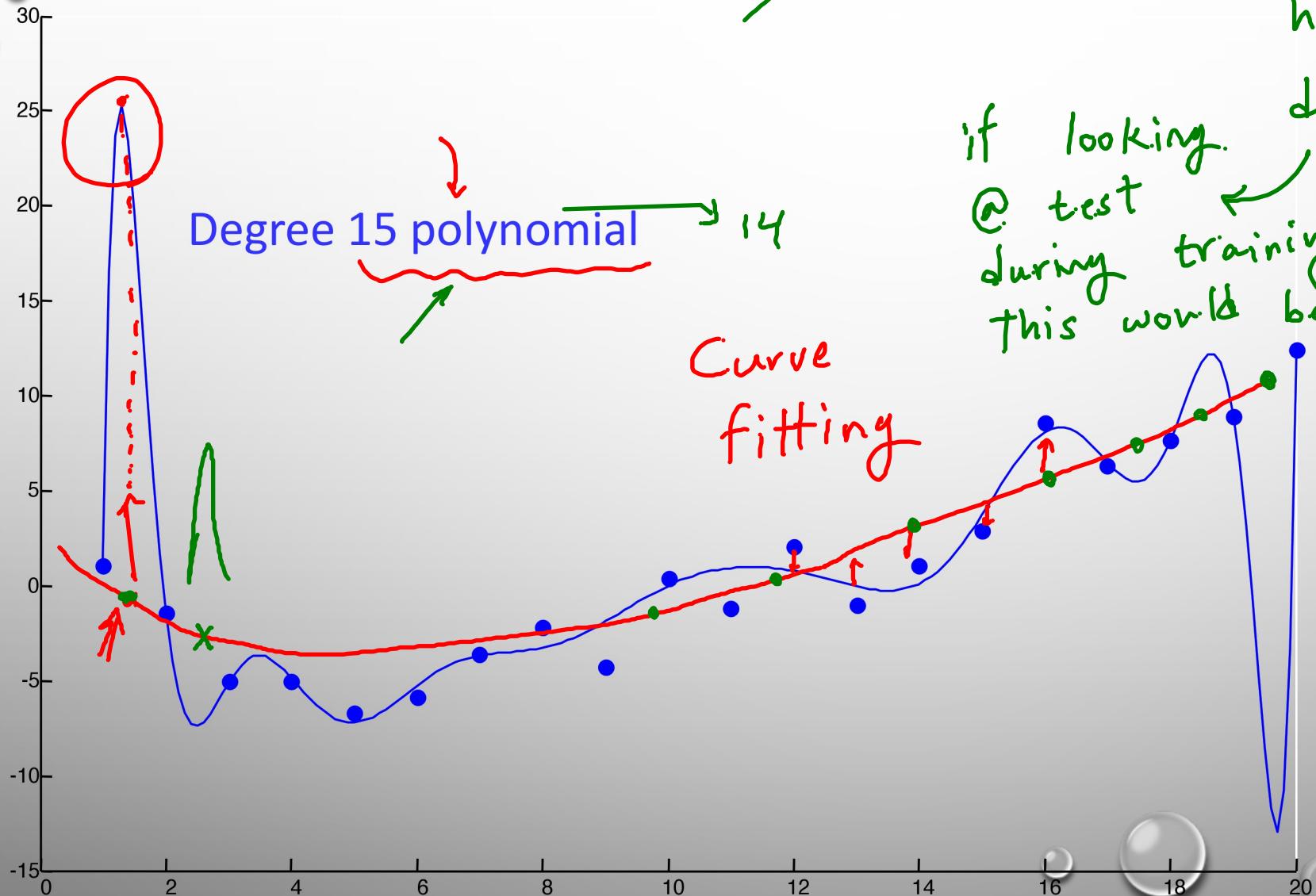
Important Concepts

- Data: labeled instances, e.g. Emails marked spam/ham
 - Training set
 - Held out set
 - Test set
- Features: attribute-value pairs which characterize each x
- Experimentation cycle
 - Learn parameters (e.g. model probabilities) on training set
 - (Tune hyperparameters on held-out set)
 - Compute accuracy of test set
 - Very important: never “peek” at the test set!
- Evaluation
 - Accuracy: fraction of instances predicted correctly
- Overfitting and generalization
 - Want a classifier which does well on *test* data
 - Overfitting: fitting the training data very closely, but not generalizing well
 - We'll investigate overfitting and generalization formally in a few lectures



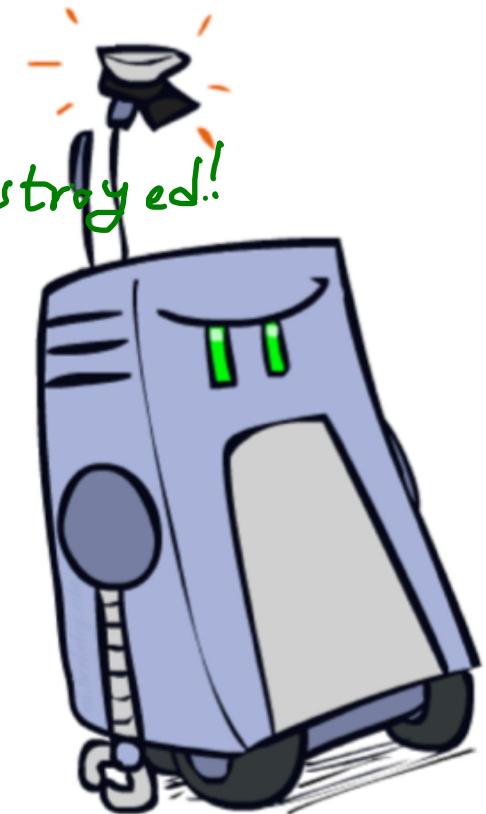
Generalization and Overfitting





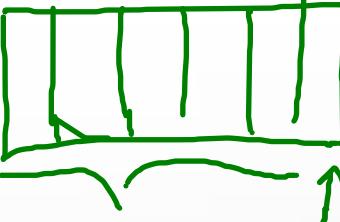
Overfitting → test set
helps us to detect overfitting.

if looking @ test during training this would be destroyed!



K-fold Cross-Validation

Example: Overfitting



$P(\text{features}, C = 2)$

$P(C = 2) = 0.1$

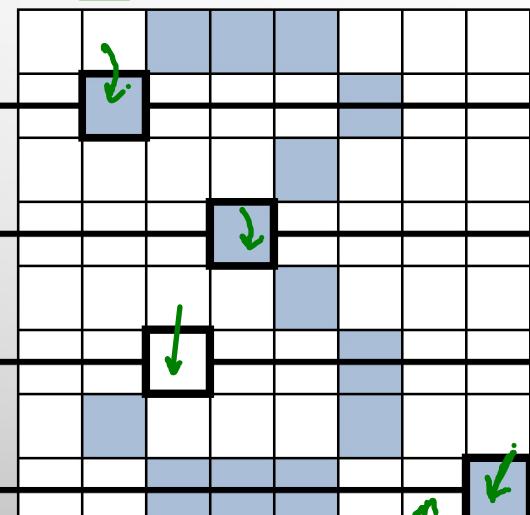
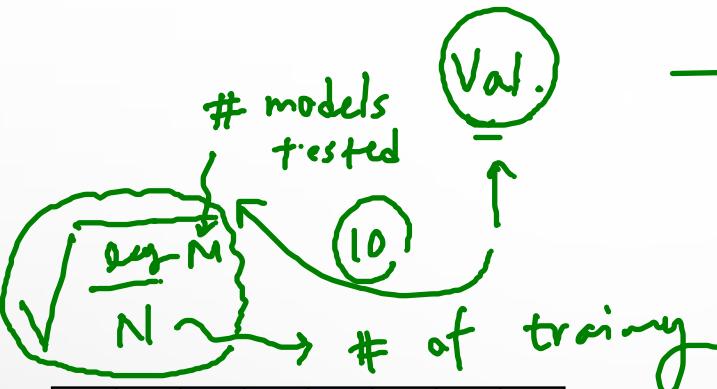
$P(\text{on}|C = 2) = 0.8$

$P(\text{on}|C = 2) = 0.1$

$P(\text{off}|C = 2) = 0.1$

$P(\text{on}|C = 2) = 0.01$

> 0



$P(\text{features}, C = 3)$

$P(C = 3) = 0.1$

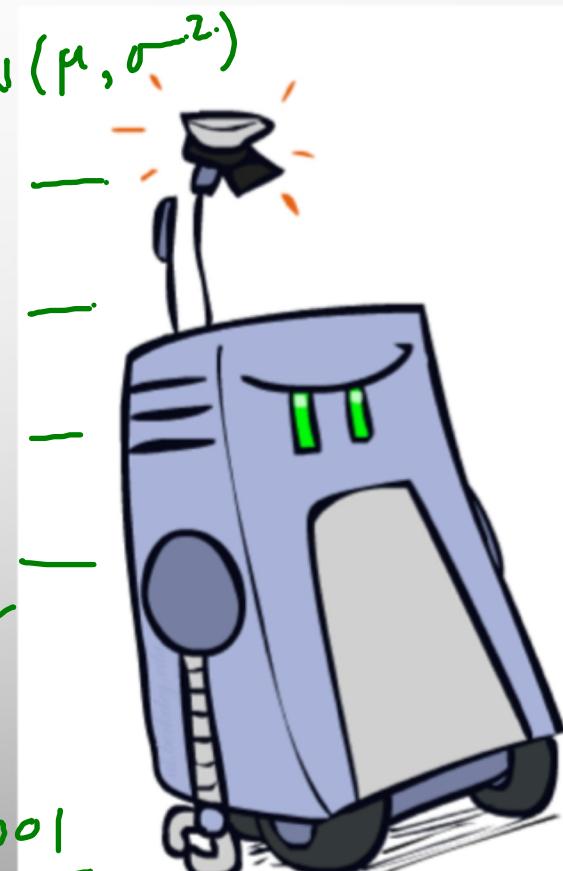
$P(\text{on}|C = 3) = 0.8$

$P(\text{on}|C = 3) = 0.9$

$P(\text{off}|C = 3) = 0.7$

$P(\text{on}|C = 3) = 0.0$

0.0001



Example: Overfitting

- Posteriors determined by *relative* probabilities (odds ratios):

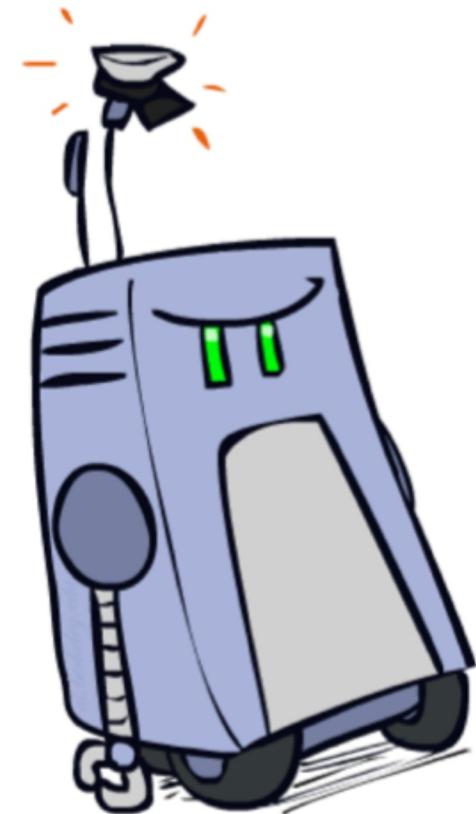
$$\frac{P(W|\text{ham})}{P(W|\text{spam})}$$

| | | |
|------------|---|-----|
| south-west | : | inf |
| nation | : | inf |
| morally | : | inf |
| nicely | : | inf |
| extent | : | inf |
| seriously | : | inf |
| ... | | |

$$\frac{P(W|\text{spam})}{P(W|\text{ham})}$$

| | | |
|------------|---|-----|
| screens | : | inf |
| minute | : | inf |
| guaranteed | : | inf |
| \$205.00 | : | inf |
| delivery | : | inf |
| signature | : | inf |
| ... | | |

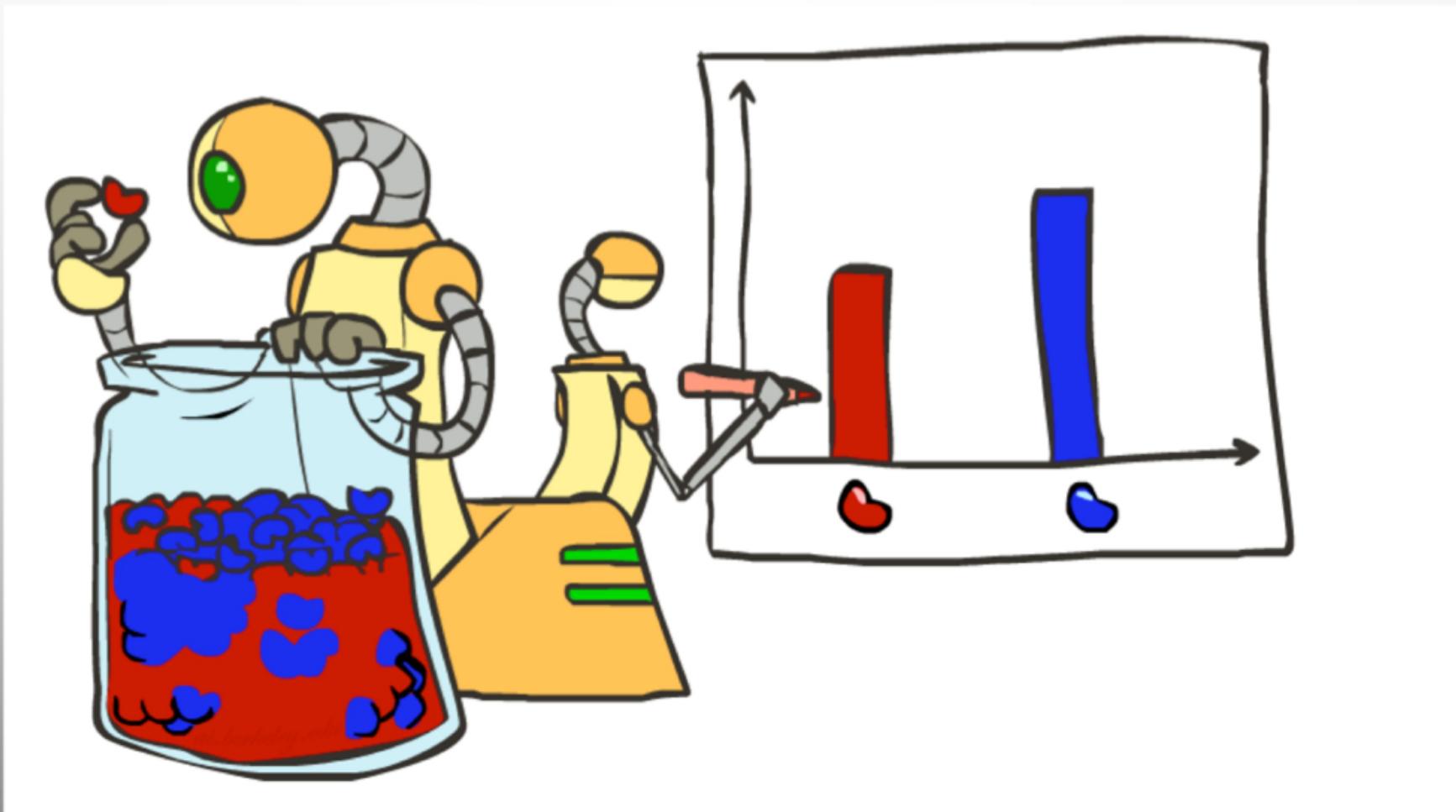
What went wrong here?



Generalization and Overfitting

- Relative frequency parameters will **overfit** the training data!
 - Just because we never saw a 3 with pixel (15,15) on during training doesn't mean we won't see it at test time
 - **Unlikely that every occurrence of "minute" is 100% spam**
 - Unlikely that every occurrence of "seriously" is 100% ham
 - What about all the words that don't occur in the training set at all?
 - **In general, we can't go around giving unseen events zero probability**
- As an extreme case, imagine using the entire email as the only feature
 - Would get the training data perfect (if deterministic labeling)
 - Wouldn't generalize at all
 - Just making the bag-of-words assumption gives us some generalization, but isn't enough
- To generalize better: we need to **smooth** or **regularize** the estimates

Parameter Estimation



Parameter Estimation

$$\max_p \sum x_i \log p + (1-x_i) \log(1-p)$$

$\hat{p} = \frac{\sum x_i}{n}$

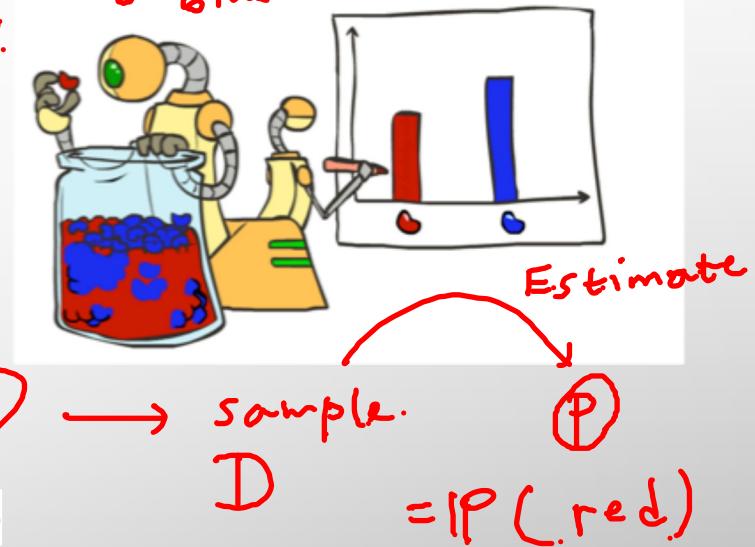
- Estimating the distribution of a random variable
- *Elicitation*: ask a human (why is this hard?)
- *Empirically*: use training data (learning!)
 - e.g.: For each outcome x , look at the *empirical rate* of that value:

$$P_{ML}(x) = \frac{\text{count}(x)}{\text{total samples}}$$

$$\lim_{n \rightarrow \infty} \hat{p} \xrightarrow{P} p$$

$$P(X_i = x) = p^x (1-p)^{1-x}$$

o blue
red



- This is the estimate that maximizes the *likelihood* of the data

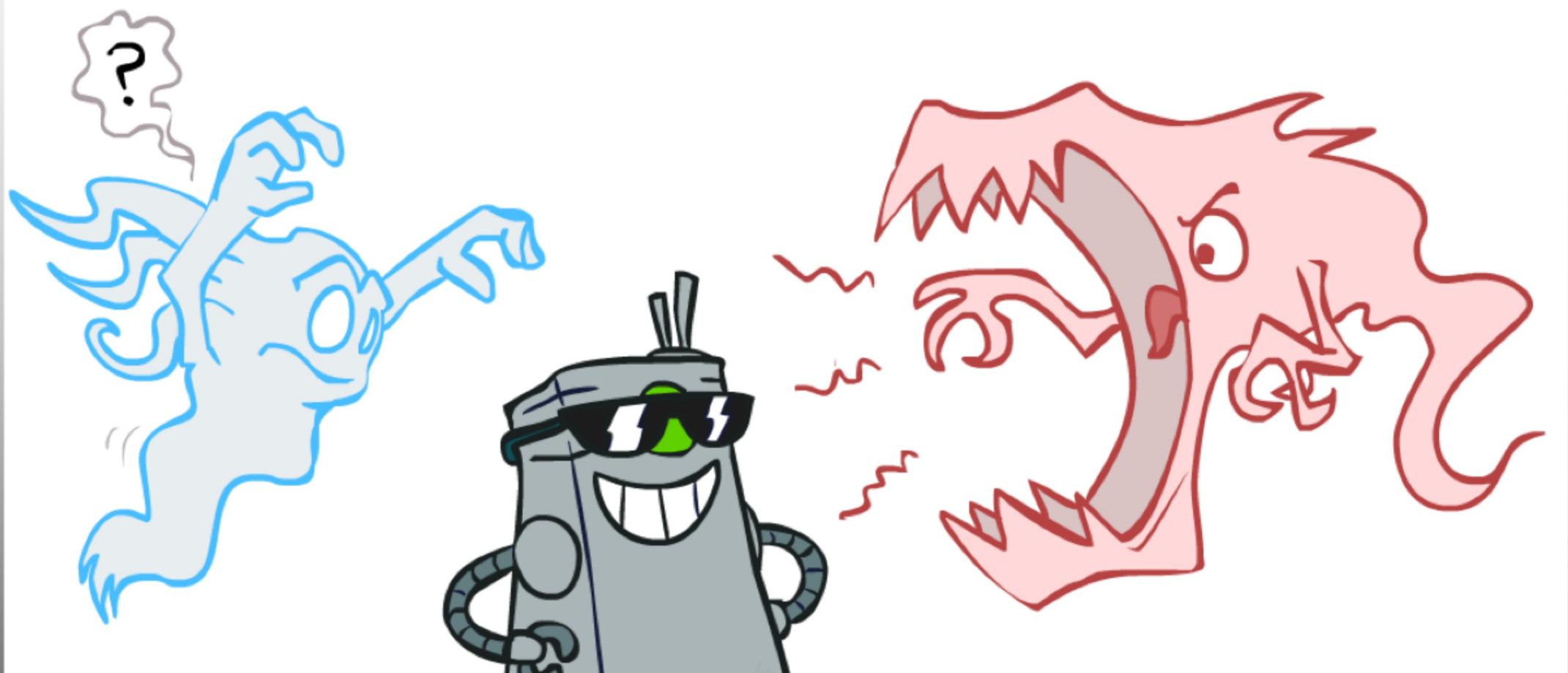
$$L(x, \theta) = \prod_i P_\theta(x_i)$$

$$\max_p \sum_{i=1}^n \log P(x_i = x_i | p)$$

maximum Likelihood

$$\max_p P(D | p) = P(x_1 | p) P(x_2 | p) \cdots P(x_n | p)$$

Smoothing



Maximum Likelihood?

- Relative frequencies are the maximum likelihood estimates

$$\theta_{ML} = \arg \max_{\theta} P(\mathbf{X}|\theta)$$

$$= \arg \max_{\theta} \prod_i P_{\theta}(X_i)$$

→ $P_{ML}(x) = \frac{\text{count}(x)}{\text{total samples}}$

Max.
A posterior est.



- Another option is to consider the most likely parameter value given the data

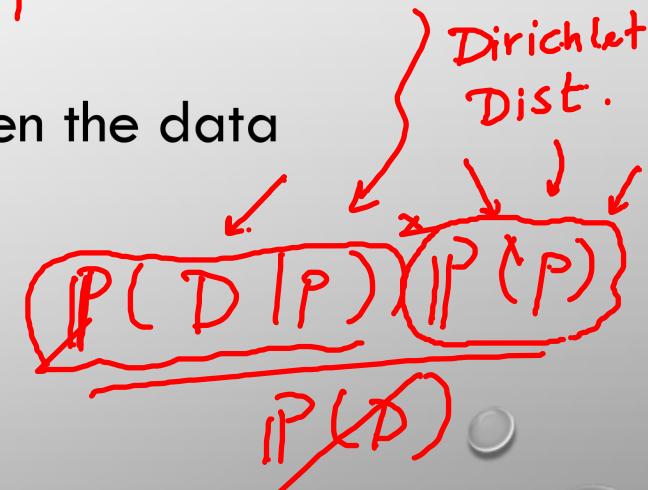
$$\theta_{MAP} = \arg \max_{\theta} P(\theta|\mathbf{X})$$

$$= \arg \max_{\theta} P(\mathbf{X}|\theta)P(\theta)/P(\mathbf{X})$$

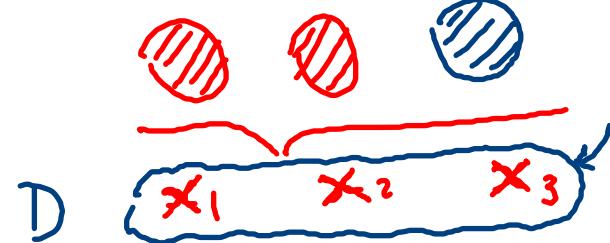
$$= \arg \max_{\theta} P(\mathbf{X}|\theta)P(\theta)$$

Solve
NB
problem?
????

$P(P)$



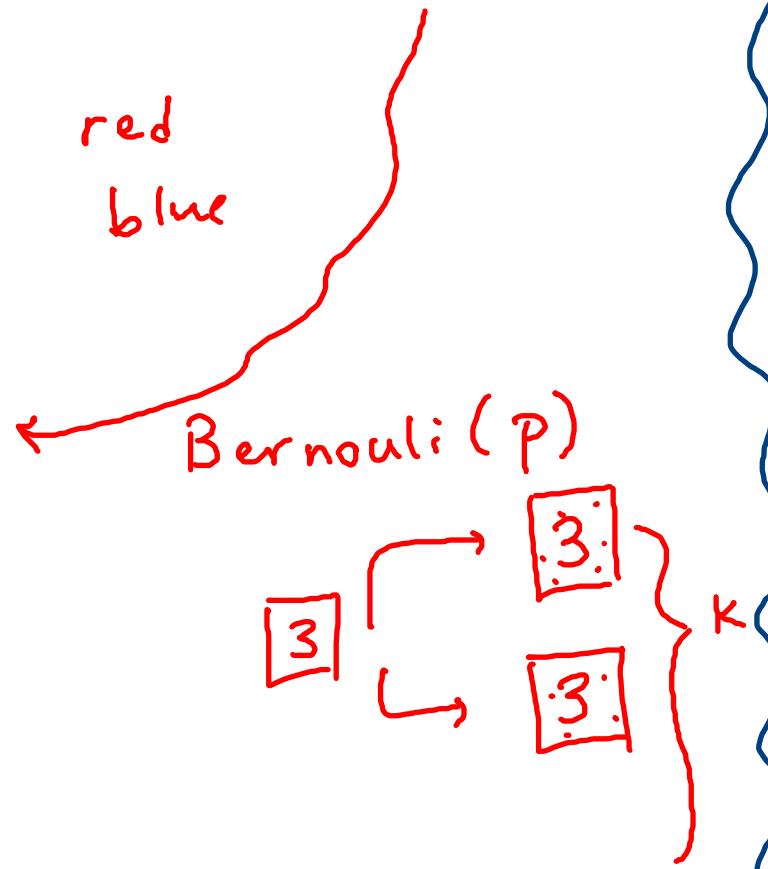
Problem Def.



$$\begin{aligned} p &= P(\text{red}) \\ 1-p &= P(\text{blue}) \end{aligned}$$

$$x_i = \begin{cases} 1 & \text{if red} \\ 0 & \text{if blue} \end{cases}$$

x_1, \dots, x_n \sim iid



Estimation:

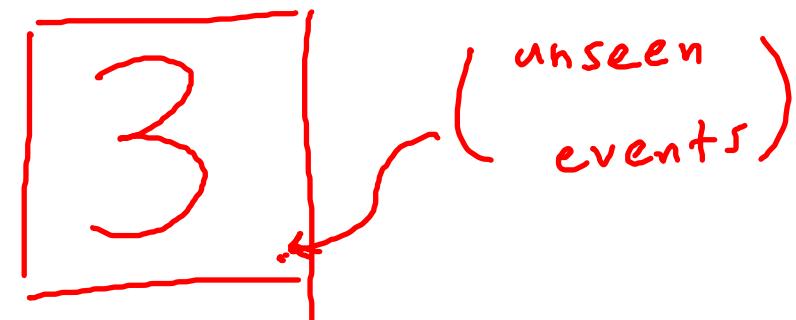
$$g(x_1, \dots, x_n) = \hat{p}$$

↑
unknown param.

$\max_p P(D|p)$
↳ Maximum Likelihood

* In practice and also in theory
if $|D| \rightarrow \infty$, $\hat{p} \xrightarrow{P} p$.

* What if $|D|$ is small.



- Sol.
- (1) increase training set size
 - (2) Augment the data set.
 - (3) Bayesian Est.

$$\text{ML. } P(X_i = x_i | P) = P^{x_i} (1-P)^{1-x_i}$$

$$P(D|P) = \prod_{i=1}^n P(X_i = x_i | P)$$

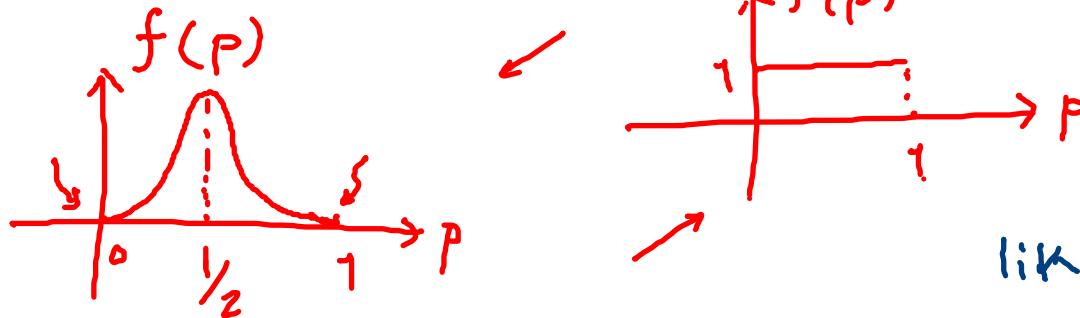
$$\begin{aligned}\log P(D|P) &= \sum_{i=1}^n P(X_i = x_i | P) \\ &= \sum_{i=1}^n x_i \log P + (1-x_i) \log 1-P\end{aligned}$$

$$\frac{\partial}{\partial P} \log P(D|P) = \sum_{i=1}^n \frac{x_i}{P} - \frac{(1-x_i)}{1-P} = 0$$

$$\Rightarrow \hat{P}_{\text{ML}} = \frac{1}{n} \sum_{i=1}^n x_i$$

prior dist.

$$P(p) \sim$$



$$P(D|p)$$

vs.

$$P(p|D)$$

is a better
choice, as it
makes more
sense to pick
the most probable

$$\underbrace{P(p|D)}_{\text{Evidence}} = \frac{\underbrace{P(D|p)}_{\text{Likelihood}} \underbrace{P(p)}_{\text{Prior}}}{P(D)}$$

$$\hat{P}_{MAP} = \operatorname{argmax}_p$$

$$\underbrace{P(D|p)}_{\text{Evidence}} \underbrace{P(p)}_{\text{Prior}} \rightarrow \hat{P}_{MAP} \stackrel{?}{=} \hat{P}_{ML}$$

Dirichlet Dist. ↗

$$\underline{P} = (\underline{P}_1, \underline{P}_2, \dots, \underline{P}_k) \rightarrow \sum_{i=1}^k \underline{P}_i = 1$$

unknown

$$f(\underline{P}) := \frac{1}{Z} \prod_{i=1}^k P_i^{\alpha_i - 1}$$

: params of
4 The Dirichlet
3

$$\underline{P} = (P, 1-P) \quad \alpha_1 = \alpha_2 = 2$$

number of colors

$$\left[\prod_{i=1}^n \left(P^{x_i} (1-P)^{1-x_i} \right) \right] \cdot [P \cdot (1-P)]$$

like likelihood

$P(D|P)$

"Conjugate prior"

$P(P)$

$K \rightarrow$ hyperparam → tune
val. data ↗

$$K \in \{0, 1, 2, \dots, 10\}$$

training data

$\hookrightarrow K=0 \rightarrow$ has the most consistency with the data → overfit

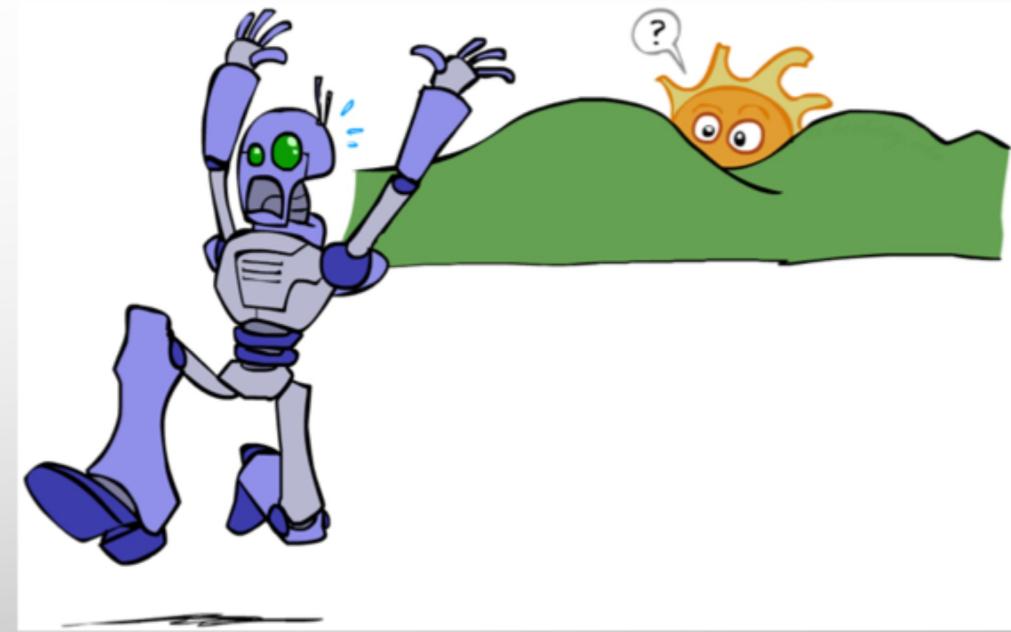
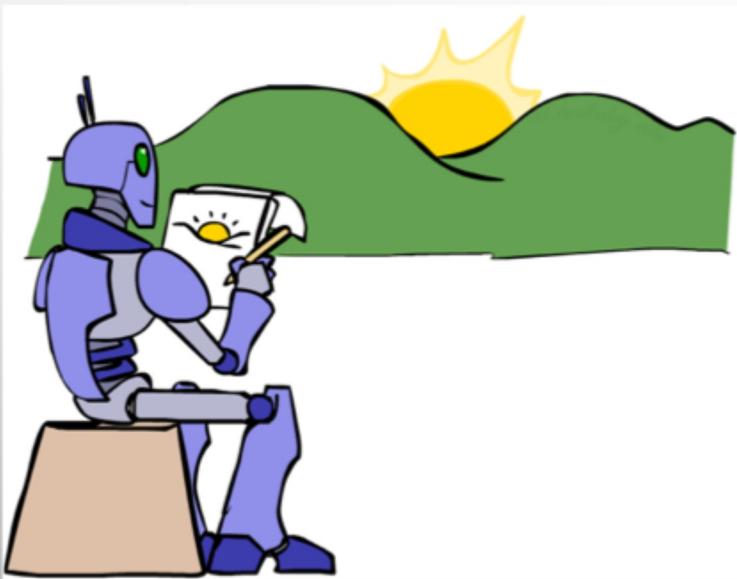


$$\hat{P}_{MAP} = \frac{\sum x_i + 1}{n + 2}$$

$K \rightarrow \infty \quad \hat{P}_{MAP} = 1/2$

Laplace Smoothing

Unseen Events

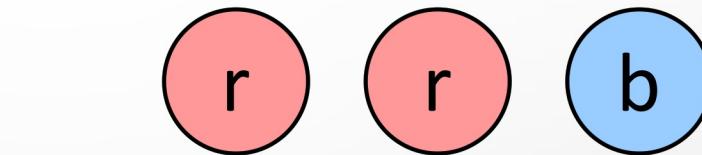


Laplace Smoothing

- Laplace's estimate:
 - Pretend you saw every outcome once more than you actually did

$$P_{LAP}(x) = \frac{c(x) + 1}{\sum_x [c(x) + 1]}$$

$$= \frac{c(x) + 1}{N + |X|}$$



$$P_{ML}(X) =$$

$$P_{LAP}(X) =$$

- Can derive this estimate with *Dirichlet priors* (See *Probabilistic Graphical Models course*)

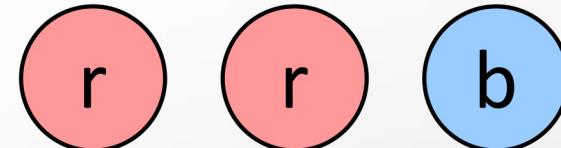
Laplace Smoothing

- Laplace's estimate (extended):
 - Pretend you saw every outcome k extra times

$$P_{LAP,k}(x) = \frac{c(x) + k}{N + k|X|}$$

- What's Laplace with k = 0?
- k is the **strength** of the prior
- Laplace for conditionals:
 - Smooth each condition independently:

$$P_{LAP,k}(x|y) = \frac{c(x, y) + k}{c(y) + k|X|}$$



$$P_{LAP,0}(X) =$$

$$P_{LAP,1}(X) =$$

$$P_{LAP,100}(X) =$$

Estimation: Linear Interpolation

$$P(x|Y)$$

- In practice, Laplace often performs poorly for $P(X|Y)$:

{
• When $|X|$ is very large → Text classification $|X| = \text{size Dic.}$
• When $|Y|$ is very large → 1000 topic

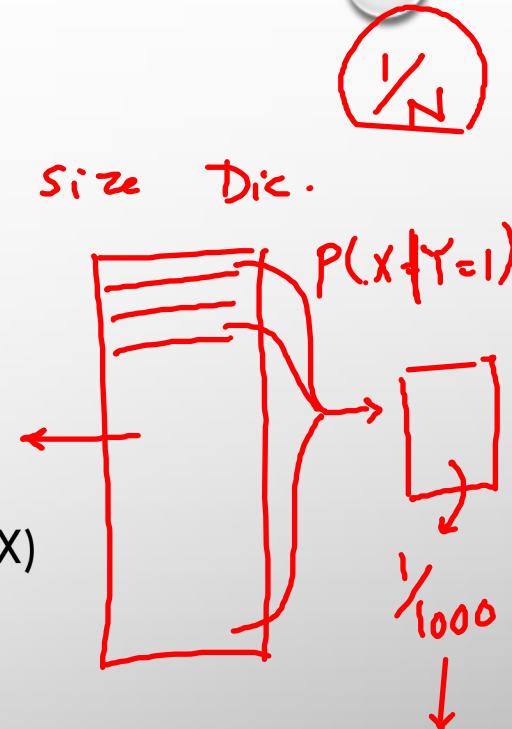
- Another option: linear interpolation

- Also get the empirical $P(X)$ from the data
- Make sure the estimate of $P(X|Y)$ isn't too different from the empirical $P(X)$

$$P_{LIN}(x|y) = \alpha \hat{P}(x|y) + (1.0 - \alpha) \hat{P}(x)$$

- What if α is 0? 1?
• See Stochastic Processes course for more interesting options of making the estimation.

noisy less noisy
 $\hat{P}(x|y)$ hyper param



Real NB: Smoothing

- For real classification problems, smoothing is critical
- New odds ratios:

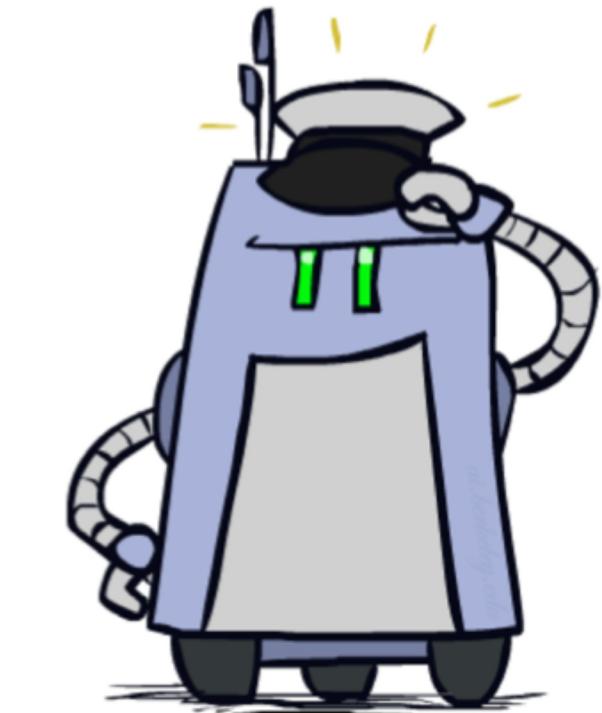
$$\frac{P(W|\text{ham})}{P(W|\text{spam})}$$

| | | |
|-----------|---|------|
| helvetica | : | 11.4 |
| seems | : | 10.8 |
| group | : | 10.2 |
| ago | : | 8.4 |
| areas | : | 8.3 |
| ... | | |

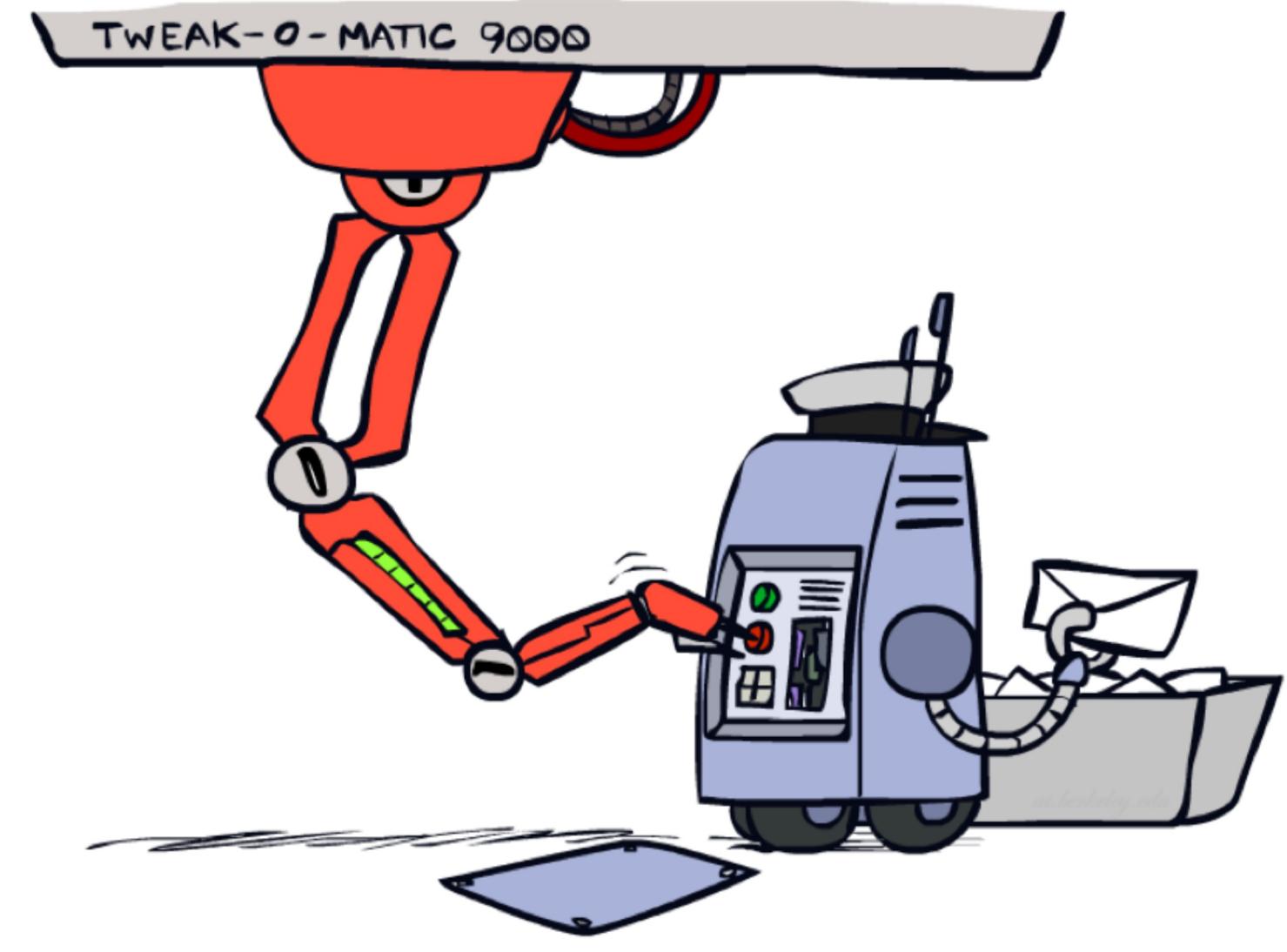
$$\frac{P(W|\text{spam})}{P(W|\text{ham})}$$

| | | |
|---------|---|------|
| verdana | : | 28.8 |
| Credit | : | 28.4 |
| ORDER | : | 27.2 |
| | : | 26.9 |
| money | : | 26.5 |
| ... | | |

Do these make more sense?



Tuning



$$P(X_i|Y)$$

| | X_i | Y | P |
|-----|-------|-----|-----|
| T | T | T | O |
| F | F | T | O |
| T | F | F | O |
| F | F | F | O |

Tuning on Held-Out Data

- Now we've got two kinds of unknowns
 - Parameters: the probabilities $P(X|Y)$, $P(Y)$
 - Hyperparameters: e.g. The amount / type of smoothing to do, k , α
- What should we learn where?
 - Learn parameters from training data
 - Tune hyperparameters on different data
 - Why?
 - For each value of the hyperparameters, train and test on the held-out data
 - Choose the best value and do a final test on the test data

