

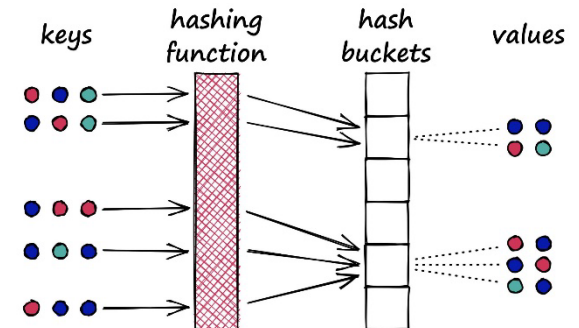
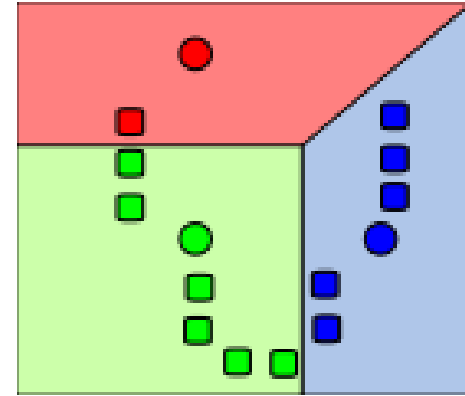


EM Algorithm

Applied Machine Learning
Derek Hoiem

Things to remember from last class

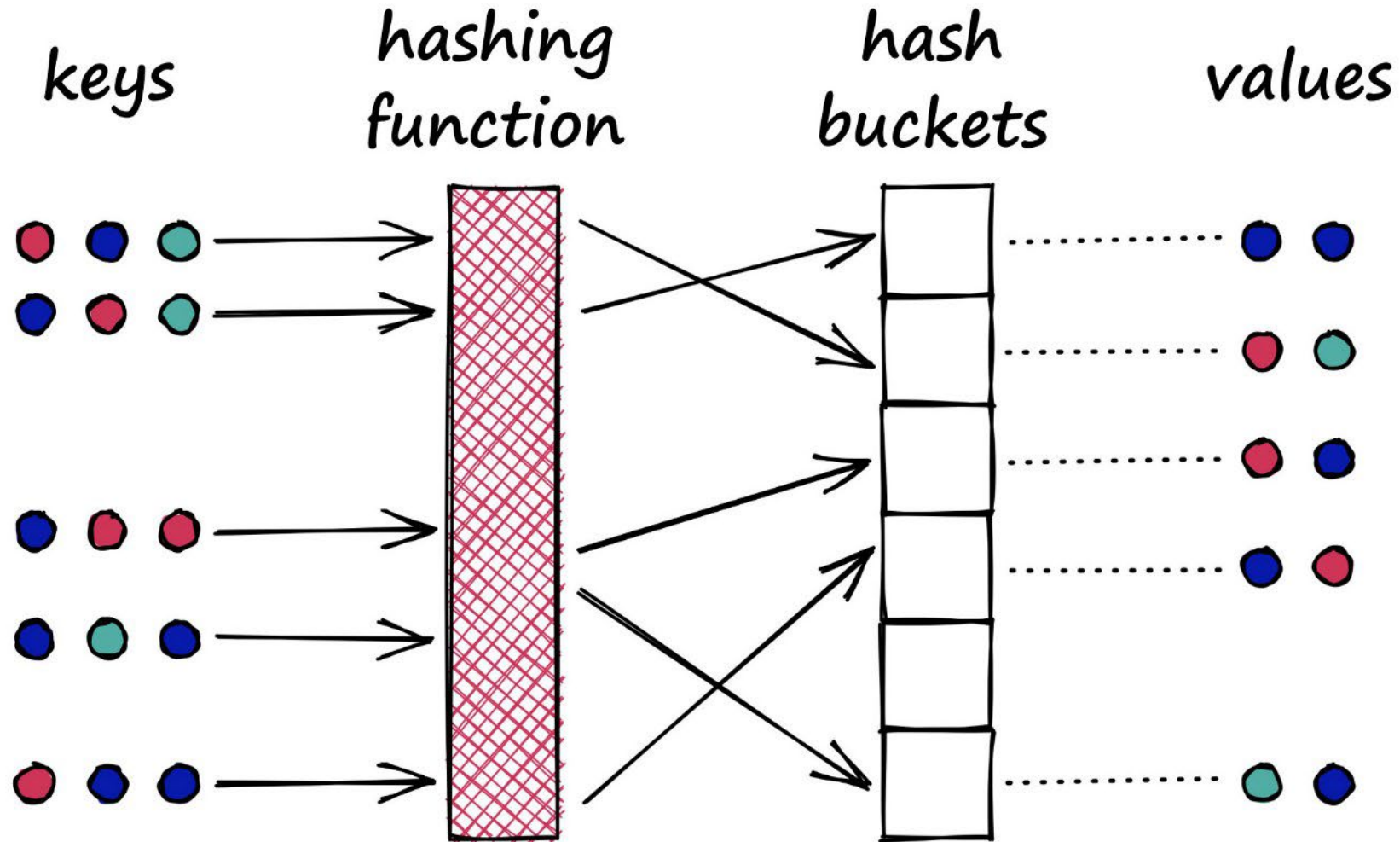
- Clustering groups similar data points
- K-means is the must-know method, but there are many others
- TF-IDF is used for similarity of tokenized documents and used with index for fast search
- *Approximate search methods like LSH can be used to find similar points quickly*
- Use highly optimized libraries like FAISS



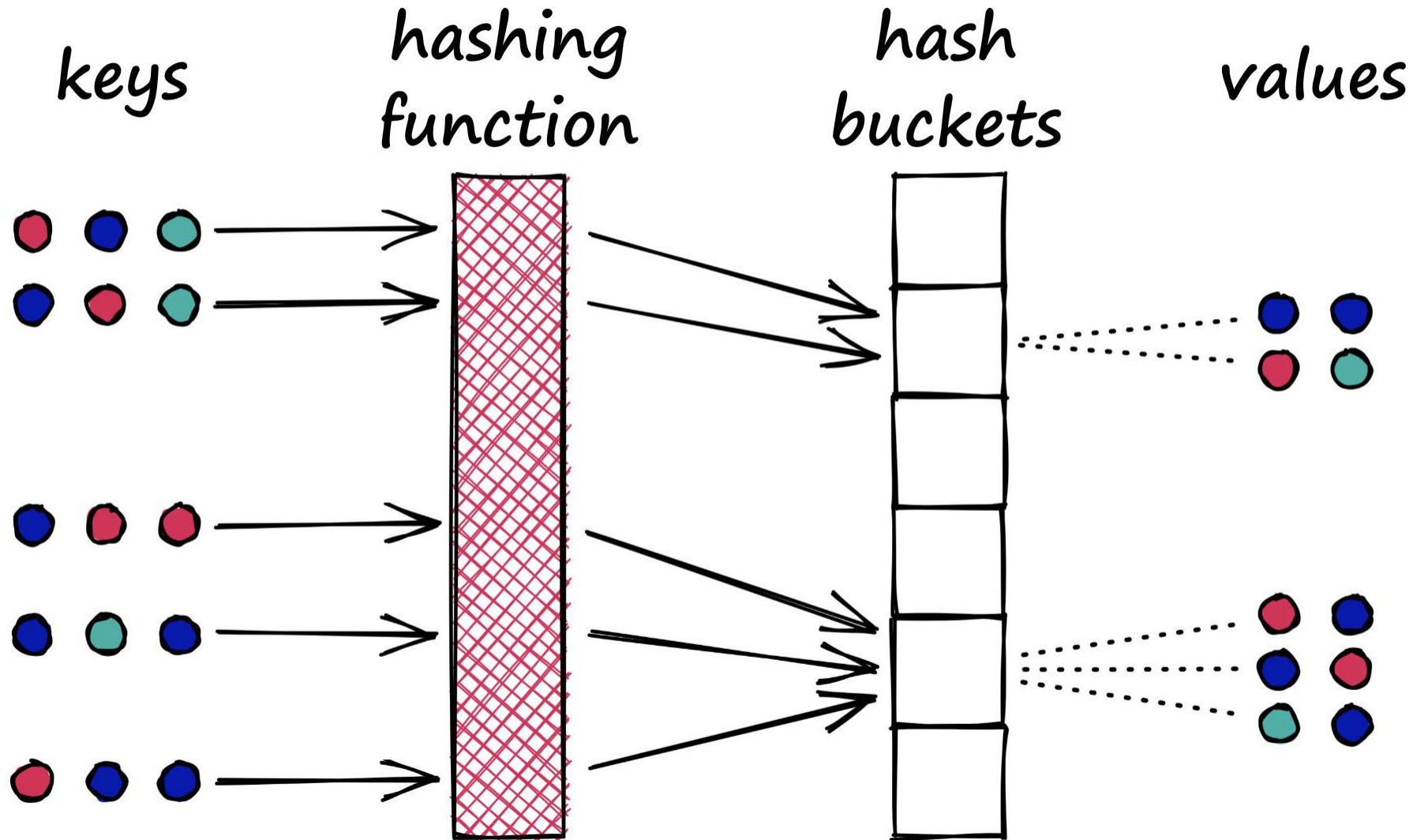
Locality Sensitive Hashing (LSH)

A fast approximate search method to return similar data points to query

A typical hash function aims to place different values in separate buckets



LSH aims to put similar keys in the *same* bucket



Basic LSH process

1. Convert each data point into an array of bits or integers, using the same conversion process/parameters for each
2. Map the arrays into buckets (e.g. with 10 bits, you have 2^{10} buckets)
 - Can use subsets of arrays to create multiple sets of buckets
3. On query, return points in the same bucket(s)
 - Can check additional buckets by flipping bits to find points within hash distances greater than 0

Random Projection LSH

Data Preparation

Given data $\{X\}$ with dimension d :

1. Create b random vectors h_b of length d `h = np.random.rand(nbits, d) - .5`
2. Convert each X_n to b bits: $X_n h^T > 0$

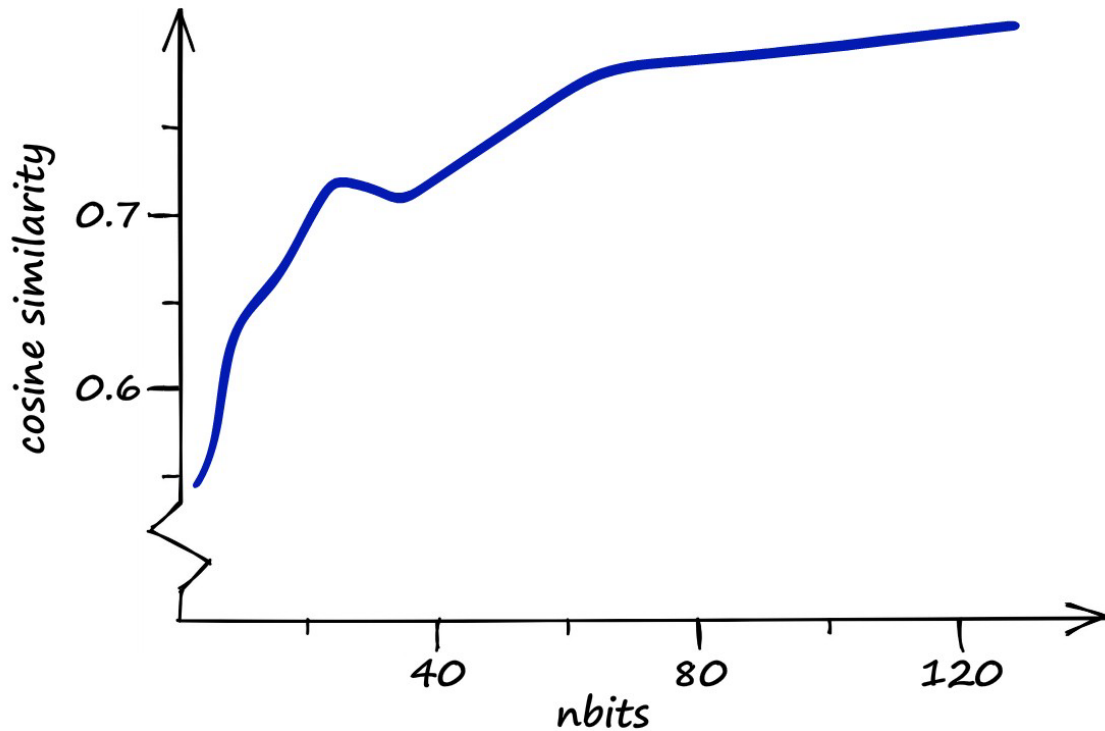
Query

1. Convert X_q to bits using h
2. Check buckets based on bit vector and similar bit vectors to return most similar data points

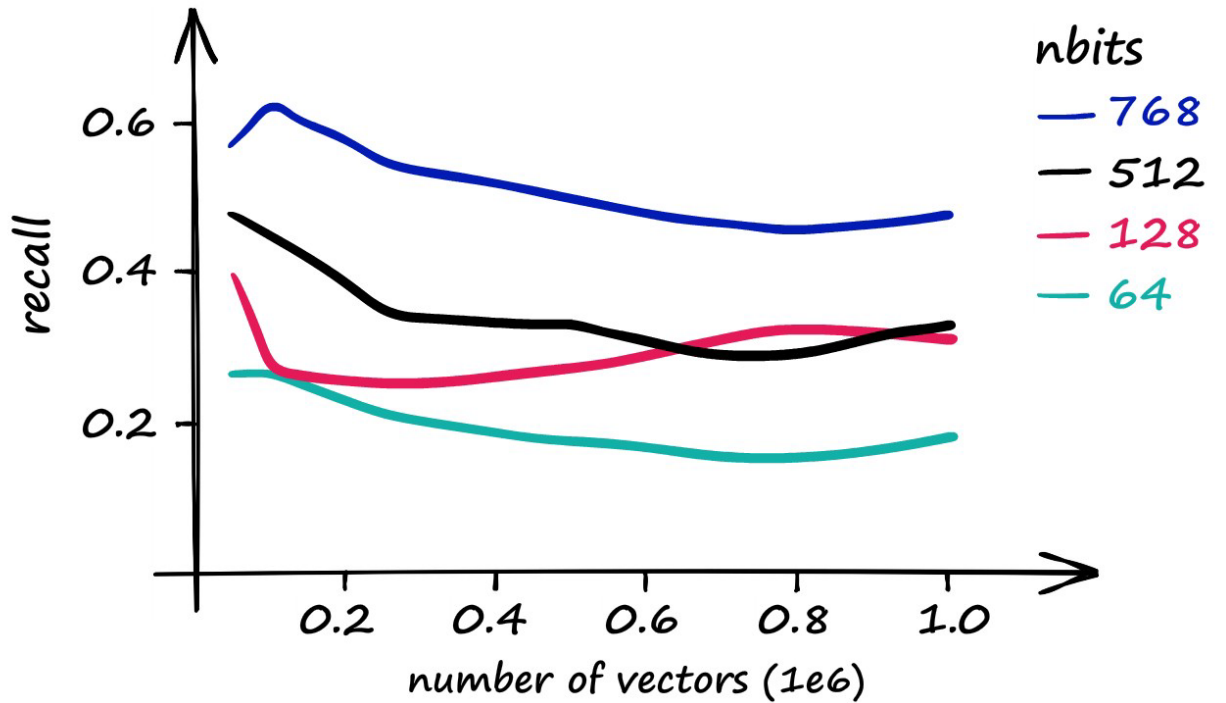
Key parameter: nbits

- Example with 1M 128-bit SIFT vectors

More bits returns more similar vectors



Similarity of LSH-returned vector

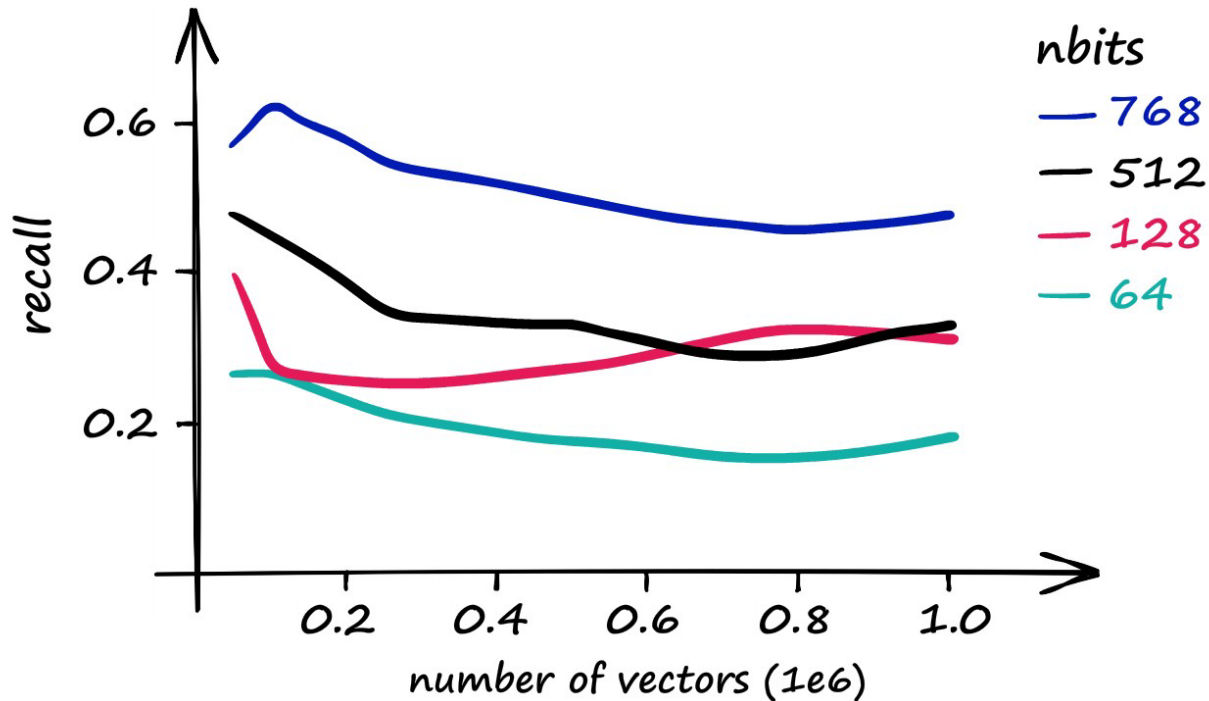


Recall vs. exact nearest neighbor

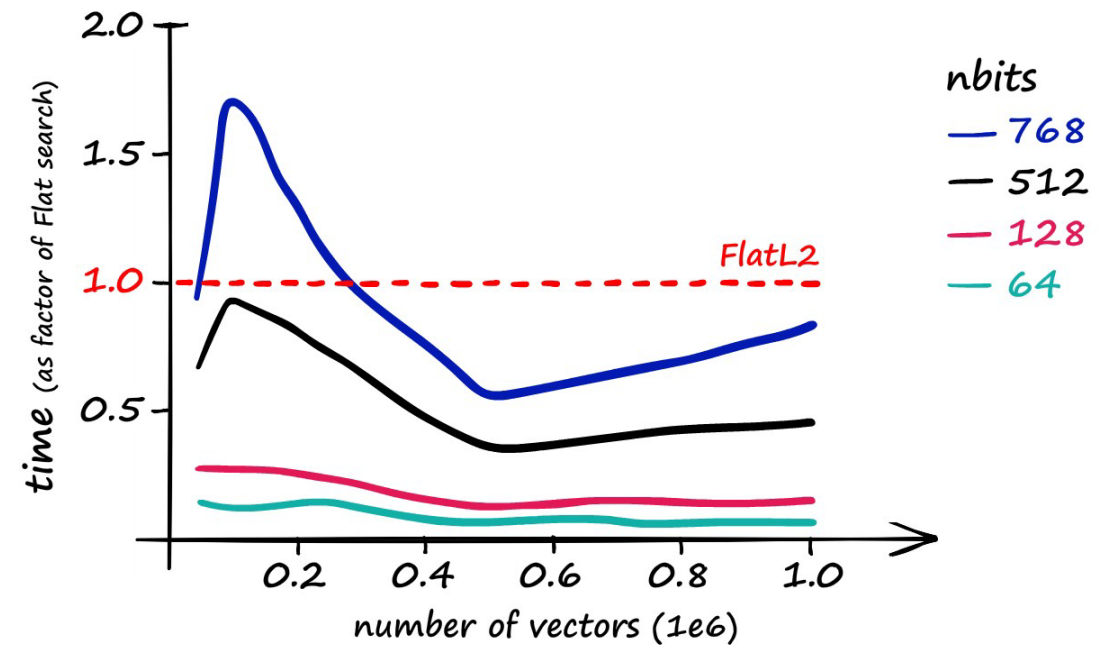
Key parameter: nbits

- Example with 1M 128-bit SIFT vectors

But more bits takes more time to query because it needs to search more buckets to find a collision



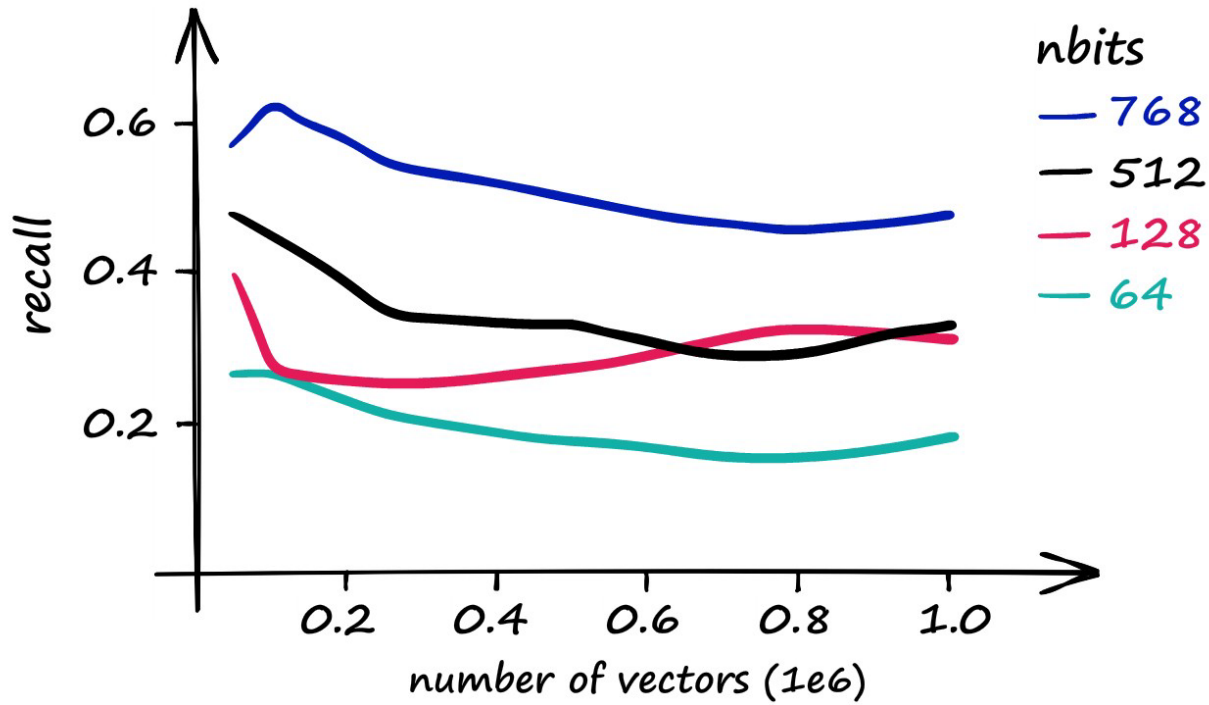
Recall vs. exact nearest neighbor



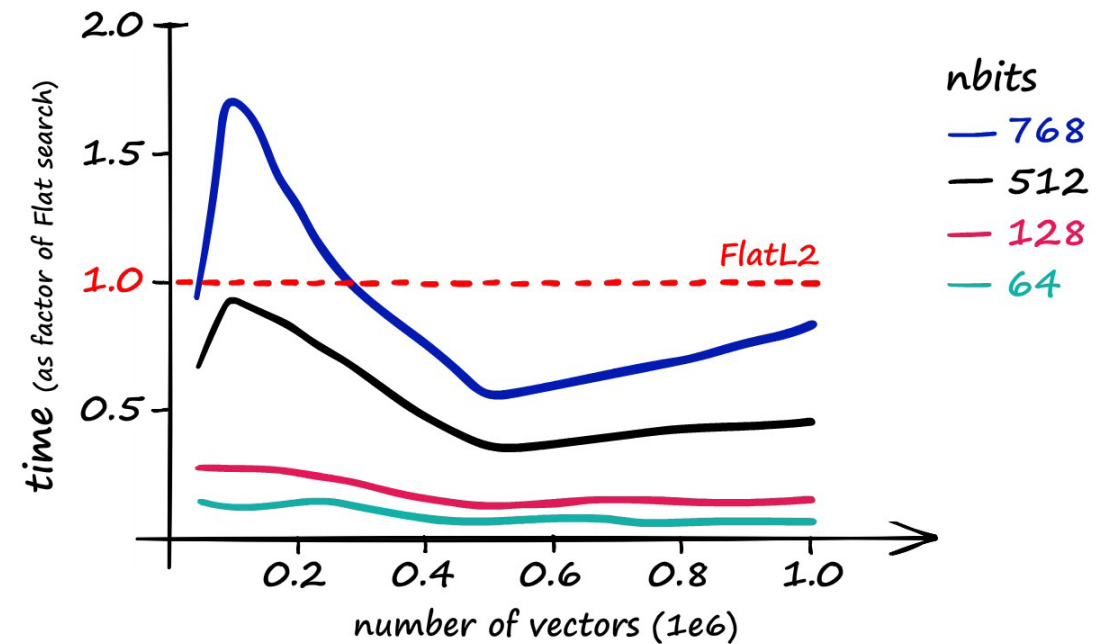
Time compared to brute force search

Key parameter: nbits

- Rule of thumb: nbits = dim is a decent choice (1 bit per feature dimension)
- Optionally, can retrieve K closest data points and then use brute force search on those



Recall vs. exact nearest neighbor



Time compared to brute force search

Nice video about LSH in faiss:

https://youtu.be/ZLfdQq_u7Eo

which is part of this very detailed and helpful post:

<https://www.pinecone.io/learn/locality-sensitive-hashing-random-projection/>

Today's Class: EM Algorithm

I will describe three problems. Think about what they have in common.

“Bad Annotators” Problem

You want to train an algorithm to predict whether a photograph is attractive. You collect annotations from Mechanical Turk. Some annotators try to give accurate ratings, but others answer randomly.

Challenge: Determine which people to trust and the average rating by accurate annotators.



Annotator
Ratings

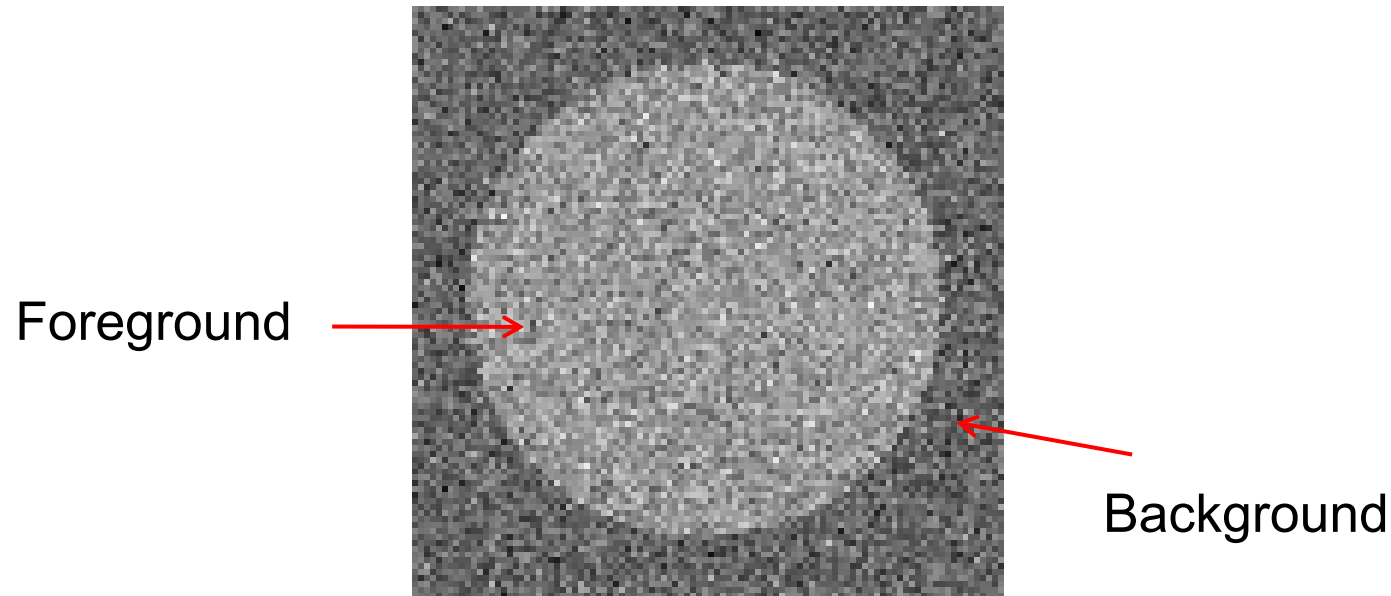
10
8
9
2
8

Photo: Jam343 (Flickr)

Foreground/Background Segmentation

You are given an image and want to assign foreground/background pixels.

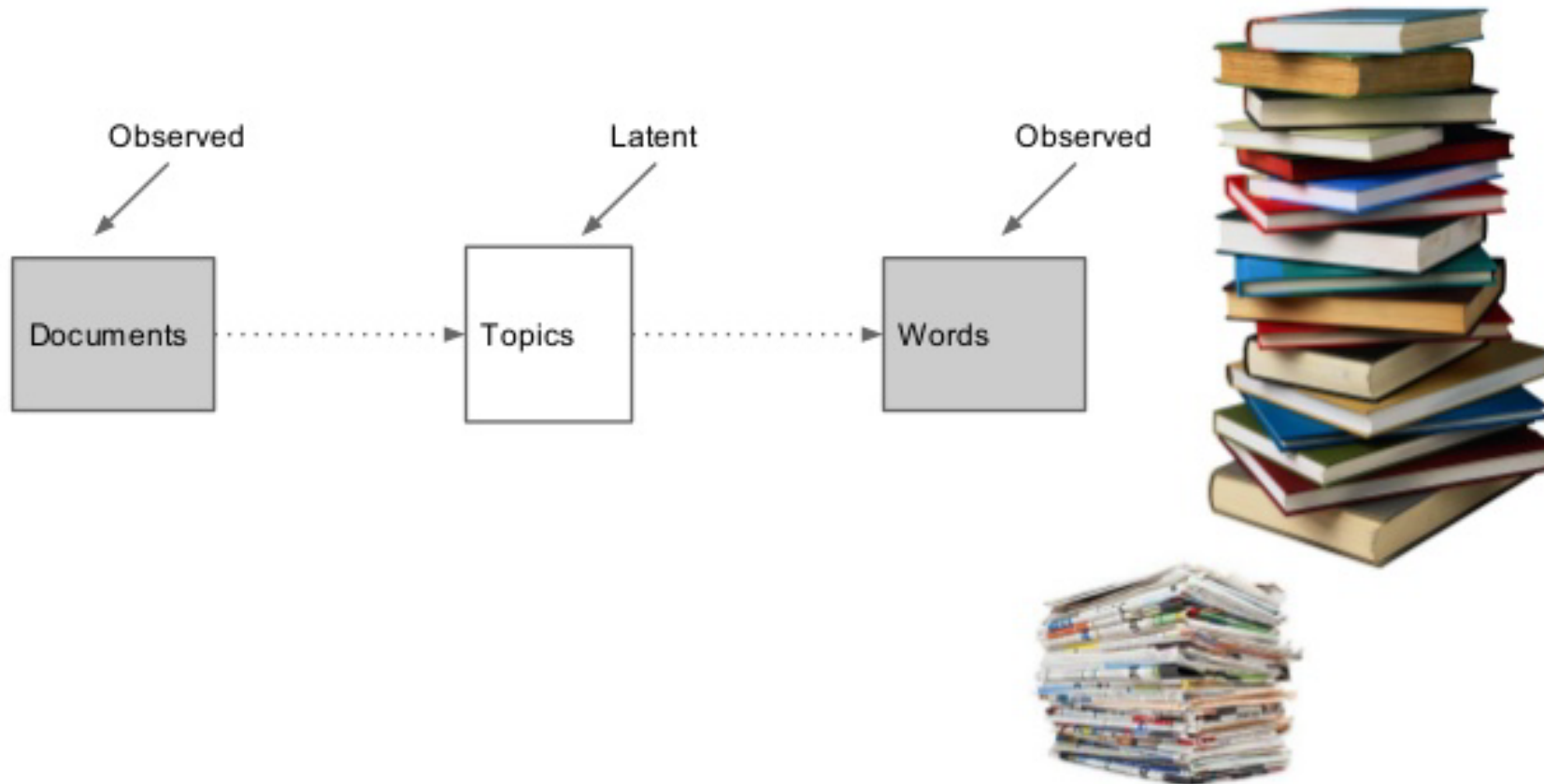
Challenge: Segment the image into figure and ground without knowing what the foreground looks like in advance.



Topic Models

Documents have a “topic” that is predictive of the words

Challenge: We don’t know what the topics are, or what distribution of words each has



What do these problems have in common?

1. We think there is some underlying factor that is not known
 - Whether annotator is good or bad
 - Whether pixel is in foreground or background
 - Topic of a document
2. We have some model for the probability of data given that underlying factor
3. But we don't know the parameters of the model

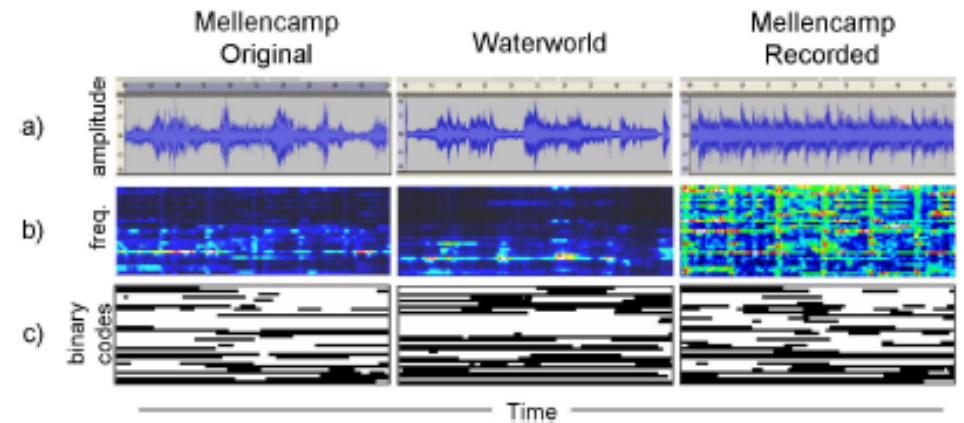
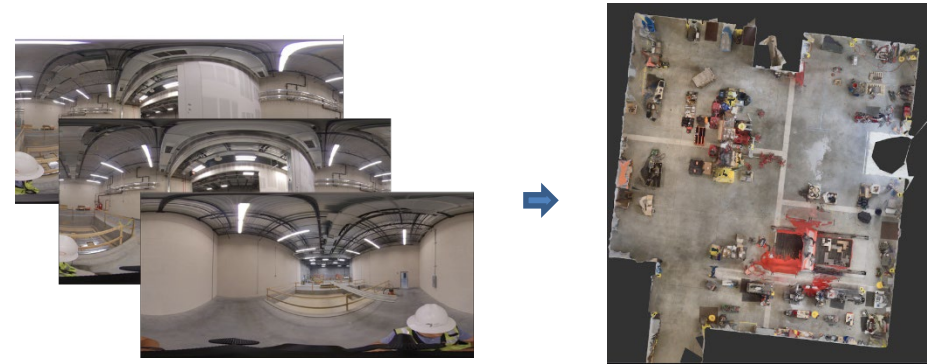
These are “missing data” or “latent variable” problems – a critical piece of information is not observed

Today's Class

- Examples of problems with hidden or latent variables
 - Untrustworthy annotators
 - Pixel segmentation
 - Topic models
- Background
 - Maximum Likelihood Estimation
 - Probabilistic Inference
- Dealing with Latent Variables (latent = hidden, not observed)
 - EM algorithm, Bad annotator problem
 - Hard EM

I have used EM in research and practice many times, e.g.

- Given multiple images and scene geometry, estimate true color of each floor map pixel
 - Latent variables: which pixels are occluded
- For an audio clip of music with background noise, which extracted sound signatures are due to music vs background noise?
 - Latent variables: whether each extracted sound signature at a given time is due to background or music
- Mixture of Gaussian probability model (next class)



Bad Annotator Problem

You hire annotators to label attractiveness of images. Some annotators do their best. Some are “bad” and assign random scores.

Goal: We want to estimate the average score of good annotators for each image

Assumptions

1. Bad annotators are always bad. Good annotators are always good.
2. For each image i , the scores from good annotators follow a Gaussian distribution $s_{ia} \sim N(\mu_i, \sigma^2)$

$$P(s_{ia} | z_a = 1) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2} \frac{(s_{ia} - \mu_i)^2}{\sigma^2}\right)$$

3. The scores from bad annotators always follow a uniform distribution

$$P(s_{ia} | z_a = 0) = 1 \quad (\text{scores range from 0 to 1})$$

Notation

$s_{ia} \in [0,1]$: score for image i by annotator a

$z_a \in \{0,1\}$: whether annotator a is good

μ_i : true mean score for image i

σ : standard deviation of true scores, same for each image

π_z : $P(z_a = 1)$, prior probability that annotator is good

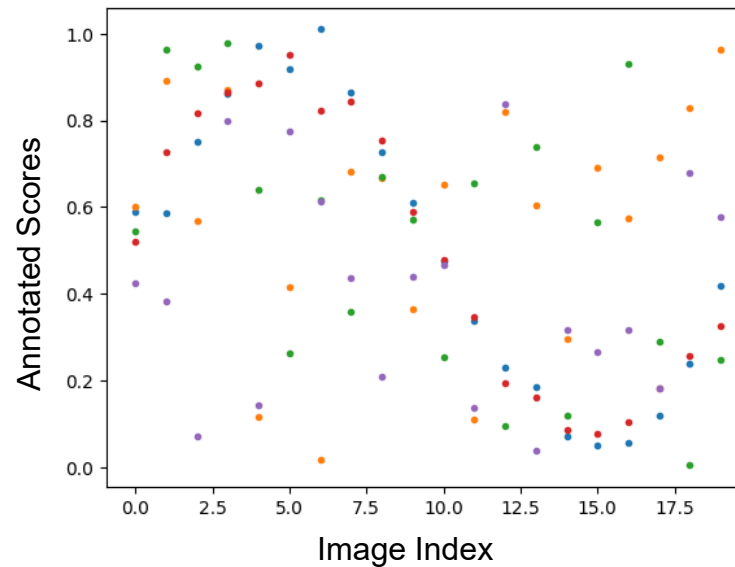
Latent Variable Problems: Bad Annotator

Challenge: Figure out which annotators are good and estimate the true mean score for each image

Three steps:

1. If we knew which annotators were good, how would we estimate the score distribution for each image?
2. Given the distribution parameters, how do we compute the likelihood that an annotator is good?
3. How can we get annotator labels and score models at once?

$i=0$	$i=1$	$i=2$...
0.59	0.88	0.7	
0.52	0.96	0.71	
0.50	0.88	0.52	
0.39	0.94	0.65	
0.65	0.69	0.46	



Maximum Likelihood Estimation

1. If we knew which annotators were good, how would we estimate the score distribution for each image?

Solve for parameters that maximize the data likelihood

$$\begin{array}{l} \text{MLE:} \\ \text{data} \swarrow \\ \underline{s}_i = \{s_{i0} \dots s_{im}\} \\ \hat{\theta} = \underset{\theta}{\operatorname{argmax}} p(\underline{s}_i | \theta) \quad \swarrow \text{parameters} \end{array}$$

Scores are independent of each other, given the model

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \prod_a p(s_{ia} | \theta)$$

Model is Gaussian

$$p(s_{ia} | \mu_i, \sigma) = \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{1}{2} \frac{(s_{ia} - \mu_i)^2}{\sigma^2}\right)$$

Solving for mean

Easier to take derivative of sum of logs than a product, and $\log(x)$ has the same max as x

To find max wrt a variable, set the partial derivative wrt that variable to 0

Do the math

Solve. M is number of "a"s

$$p(s_{ia} | \mu_i, \sigma) = \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{1}{2} \frac{(s_{ia} - \mu_i)^2}{\sigma^2}\right)$$

$$\operatorname{argmax}_{\theta} \prod_a p(s_{ia} | \theta)$$

$$\frac{\partial}{\partial \mu} \sum_a \left[\log \frac{1}{\sqrt{2\pi}} + \log \frac{1}{\sigma} - \frac{1}{2} \frac{(s_{ia} - \mu_i)^2}{\sigma^2} \right] = 0$$

$$\sum_a (s_{ia} - \mu_i) / \sigma^2 = 0 \Rightarrow \sum_a s_{ia} - \sum_a \mu_i = 0$$

$$\mu_i = \sum_a s_{ia} / M$$

Solving for standard deviation

Now take partial derivative wrt sigma. Since sigma is the same for all images, we are now summing over all images and annotations

$$\frac{\partial}{\partial \sigma} \sum_i \sum_a \left[\log \frac{1}{\sqrt{2\pi}} + \log \frac{1}{\sigma} - \frac{1}{2} \frac{(s_{ia} - \mu_i)^2}{\sigma^2} \right] = 0$$

$$\sum_i \sum_a \left[-\frac{1}{\sigma} + \frac{(s_{ia} - \mu_i)^2}{\sigma^3} \right] = 0 \Rightarrow \sum_{i,a} -\sigma^{-2} + \sum_{i,a} (s_{ia} - \mu_i)^2 = 0$$

Solution is average squared difference from mean

$$\sigma^2 = \sum_{i,a} (s_{ia} - \mu_i)^2 / (M \cdot N)$$

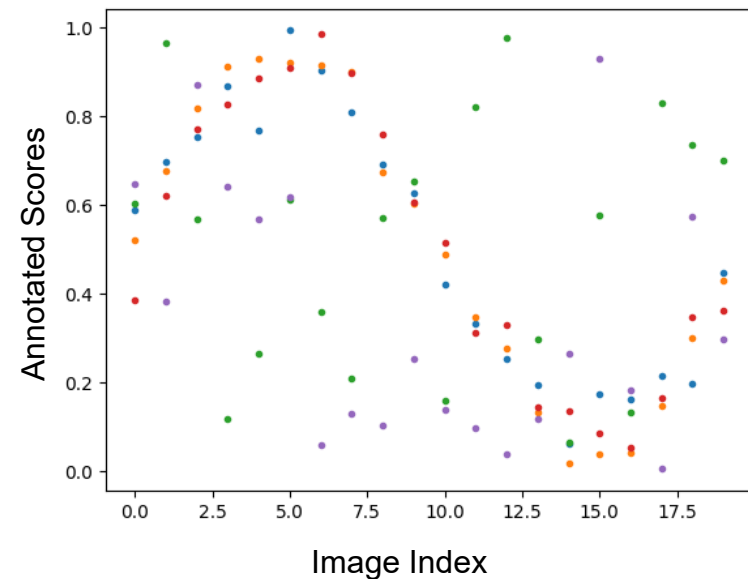
Applying MLE in code

$$\mu_i = \sum_a s_{ia} / M$$

$$\sigma^2 = \sum_{i,a} (s_{ia} - \mu_i)^2 / (M \cdot N)$$

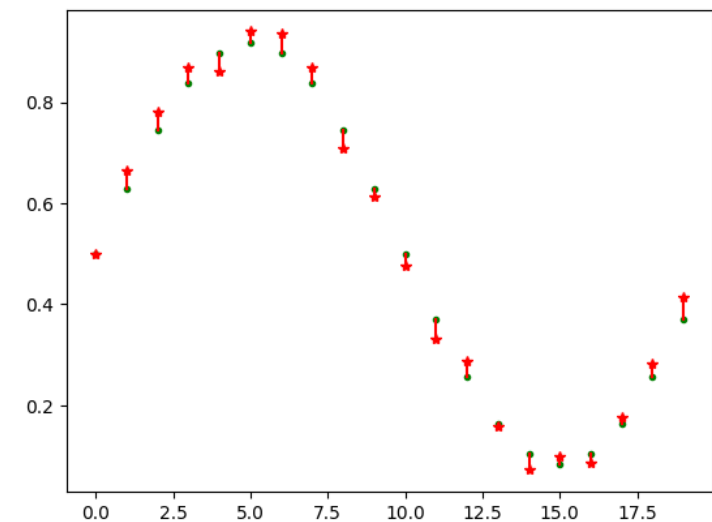
```
# initialize by assuming that all scores are good
score_mean = scores.mean(axis=1).reshape((len(scores), 1)) # mu_i
score_std = np.sqrt(np.sum((scores-score_mean)**2, axis=None)/N/M) # sigma
```

Assuming all the scores are good (when they are not)



(Green = true; red = prediction)

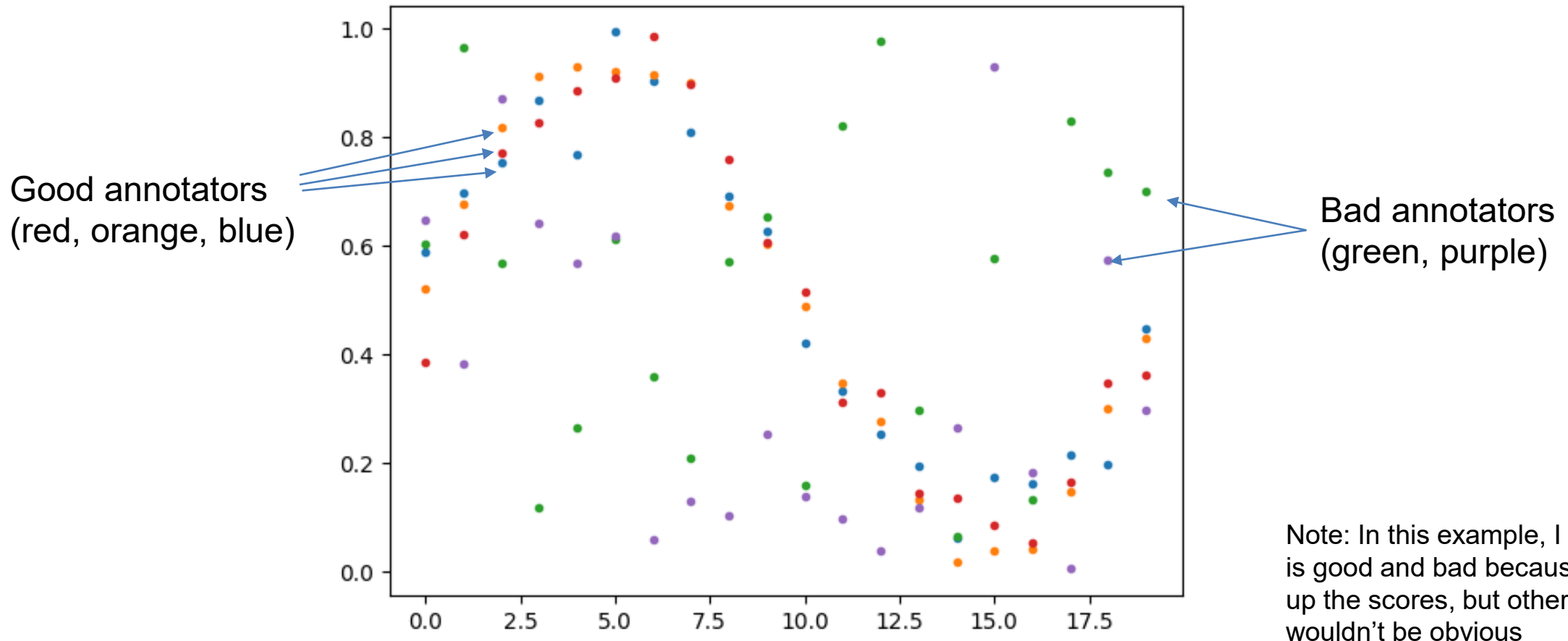
If we knew somehow which annotators are good



(Green = true; red = prediction)

Probabilistic Inference

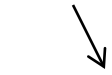
2. Given the distribution parameters, how do we compute the likelihood that an annotator is good?



Probabilistic Inference (general case)

Given the model parameters, compute the likelihood that a particular model generated a sample

component or label



$$p(z_n = m \mid x_n, \theta)$$

z_n is the unknown label of data point x_n

General strategy: We know $p(x_n \mid z_n = 0, \theta)$ and $p(x_n \mid z_n = 1, \theta)$ and $p(z_n = 1 \mid \theta)$

We want to know $p(z_n = m \mid x_n, \theta)$

Use probability rules to get from what we know to what we want to know.

Probabilistic Inference (general case)

Given the model parameters, compute the likelihood that a particular model generated a sample

component or label

$$p(z_n = m | x_n, \theta) = \frac{p(z_n = m, x_n | \theta_m)}{p(x_n | \theta)}$$

Rule of conditional probability

Probabilistic Inference (general case)

Given the model parameters, compute the likelihood that a particular model generated a sample

component or label

$$p(z_n = m | x_n, \theta) = \frac{p(z_n = m, x_n | \theta_m)}{p(x_n | \theta)}$$

$$= \frac{p(z_n = m, x_n | \theta_m)}{\sum_k p(z_n = k, x_n | \theta_k)}$$

Law of total probability

Probabilistic Inference (general case)

Given the model parameters, compute the likelihood that a particular model generated a sample

component or label

$$\begin{aligned} p(z_n = m | x_n, \theta) &= \frac{p(z_n = m, x_n | \theta_m)}{p(x_n | \theta)} \\ &= \frac{p(z_n = m, x_n | \theta_m)}{\sum_k p(z_n = k, x_n | \theta_k)} \\ &= \frac{p(x_n | z_n = m, \theta_m) p(z_n = m | \theta_m)}{\sum_k p(x_n | z_n = k, \theta_k) p(z_n = k | \theta_k)} \end{aligned}$$

Chain rule of probability

Example: Inference for Annotator Labels

$$p(z_a = 1 | s_a, \theta) = \frac{p(s_a | z_a = 1, \theta) p(z_a = 1, \theta)}{p(s_a | z_a = 1, \theta) p(z_a = 1, \theta) + p(s_a | z_a = 0, \theta) p(z_a = 0, \theta)}$$

$$p(s_a | z_a = 1, \theta) = \prod_i N(s_{ia}, \mu_i, \sigma) \text{ (normal pdf)}$$

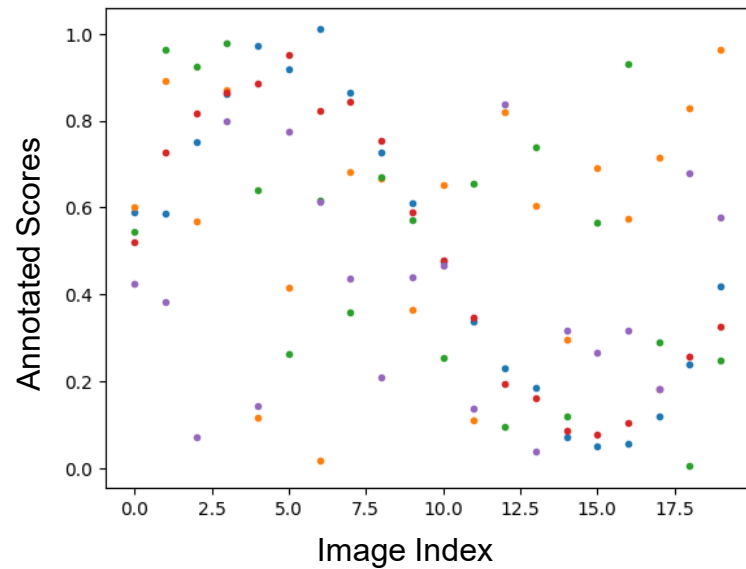
$$p(s_a | z_a = 0, \theta) = \prod_i 1 \text{ (uniform)}$$

$$p(z_a = 1, \theta) = \pi_z \text{ (prior)}$$

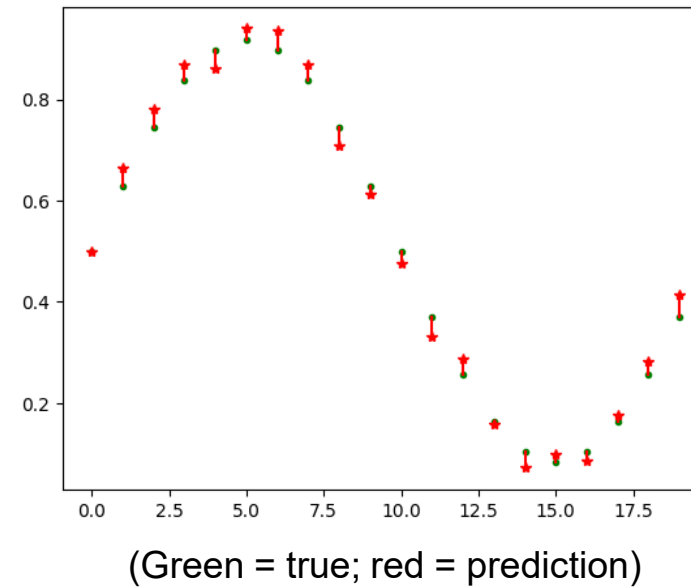
```
p_good = np.zeros((5,1)) # w_a = P(z_a=1 | scores, theta_t)
for a in range(M):
    p_s_good = pz # P(s_ia | z=1, mu_i, std)P(z_a=1)
    p_s_bad = 1-pz # P(s_ia | z=0)P(z_a=0)
    for i in range(N):
        p_s_good *= 1/np.sqrt(2*np.pi)/score_std * np.exp(-1/2 * (scores[i,a]-score_mean[i])**2/score_std**2)
        p_s_bad *= 1 # uniform in range [0, 1]
    p_good[a] = p_s_good / (p_s_good + p_s_bad)
```

Dealing with Latent Variables

3. How can we get annotator labels and score models at once?



Estimated scores



Good annotators: 0, 1, 3

[2 minute break]

- After the break:
 - Intuitive solution to solving for parameters and latent variables
 - Derive that the intuitive solution is correct
 - Demo for the untrustworthy annotator problem

Simple solution

1. Initialize parameters
2. Compute the probability of each hidden variable given the current parameters
3. Compute new parameters for each model, weighted by likelihood of hidden variables
4. Repeat 2-3 until convergence

Annotator Problem: Simple Solution

1. Initialize parameters

- Estimate parameters assuming all annotators are good

2. Compute likelihood of hidden variables for current parameters

$$w_a = p(z_a = 1 | \mathbf{s}_a, \boldsymbol{\mu}, \sigma, \pi_z)$$

3. Estimate new parameters for each model, weighted by likelihood

$$\hat{\mu}_i = \sum_a w_a s_{ia} / \sum_a w_a$$

$$\hat{\sigma}^2 = \sum_{i,a} (s_{ia} - \mu_i)^2 w_a / \sum_{i,a} w_a$$

$$\pi_z = \sum_a w_a / M$$

Expectation Maximization (EM) Algorithm

$$\text{Goal: } \hat{\theta} = \operatorname{argmax}_{\theta} \log \left(\sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z} | \theta) \right)$$

Log of sums is intractable

Jensen's Inequality

$$f(\mathbb{E}[X]) \geq \mathbb{E}[f(X)]$$

for concave functions $f(x)$

So we maximize the lower bound by maximizing the log of sums instead of sum of logs!

Expectation Maximization (EM) Algorithm

$$\text{Goal: } \hat{\theta} = \operatorname{argmax}_{\theta} \log \left(\sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z} | \theta) \right)$$

1. E-step: compute

$$E_{\mathbf{z}|\mathbf{x}, \theta^{(t)}} [\log(p(\mathbf{x}, \mathbf{z} | \theta))] = \sum_{\mathbf{z}} \log(p(\mathbf{x}, \mathbf{z} | \theta)) p(\mathbf{z} | \mathbf{x}, \theta^{(t)})$$

2. M-step: solve

$$\theta^{(t+1)} = \operatorname{argmax}_{\theta} \sum_{\mathbf{z}} \log(p(\mathbf{x}, \mathbf{z} | \theta)) p(\mathbf{z} | \mathbf{x}, \theta^{(t)})$$

Expectation Maximization (EM) Algorithm

log of expectation of $p(\mathbf{x}|\mathbf{z})$ over $p(\mathbf{z})$

$$\text{Goal: } \hat{\theta} = \underset{\theta}{\operatorname{argmax}} \log \left(\sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z} | \theta) \right) \quad f(\mathbb{E}[X]) \geq \mathbb{E}[f(X)]$$

1. E-step: compute $\text{expectation of log of } P(\mathbf{x}|\mathbf{z}) \text{ over estimated } P(\mathbf{z})$

$$\mathbb{E}_{\mathbf{z}|\mathbf{x}, \theta^{(t)}} [\log(p(\mathbf{x}, \mathbf{z} | \theta))] = \sum_{\mathbf{z}} \log(p(\mathbf{x}, \mathbf{z} | \theta)) p(\mathbf{z} | \mathbf{x}, \theta^{(t)})$$

2. M-step: solve

$$\theta^{(t+1)} = \underset{\theta}{\operatorname{argmax}} \sum_{\mathbf{z}} \log(p(\mathbf{x}, \mathbf{z} | \theta)) p(\mathbf{z} | \mathbf{x}, \theta^{(t)})$$

EM for Annotator Problem: **E-Step**

$$E\text{-step: } E_{z|s, \theta^{(k)}}(\log P(s, z | \theta)) = \sum_z \log P(s, z | \theta) P(z | s, \theta^{(k)})$$

μ, σ^2, π_z

Use $\hat{\theta}$ as current estimate of θ

$$\sum_z \log P(s, z | \hat{\theta}) P(z | s, \hat{\theta})$$

These sums over all combinations of z and joint probabilities are complicated!

First, let's expand the joint probabilities

EM for Annotator Problem: E-Step

Write joint as product of individual

Log of product is sum of logs

$$E\text{-step: } E_{z|s, \theta^{(k)}}(\log P(s, z | \theta)) = \sum_z \log P(s, z | \theta) P(z | s, \theta^{(k)})$$

Use $\hat{\theta}$ as current estimate of θ

$$\sum_z \log P(s, z | \theta) P(z | s, \hat{\theta})$$

$$= \sum_z \left[\log \prod_{i,a} P(s_{i,a} | z_a, \theta) P(z_a | \theta) \right] P(z | s, \hat{\theta})$$

$$= \sum_z \left[\sum_{i,a} \log (P(s_{i,a} | z_a, \theta) P(z_a | \theta)) \right] P(z | s, \hat{\theta})$$

Now I want to deal with this joint z and complex sum over z

EM for Annotator Problem: E-Step

$$E\text{-step: } E_{z|s, \theta^{(k)}}(\log P(s, z | \theta)) = \sum_z \log P(s, z | \theta) P(z | s, \theta^{(k)})$$

Use $\hat{\theta}$ as current estimate of θ

$$\sum_z \log P(s, z | \theta) P(z | s, \hat{\theta})$$

$$= \sum_z \left[\log \prod_{i,a} P(s_{i,a} | z_a, \theta) P(z_a | \theta) \right] P(z | s, \hat{\theta})$$

$$= \sum_z \left[\sum_{i,a} \log (P(s_{i,a} | z_a, \theta) P(z_a | \theta)) \right] P(z | s, \hat{\theta})$$

$$= \sum_{i,a} \sum_z \left[\log (P(s_{i,a} | z_a, \theta) P(z_a | \theta)) \right] P(z | s, \hat{\theta})$$

$$= \sum_{i,a} \sum_{z_0} \dots \sum_{z_m} \left[\log (P(s_{i,a} | z_a, \theta) P(z_a | \theta)) \right] P(z | s, \hat{\theta})$$

$$= \sum_{i,a} \sum_{z_a} \left[\log P(s_{i,a} | z_a, \theta) + \log P(z_a | \theta) \right] \sum_{z_{a'}} P(z | s, \hat{\theta})$$

Rearrange sums

Make it clear that sum over z is a series of sums

Rearrange again, pulling out the sum over z_a

This lets me marginalize out the other z 's!

EM for Annotator Problem: E-Step

$$E\text{-step: } E_{z|s, \theta^{(k)}}(\log P(s, z | \theta)) = \sum_z \log P(s, z | \theta) P(z | s, \theta^{(k)})$$

Use $\hat{\theta}$ as current estimate of θ

$$\sum_z \log P(s, z | \hat{\theta}) P(z | s, \hat{\theta})$$

$$= \sum_z \left[\log \prod_{i,a} P(s_{i,a} | z_a, \theta) P(z_a | \theta) \right] P(z | s, \hat{\theta})$$

$$= \sum_z \left[\sum_{i,a} \log (P(s_{i,a} | z_a, \theta) P(z_a | \theta)) \right] P(z | s, \hat{\theta})$$

$$= \sum_{i,a} \sum_z \left[\right]$$

$$= \sum_{i,a} \sum_{z_0} \dots \sum_{z_m} \left[\right]$$

$$= \sum_{i,a} \sum_{z_a} \left[\log P(s_{i,a} | z_a, \theta) + \log P(z_a | \theta) \right] \sum_{z \neq z_a} P(z | s, \hat{\theta})$$

$$= \sum_{i,a} \left[\log \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2} \frac{(s_{i,a} - \mu_i)^2}{\sigma^2}} + \log \pi_z \right] P(z_a = 1 | s, \hat{\theta}) + \left[\log(1) + \log(1 - \pi_z) \right] P(z_a = 0 | s, \hat{\theta})$$

$$P(z_a = 1 | s, \hat{\theta}) = P(z_a = 1 | s_a, \hat{\theta})$$

$$= P(z_a = 1, s_a | \hat{\theta}) / P(s_a | \hat{\theta})$$

$$= \frac{P(s_a | z_a = 1, \hat{\theta}) P(z_a = 1 | \hat{\theta})}{P(s_a | z_a = 1, \hat{\theta}) P(z_a = 1 | \hat{\theta}) + P(s_a | z_a = 0, \hat{\theta}) P(z_a = 0 | \hat{\theta})}$$

Now, plug in the functions for score probabilities of good and bad annotators

And calculate how to get $P(z|s)$, as we did before

EM for Annotator

Problem: E-Step

$$\begin{aligned}P(z_a=1 | \underline{s}, \hat{\theta}) &= P(z_a=1 | \underline{s}_a, \hat{\theta}) \\&= P(z_a=1, \underline{s}_a | \hat{\theta}) / P(\underline{s}_a | \hat{\theta}) \\&= P(\underline{s}_a | z_a=1, \hat{\theta}) P(z_a=1 | \hat{\theta}) / [P(\underline{s}_a | z_a=1, \hat{\theta}) P(z_a=1 | \hat{\theta}) \\&\quad + P(\underline{s}_a | z_a=0, \hat{\theta}) P(z_a=0 | \hat{\theta})]\end{aligned}$$

Writing out the
inference
computations

$$P(\underline{s}_a | z_a=1, \hat{\theta}) = \prod_i P(s_{ia} | z_a=1, \hat{\theta}) = \prod_i N(s_{ia}; \hat{\mu}_i, \hat{\sigma}_i^2)$$

$$P(z_a=1 | \hat{\theta}) = \pi_{1z} \quad P(z_a=0 | \hat{\theta}) = 1 - \pi_{1z}$$

$$P(\underline{s}_a | z_a=0, \hat{\theta}) = \prod_i P(s_{ia} | z_a=0, \hat{\theta}) = 1$$

$$\text{Let } P(z_a=1 | \underline{s}, \hat{\theta}) = w_a$$

EM for Annotator Problem: M-Step

$$\text{Let } P(z_a=1 | S, \hat{\theta}) = w_a$$

Calculate parameters that maximize the expression from the E-step, given our current estimates of $P(z|s)$

This is very similar to the MLE derivation

M-Step: find parameters that maximize expected log likelihood

$$\frac{\partial}{\partial \mu_i} \sum_{i,a} \left[\log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{2} \frac{(s_{ia} - \mu_i)^2}{\sigma^2} + \log \pi_z \right] w_a + \log(1 - \pi_z)(1 - w_a) = 0$$

$$\sum_a \left[\frac{(s_{ia} - \mu_i)}{\sigma^2} \right] w_a = 0 \Rightarrow \hat{\mu}_i = \frac{\sum_a w_a s_{ia}}{\sum_a w_a}$$

$$\frac{\partial}{\partial \sigma} \sum_{i,a} \left[\log \frac{1}{\sigma} - \frac{1}{2} \frac{(s_{ia} - \mu_i)^2}{\sigma^2} \right] w_a = 0$$

$$\sum_{i,a} \left[-\frac{1}{\sigma} + \frac{(s_{ia} - \mu_i)^2}{\sigma^3} \right] w_a = 0 \Rightarrow \sum_{i,a} \left[-\sigma^2 + (s_{ia} - \mu_i)^2 \right] w_a = 0$$

$$\sigma^2 = \frac{\sum_{i,a} (s_{ia} - \mu_i)^2 w_a}{\sum_{i,a} w_a}$$

$$\frac{\partial}{\partial \pi_z} \sum_{i,a} \log \pi_z w_a + \log(1 - \pi_z)(1 - w_a) = 0$$

$$\sum_{i,a} \frac{1}{\pi_z} w_a - \frac{1}{1 - \pi_z} (1 - w_a) = 0 \Rightarrow \sum_{i,a} (1 - \pi_z) w_a - \pi_z (1 - w_a) = 0$$

$$\sum_{i,a} [w_a - \pi_z w_a - \pi_z + \pi_z w_a] = 0$$

$$\pi_z = \frac{\sum_{i,a} w_a}{\sum_{i,a} 1} = \frac{\sum_a w_a}{M}$$

$\hat{\mu}_i$

$\hat{\sigma}$

$\hat{\pi}_z$

EM Annotator Problem Demo

<https://colab.research.google.com/drive/1sutnFg-xe-ljgiY8qAJt5USf2MEBZS2L?usp=sharing>

EM Algorithm

- Maximizes a lower bound on the data likelihood at each iteration
- Each step increases the data likelihood
 - Converges to *local maximum*
- Common tricks to derivation
 - Find terms that sum or integrate to 1
 - Lagrange multiplier to deal with constraints
- Although the derivation is long, it pretty much always boils down to iteratively
 1. Estimating likelihood of latent variables given parameters
 2. Computing estimates of parameters that are weighted by the latent variable likelihoods

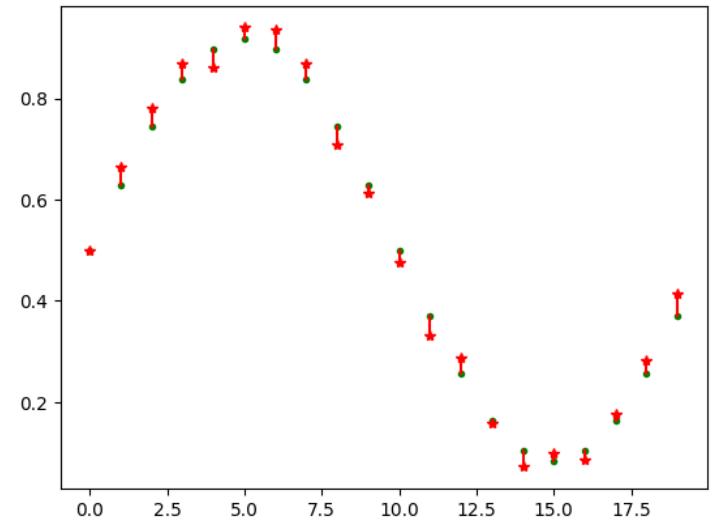
“Hard EM”

- Same as EM except compute z^* as most likely values for hidden variables
- K-means is an example
- Advantages
 - Simpler: can be applied when cannot derive EM
 - Sometimes works better if you want to make hard predictions at the end
- But
 - Generally, pdf parameters are not as accurate as EM

What to remember

- EM is a widely applicable algorithm to solve for latent variables and parameters that make the observed data likely
- While derivation is long and somewhat complicated, the application is simple
- We'll see other examples of EM use in mixture of Gaussian and topic models

Estimated scores



(Green = true; red = prediction)

Good annotators: 0, 1, 3

Next class

- Estimating probability density functions