

# Clustering

CS 229: Machine Learning

Emily Fox

Stanford University

February 28, 2024

©2024 Emily Fox

1

Motivating clustering approaches

©2024 Emily Fox

CS 229: Machine Learning

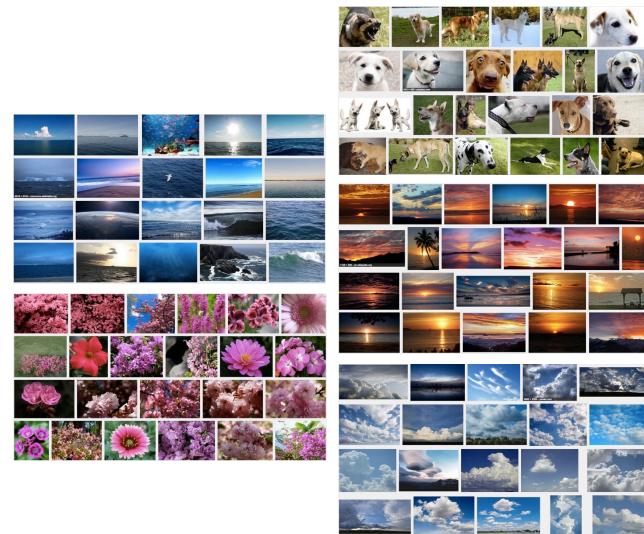
2

1

# Clustering images

Discover groups of similar images:

- Ocean
- Pink flower
- Dog
- Sunset
- Clouds
- ...

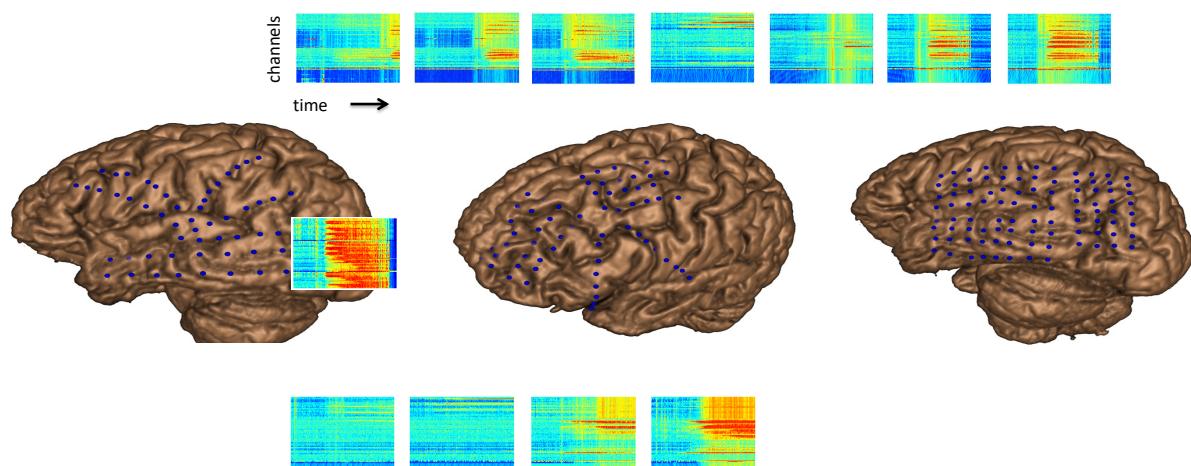


©2024 Emily Fox

CS 229: Machine Learning

3

# Characterizing patients and seizure types

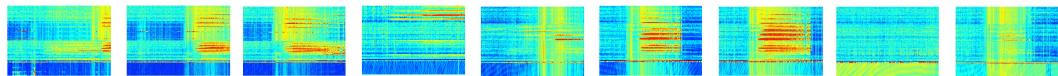


©2024 Emily Fox

CS 229: Machine Learning

4

# Characterizing patients and seizure types



©2024 Emily Fox

CS 229: Machine Learning

5

Clustering: An unsupervised learning task

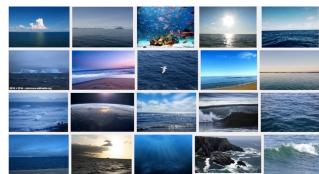
©2024 Emily Fox

CS 229: Machine Learning

6

3

## What if labels are known for training data?



Ocean



Dog



Sunset



Pink flower



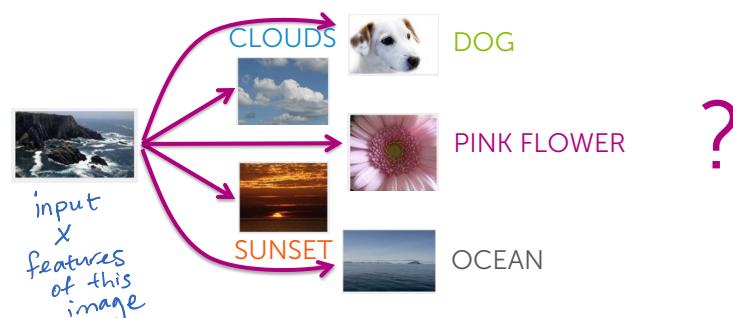
Clouds

©2024 Emily Fox

CS 229: Machine Learning

7

## Multiclass classification problem



Example of  
**supervised learning**

©2024 Emily Fox

CS 229: Machine Learning

8

# Clustering

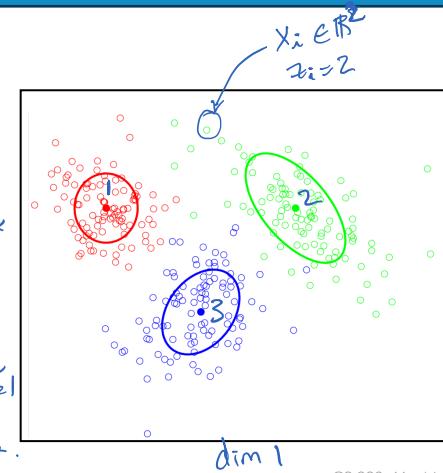
No labels provided

...uncover cluster structure from input alone

**Input:** image features  $\mathbf{x}_i$

**Output:** cluster labels  $z_i$

An unsupervised learning task



©2024 Emily Fox

CS 229: Machine Learning

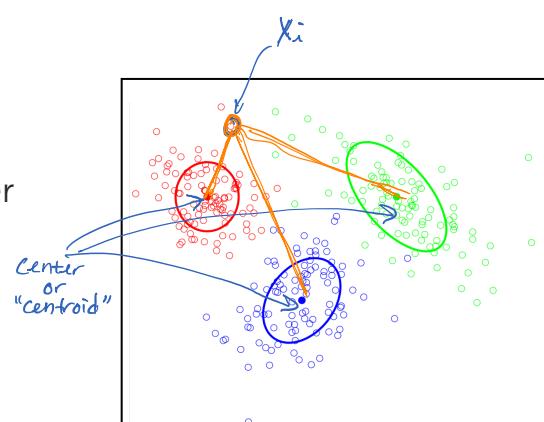
9

## What defines a cluster?

Cluster defined by center & shape/spread

Assign observation  $\mathbf{x}_i$  to cluster k if

- Score under cluster k is higher than under others
- For simplicity, often define score as distance to cluster center (ignoring shape)



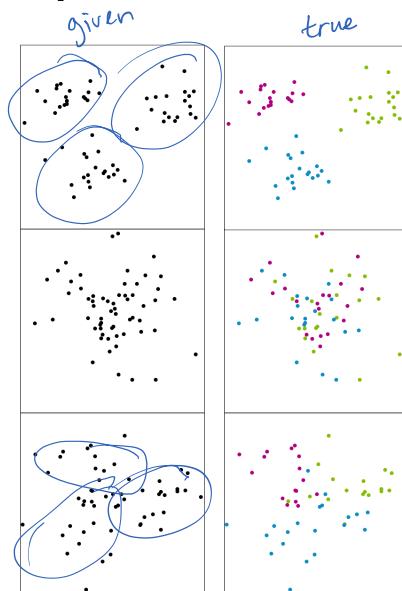
©2024 Emily Fox

CS 229: Machine Learning

10

## Hope for unsupervised learning

Easy



Impossible

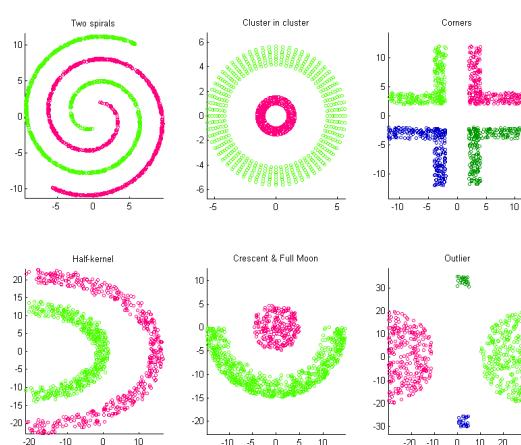
In between

©2024 Emily Fox

CS 229: Machine Learning

11

## Other (challenging!) clusters to discover...

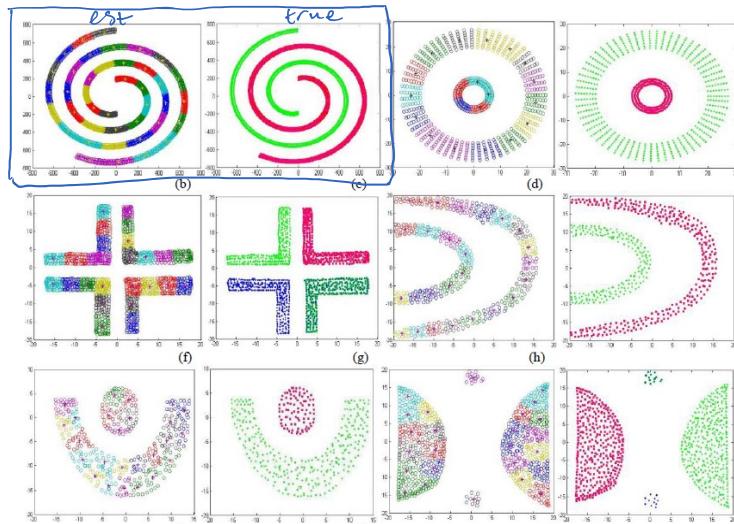


©2024 Emily Fox

CS 229: Machine Learning

12

## Other (challenging!) clusters to discover...



©2024 Emily Fox

CS 229: Machine Learning

13

## k-means: A clustering algorithm

©2024 Emily Fox

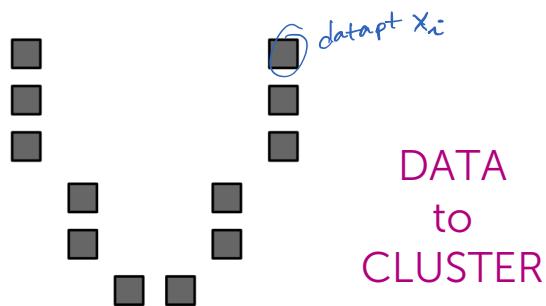
CS 229: Machine Learning

14

## k-means

Assume

- Score= distance to cluster center (smaller better)



©2024 Emily Fox

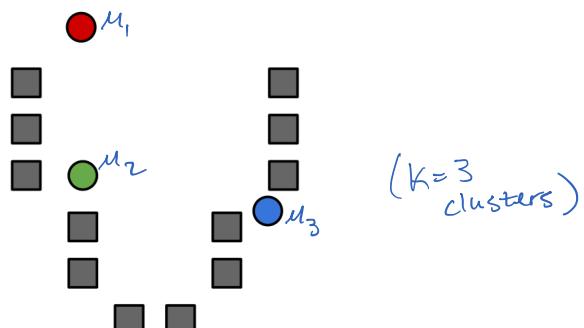
CS 229: Machine Learning

15

## k-means algorithm

0. Initialize cluster centers

$$\mu_1, \mu_2, \dots, \mu_k$$



©2024 Emily Fox

CS 229: Machine Learning

16

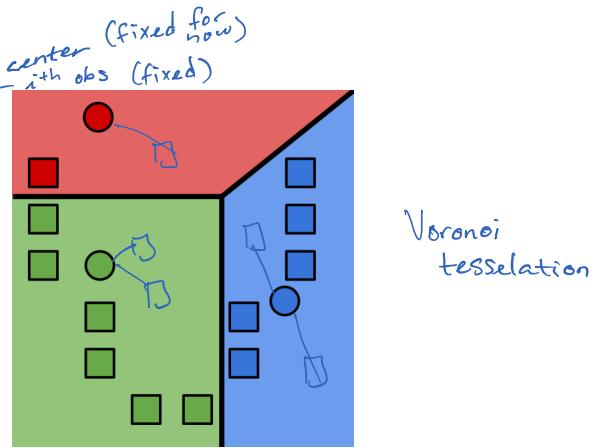
## k-means algorithm

0. Initialize cluster centers

1. Assign observations to closest cluster center

$$z_i \leftarrow \arg \min_j \|\mu_j - \mathbf{x}_i\|_2^2$$

Inferred label for obs i, whereas supervised learning has given label  $y_i$   
 Select the cluster index closest to  $\mathbf{x}_i$



©2024 Emily Fox

CS 229: Machine Learning

17

## k-means algorithm

0. Initialize cluster centers

1. Assign observations to closest cluster center

2. Revise cluster centers as mean of assigned observations

$$\mu_j = \frac{1}{n_j} \sum_{i:z_i=j} \mathbf{x}_i$$

center of mass per cluster  
 total # of obs in cluster j  
 all obs. assigned to cluster j

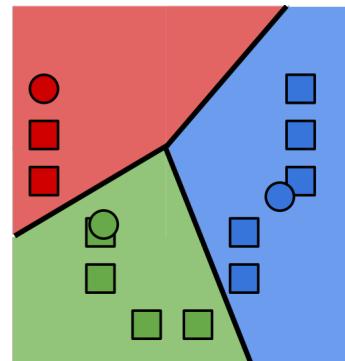
©2024 Emily Fox

CS 229: Machine Learning

18

## k-means algorithm

0. Initialize cluster centers
1. Assign observations to closest cluster center
2. Revise cluster centers as mean of assigned observations
3. Repeat 1.+2. until convergence



©2024 Emily Fox

CS 229: Machine Learning

19

## k-means as coordinate descent

©2024 Emily Fox

CS 229: Machine Learning

20

## Recall: Coordinate descent

Goal: Minimize some function  $g$

$$g(\mathbf{w}) = g(w_0, w_1, \dots, w_D)$$

$$\min_{\mathbf{w}} g(\mathbf{w})$$

when keeping others fixed

Often, hard to find minimum for all coordinates, but **easy for each coordinate**

**Coordinate descent:**

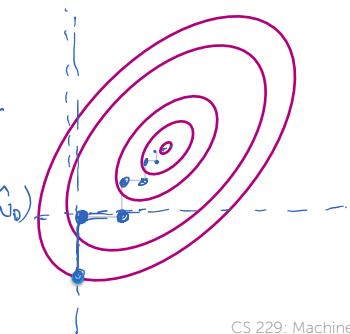
Initialize  $\hat{\mathbf{w}} = 0$  (or smartly...)

while not converged

pick a coordinate  $j$

$$\hat{w}_j \leftarrow \min_{\hat{w}} g(\hat{w}_0, \hat{w}_1, \dots, \hat{w}_{j-1}, \hat{w}, \hat{w}_{j+1}, \dots, \hat{w}_D)$$

©2024 Emily Fox



CS 229: Machine Learning

21

## k-means as a coordinate descent algorithm

1. Assign observations to closest cluster center

$$z_i \leftarrow \arg \min_j \|\mu_j - \mathbf{x}_i\|_2^2$$

2. Revise cluster centers as mean of assigned observations

$$\mu_j = \frac{1}{n_j} \sum_{i:z_i=j} \mathbf{x}_i \quad \text{equivalent to}$$

$$\mu_j \leftarrow \arg \min_{\mu} \sum_{i:z_i=j} \|\mu - \mathbf{x}_i\|_2^2$$

©2024 Emily Fox

CS 229: Machine Learning

22

## k-means as a coordinate descent algorithm

1. Assign observations to closest cluster center

$$z_i \leftarrow \arg \min_j \|\mu_j - \mathbf{x}_i\|_2^2$$

of  $\{w_0, w_1, \dots, w_D\}$   
 $\{x_1, x_2, x_3\}$

2. Revise cluster centers as mean of assigned observations

$$\mu_j \leftarrow \arg \min_{\mu} \sum_{i:z_i=j} \|\mu - \mathbf{x}_i\|_2^2$$

Overall objective:  
 $\min_{\{z_i\}, \{\mu_j\}} \sum_{i=1}^N \|\mu_{z_i} - \mathbf{x}_i\|_2^2$

~~$\min_{\{z_i\}, \{\mu_j\}} \sum_{j=1}^K \sum_{i:z_i=j} \|\mu_j - \mathbf{x}_i\|_2^2$~~

©2024 Emily Fox

CS 229: Machine Learning

23

## k-means as a coordinate descent algorithm

1. Assign observations to closest cluster center

$$z_i \leftarrow \arg \min_j \|\mu_j - \mathbf{x}_i\|_2^2$$

2. Revise cluster centers as mean of assigned observations

$$\mu_j \leftarrow \arg \min_{\mu} \sum_{i:z_i=j} \|\mu - \mathbf{x}_i\|_2^2$$

Alternating minimization  
 1. ( $z$  given  $\mu$ ) and 2. ( $\mu$  given  $z$ )  
**= coordinate descent**

©2024 Emily Fox

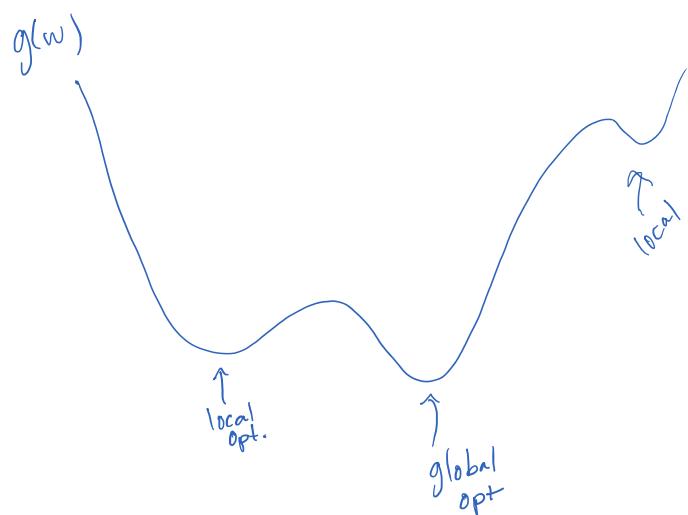
CS 229: Machine Learning

24

## Convergence of k-means

Converges to:

- Global optimum ~~Global optimum~~
- Local optimum Local optimum
- neither ~~neither~~

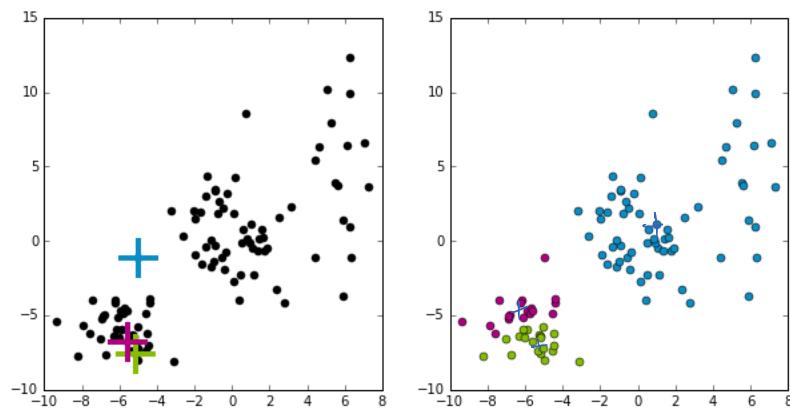


©2024 Emily Fox

CS 229: Machine Learning

25

## Convergence of k-means to local mode

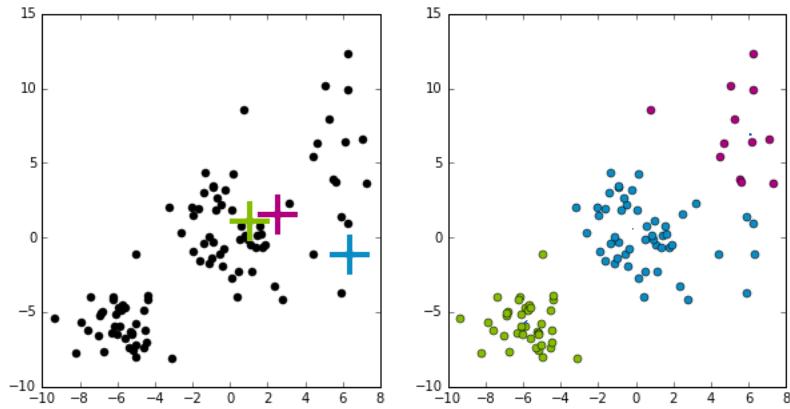


©2024 Emily Fox

CS 229: Machine Learning

26

## Convergence of k-means to local mode

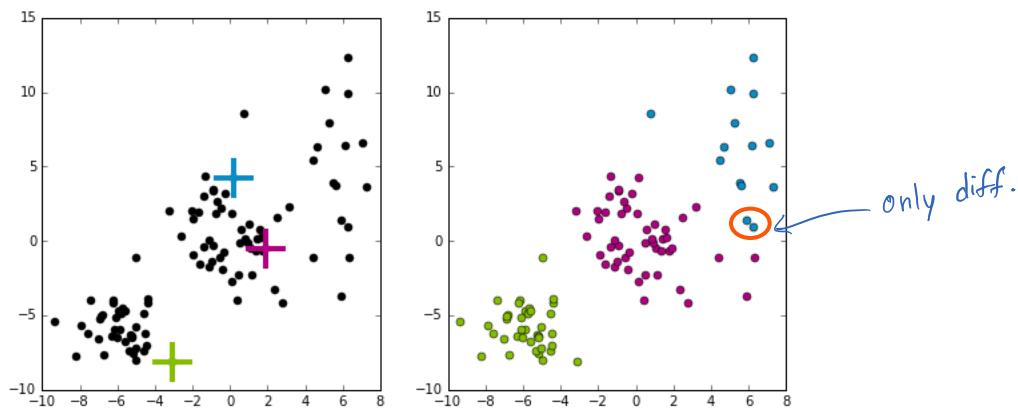


©2024 Emily Fox

CS 229: Machine Learning

27

## Convergence of k-means to local mode



©2024 Emily Fox

CS 229: Machine Learning

28

## Smart initialization with k-means++

©2024 Emily Fox

CS 229: Machine Learning

29

## k-means++ overview

Initialization of k-means algorithm is critical to quality of local optima found

### Smart initialization:

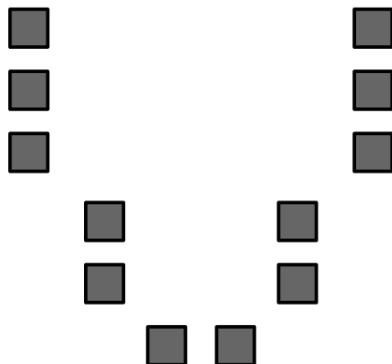
1. Choose first cluster center uniformly at random from data points
2. For each obs  $\mathbf{x}$ , compute distance  $d(\mathbf{x})$  to nearest cluster center
3. Choose new cluster center from amongst data points, with probability of  $\mathbf{x}$  being chosen proportional to  $d(\mathbf{x})^2$
4. Repeat Steps 2 and 3 until  $k$  centers have been chosen

©2024 Emily Fox

CS 229: Machine Learning

30

## k-means++ visualized

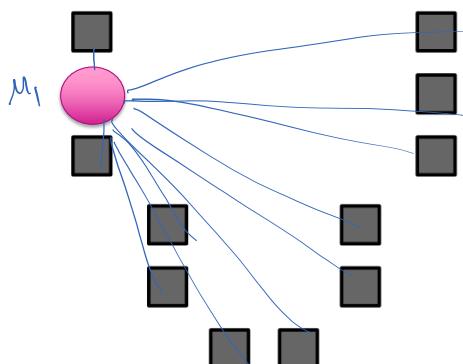


©2024 Emily Fox

CS 229: Machine Learning

31

## k-means++ visualized

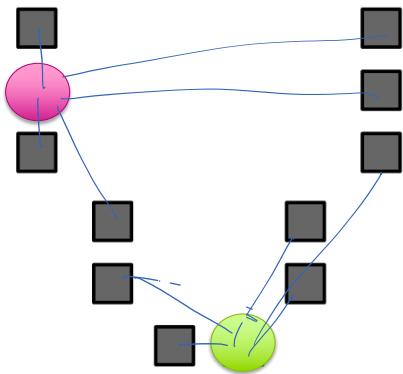


©2024 Emily Fox

CS 229: Machine Learning

32

## k-means++ visualized

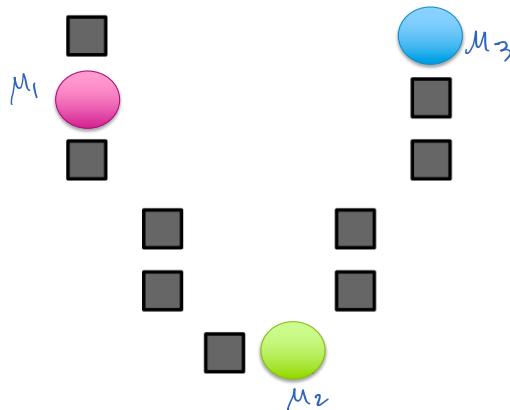


©2024 Emily Fox

CS 229: Machine Learning

33

## k-means++ visualized



©2024 Emily Fox

CS 229: Machine Learning

34

## k-means++ pros/cons

Computationally costly relative to random initialization, but the subsequent k-means often converges more rapidly

Tends to **improve quality of local optimum** and **lower runtime**

©2024 Emily Fox

CS 229: Machine Learning

35

Assessing quality of the clustering  
and choosing the # of clusters

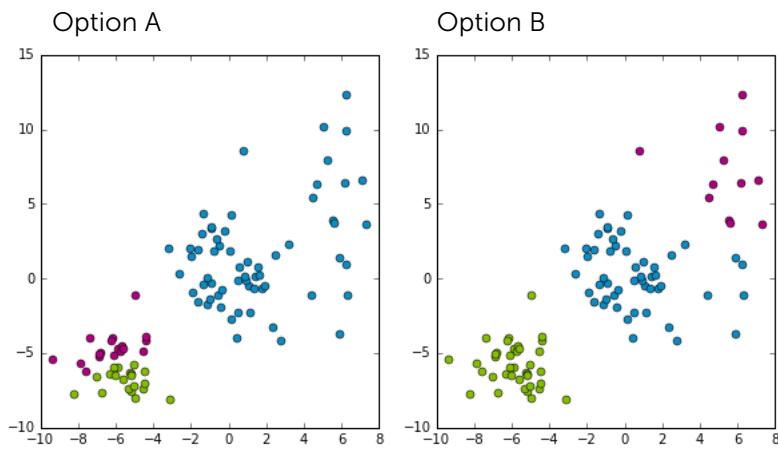
©2024 Emily Fox

CS 229: Machine Learning

36

18

## Which clustering does k-means prefer?

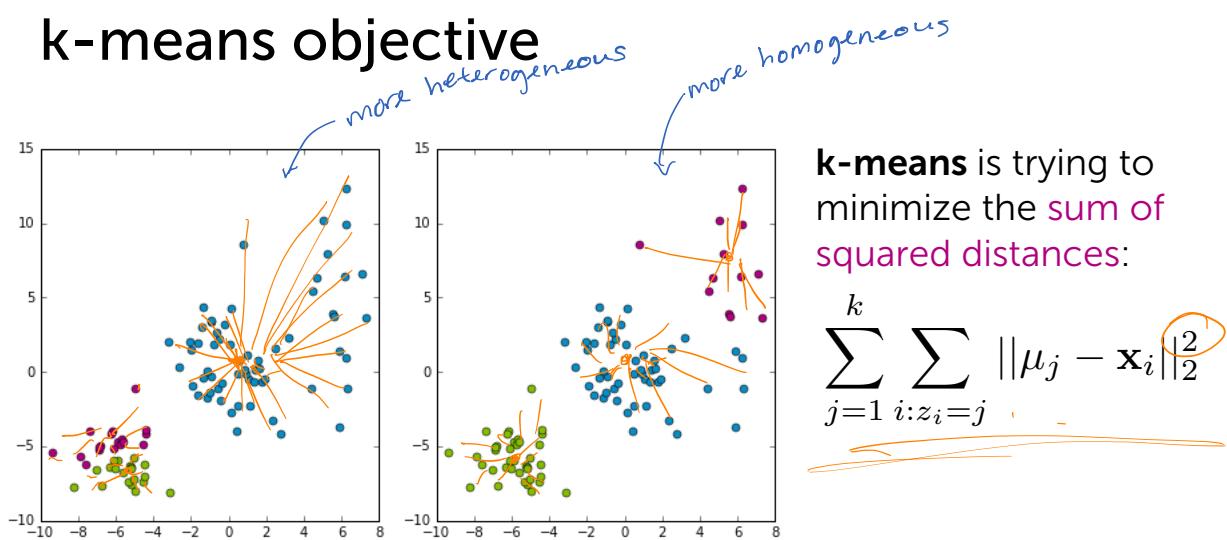


©2024 Emily Fox

CS 229: Machine Learning

37

## k-means objective

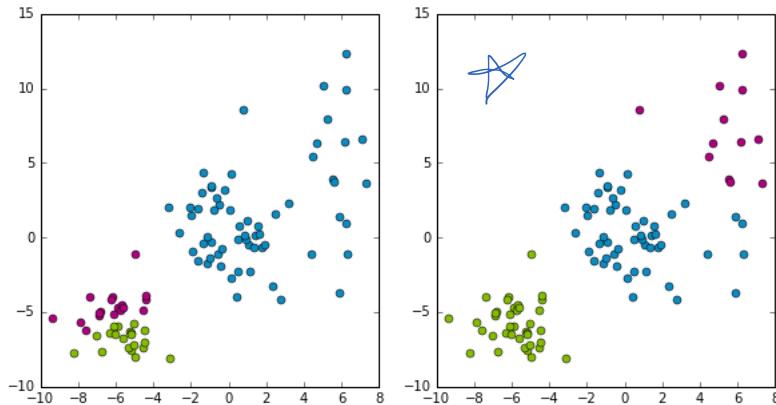


©2024 Emily Fox

CS 229: Machine Learning

38

## Cluster heterogeneity



Measure of **quality** of given clustering:

$$\sum_{j=1}^k \sum_{i:z_i=j} \|\mu_j - \mathbf{x}_i\|_2^2$$

Lower is better!

©2024 Emily Fox

CS 229: Machine Learning

39

## What happens as k increases?

Can refine clusters more and more to the data → **overfitting!**

**Extreme case** of  $k=N$ :

- can set each cluster center equal to datapoint
- heterogeneity =  $O$ ! (all distances to nearest cluster center are =  $O$ )

Lowest possible cluster heterogeneity decreases with increasing k

©2024 Emily Fox

CS 229: Machine Learning

40

## How to choose k?



©2024 Emily Fox

CS 229: Machine Learning

41

## Limitations & failure modes of k-means

©2024 Emily Fox

CS 229: Machine Learning

42

# Uncertainty in cluster assignments

Cluster 1

Cluster 2

Cluster 3

Cluster 4

Hard assignments don't tell full story

Slightly closer to Cluster 4 than Cluster 2, but count fully for Cluster 4?

©2024 Emily Fox

CS 229: Machine Learning

43

# Other limitations of k-means

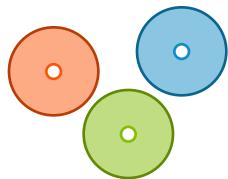
Assign observations to closest cluster center

$$z_i \leftarrow \arg \min_j \|\mu_j - \mathbf{x}_i\|_2^2$$

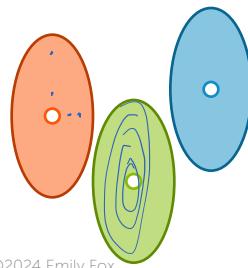
Can use weighted Euclidean, but requires *known* weights

Only center matters

Equivalent to assuming spherically symmetric clusters



Still assumes all clusters have the same axis-aligned ellipses

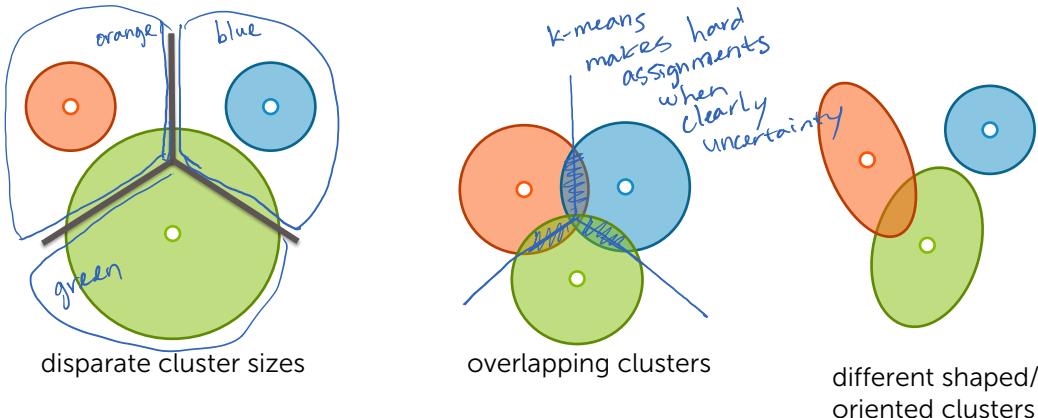


©2024 Emily Fox

CS 229: Machine Learning

44

## Failure modes of k-means



©2024 Emily Fox

CS 229: Machine Learning

45

## Motivates probabilistic model: Mixture model

- Provides **soft assignments** of observations to clusters (uncertainty in assignment)
  - e.g., 54% chance image is clouds, 45% sunset, 1% dog, and 0% pink flower
- Accounts for cluster shapes not just centers
- Enables learning weightings of dimensions
  - e.g., how much to weight each feature when computing cluster assignment

©2024 Emily Fox

CS 229: Machine Learning

46

## Summary for k-means

©2024 Emily Fox

CS 229: Machine Learning

47

## What you can do now...

- Describe the input (unlabeled observations) and output (labels) of a clustering algorithm
- Determine whether a task is supervised or unsupervised
- Cluster observations using k-means
- Cast k-means as a coordinate descent algorithm and discuss convergence
- Describe potential applications of clustering

©2024 Emily Fox

CS 229: Machine Learning

48

24

# Mixture Models: Model-Based Clustering

CS 229: Machine Learning

Emily Fox

Stanford University

February 28, 2024

©2024 Emily Fox

49

Again: Clustering images

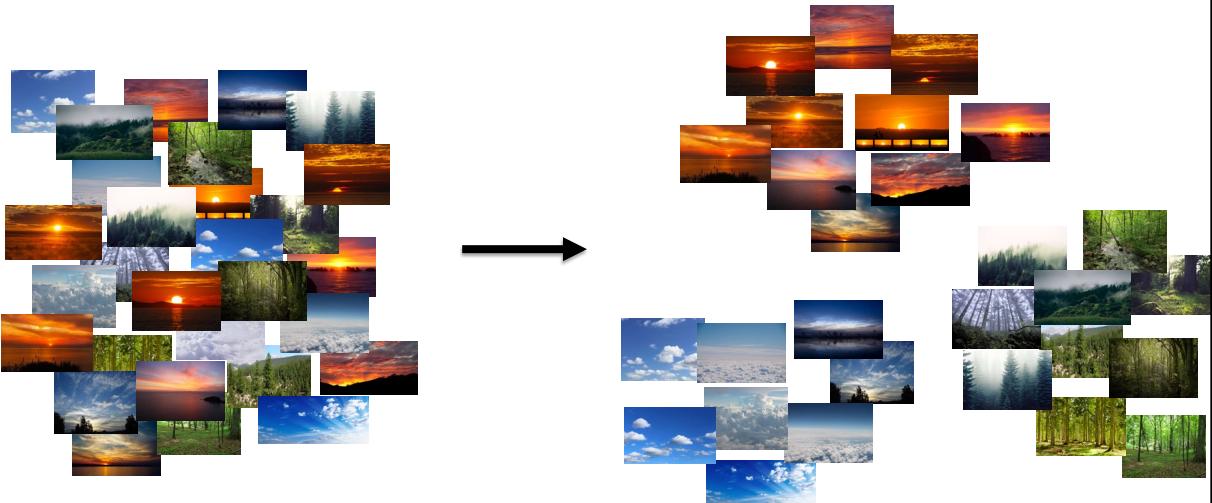
©2024 Emily Fox

CS 229: Machine Learning

50

25

## Case study: Clustering images



©2024 Emily Fox

CS 229: Machine Learning

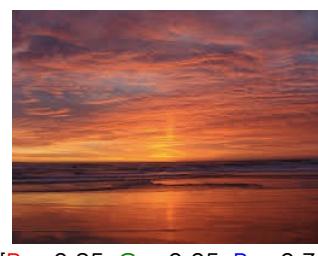
51

## Simple image representation

Consider average red, green, blue pixel intensities



[R = 0.05, G = 0.7, B = 0.9]



[R = 0.85, G = 0.05, B = 0.35]



[R = 0.02, G = 0.95, B = 0.4]

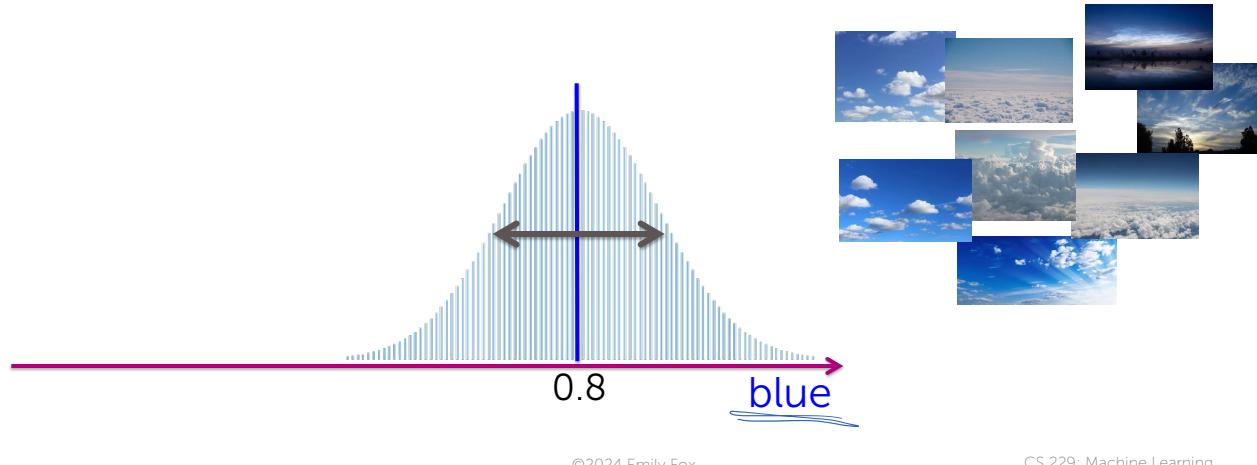
©2024 Emily Fox

CS 229: Machine Learning

52

## Distribution over all **cloud** images

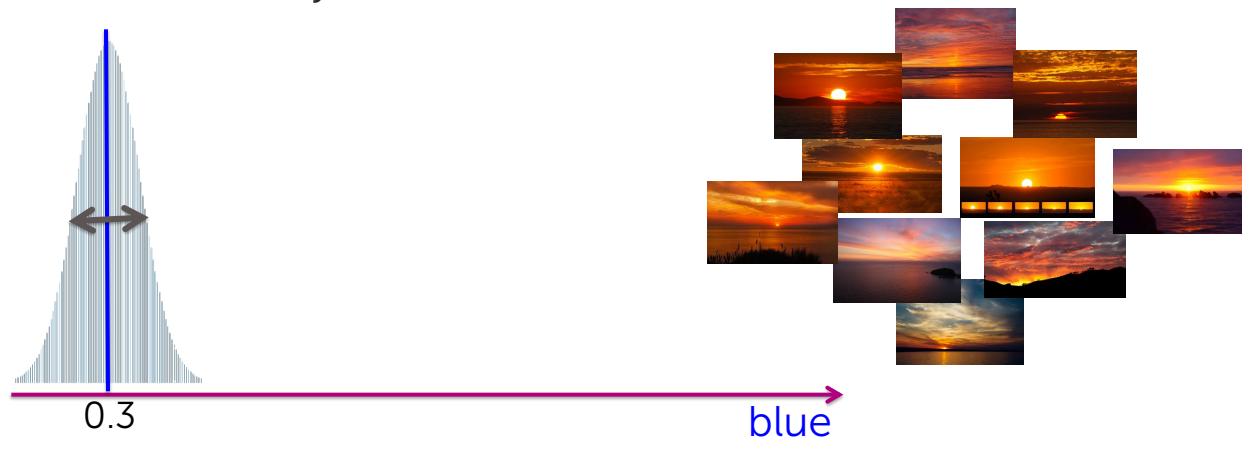
Let's look at just the **blue** dimension



53

## Distribution over all **sunset** images

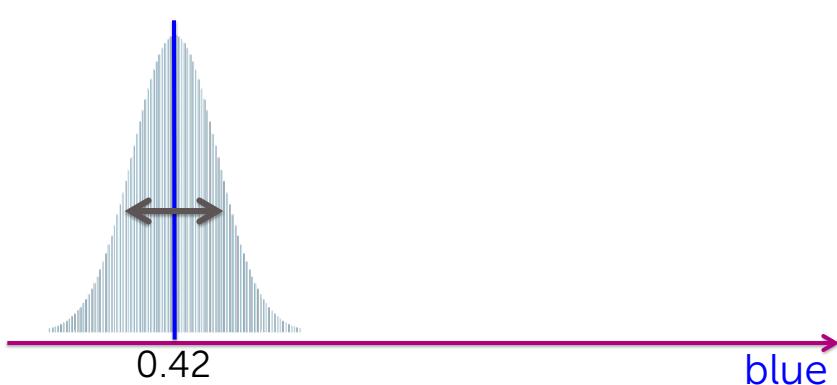
Let's look at just the **blue** dimension



54

## Distribution over all forest images

Let's look at just the **blue** dimension

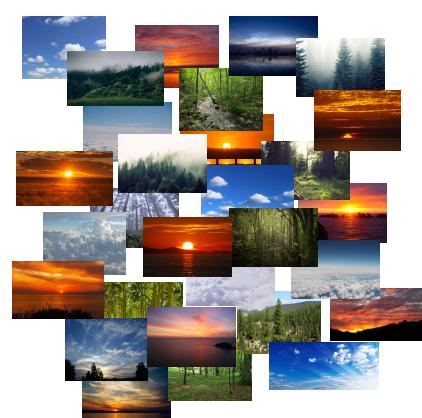
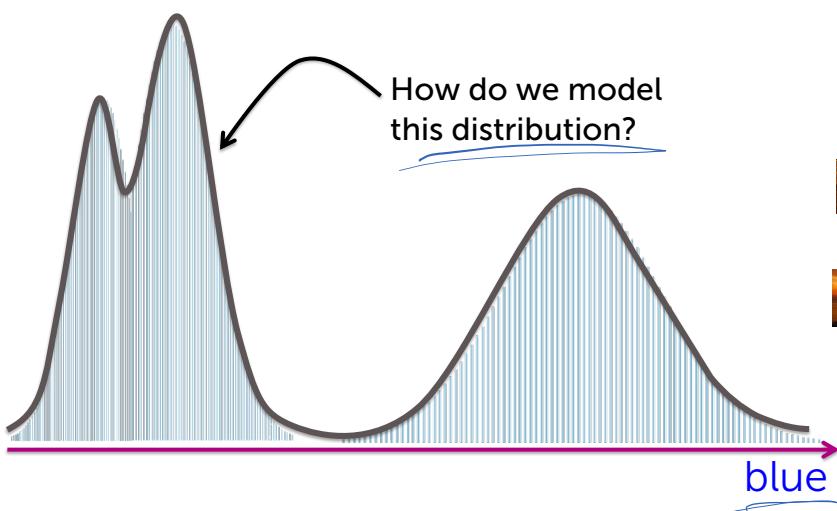


©2024 Emily Fox

CS 229: Machine Learning

55

## Distribution over **all** images



©2024 Emily Fox

CS 229: Machine Learning

56

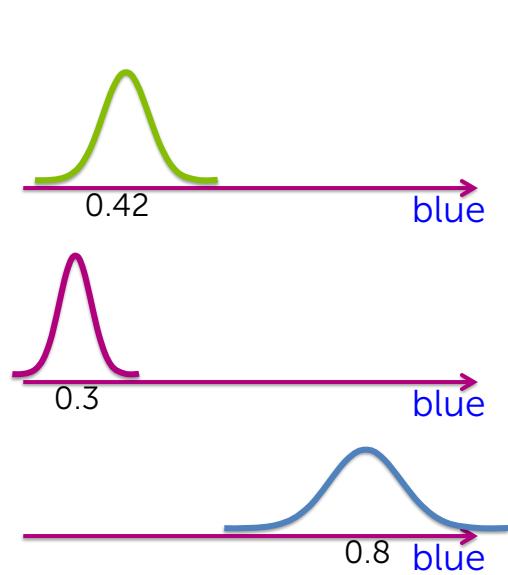
## Mixture of Gaussians

©2024 Emily Fox

CS 229: Machine Learning

57

## Model as Gaussian per cluster



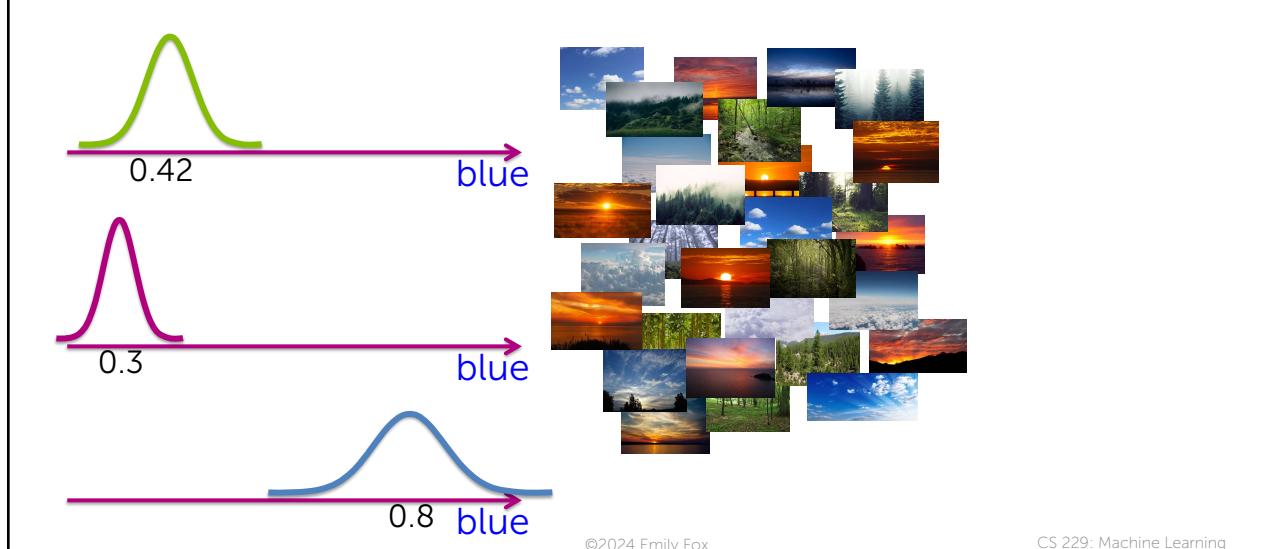
©2024 Emily Fox

CS 229: Machine Learning

58

29

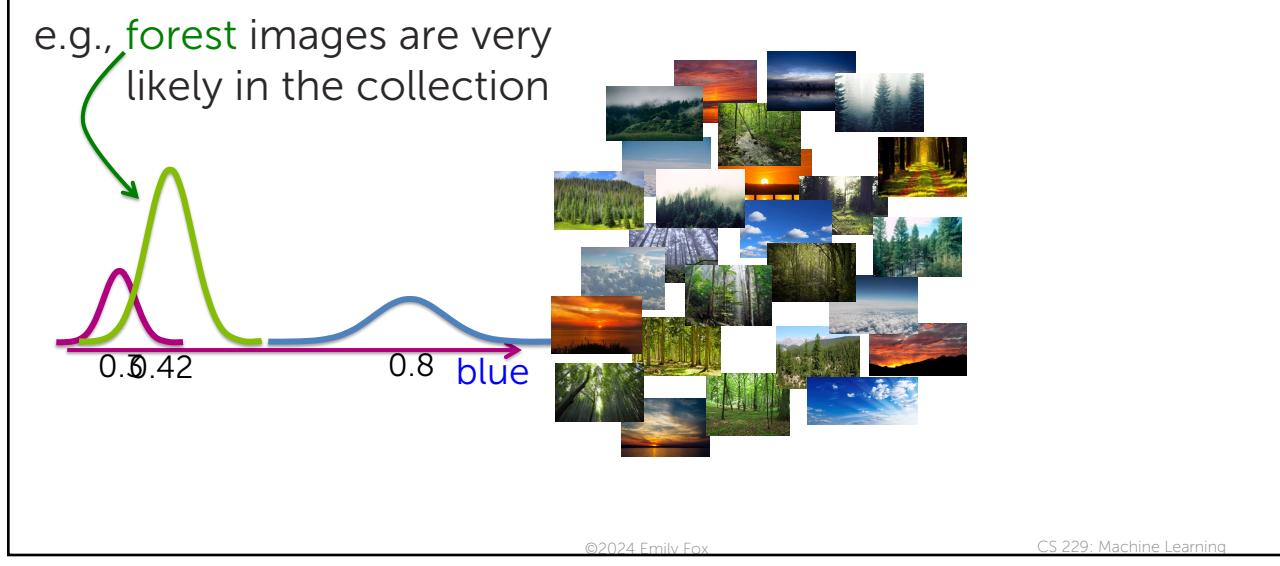
## Model of jumble of unlabeled images



59

## What if image types not equally represented?

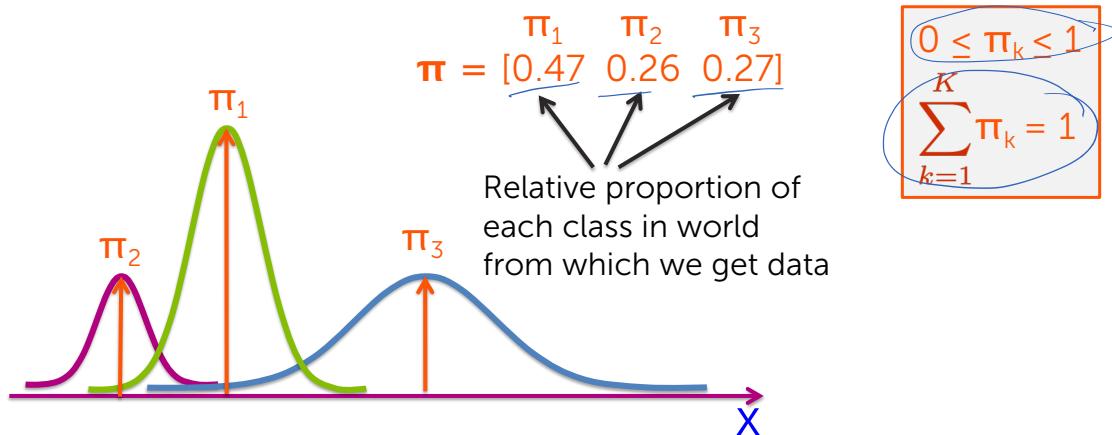
e.g., **forest** images are very likely in the collection



60

## Combination of weighted Gaussians

Associate a weight  $\pi_k$  with each Gaussian component



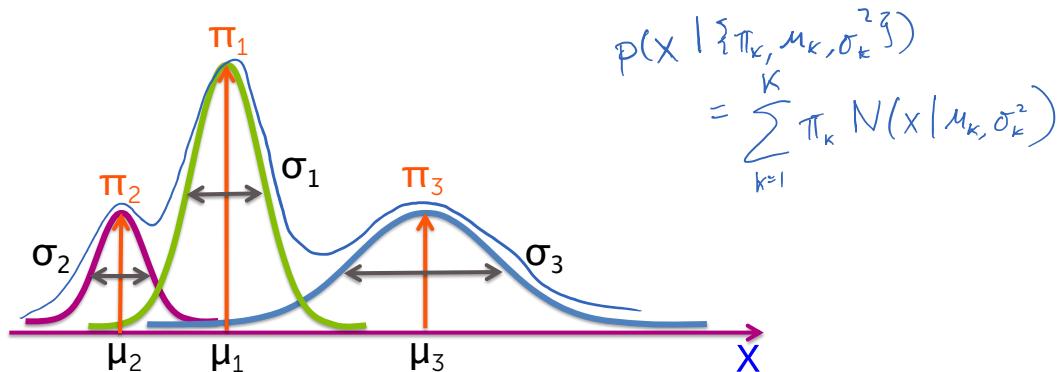
©2024 Emily Fox

CS 229: Machine Learning

61

## Mixture of Gaussians (1D)

Each mixture component represents a unique cluster specified by:  $\{\pi_k, \mu_k, \sigma_k^2\}$

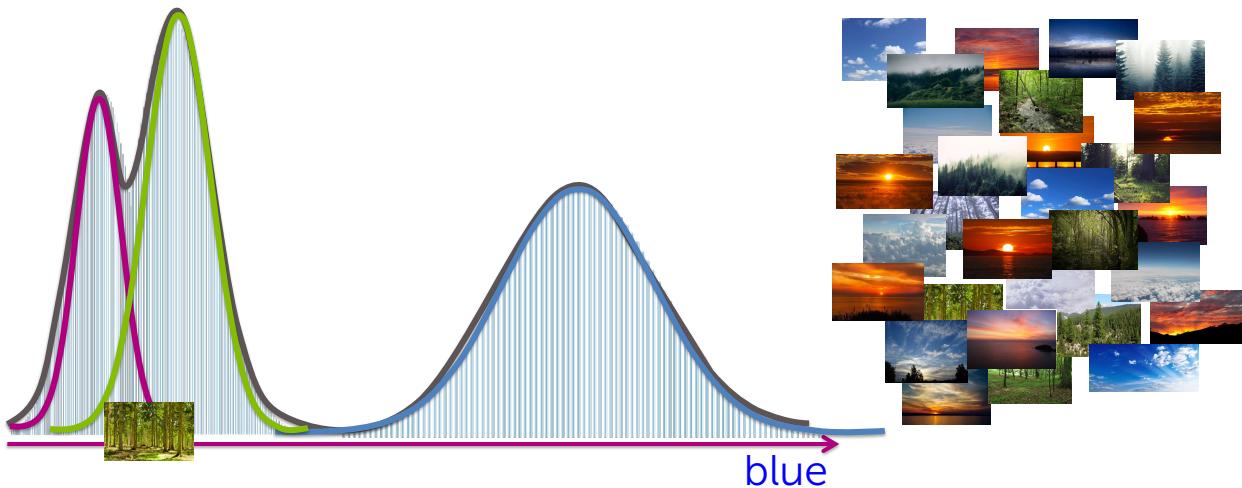


©2024 Emily Fox

CS 229: Machine Learning

62

## Distribution over **all** images



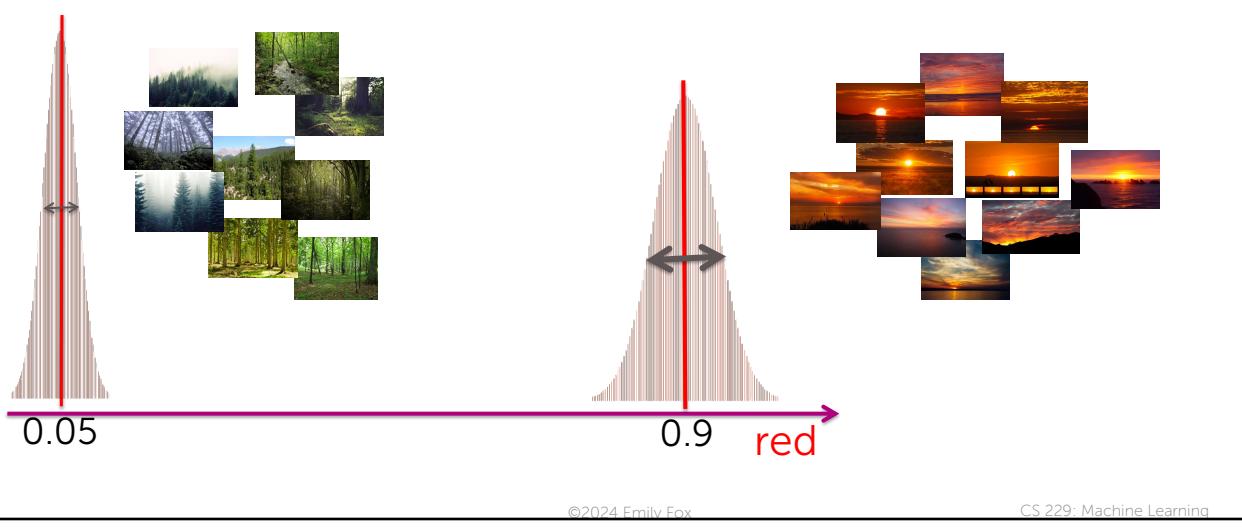
©2024 Emily Fox

CS 229: Machine Learning

63

## Can be distinguished along other dim

Now look at the **red** dimension

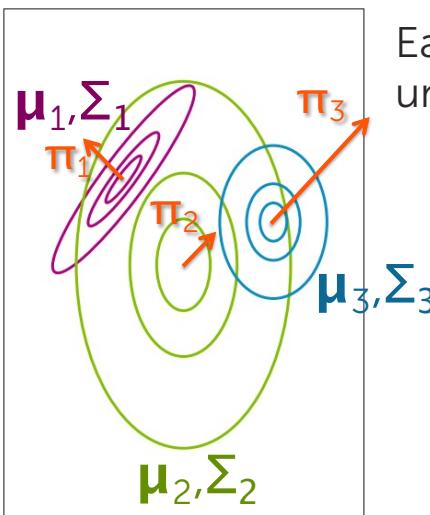


©2024 Emily Fox

CS 229: Machine Learning

64

## Mixture of Gaussians (general)



Each mixture component represents a unique cluster specified by:  $\{\pi_k, \mu_k, \Sigma_k\}$

$$\begin{aligned} p(x | \{\pi_k, \mu_k, \Sigma_k\}) \\ = \sum_{k=1}^K \pi_k N(x | \mu_k, \Sigma_k) \end{aligned}$$

©2024 Emily Fox

CS 229: Machine Learning

65

## According to the model...

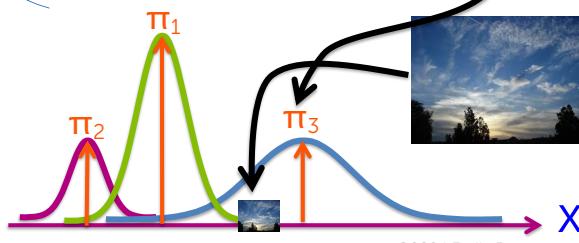
Without observing the image content, what's the probability it's from cluster k? (e.g., prob. of seeing "clouds" image)

$$p(z_i = k) = \pi_k$$

not cond.  $x_i$

Given observation  $x_i$  is from cluster k, what's the likelihood of seeing  $x_i$ ? (e.g., just look at distribution for "clouds")

$$p(x_i | z_i = k, \mu_k, \Sigma_k) = N(x_i | \mu_k, \Sigma_k)$$



©2024 Emily Fox

CS 229: Machine Learning

66

## Summary for mixture models

©2024 Emily Fox

CS 229: Machine Learning

67

## What you can do now...

- Interpret a probabilistic model-based approach to clustering using mixture models
- Describe model parameters
- Motivate the utility of soft assignments and describe what they represent
- Compare and contrast mixtures of Gaussians and k-means

©2024 Emily Fox

CS 229: Machine Learning

68

34

## Background: Gaussian distributions

OPTIONAL

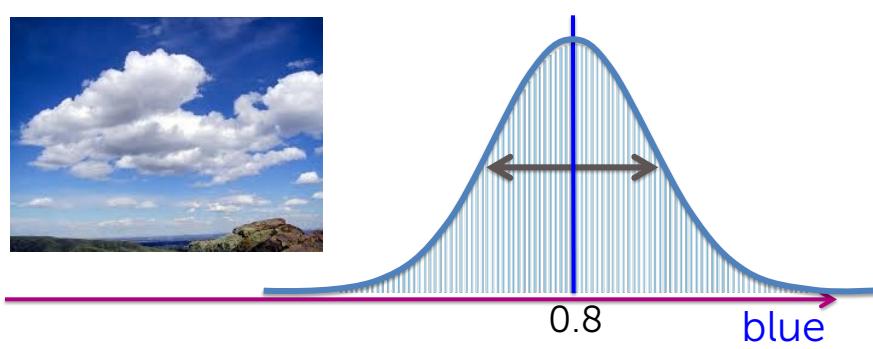
©2024 Emily Fox

CS 229: Machine Learning

69

## Model for a given image type

For **each dim** of the [R, G, B] vector, and **each image type**,  
assume a **Gaussian distribution** over color intensity



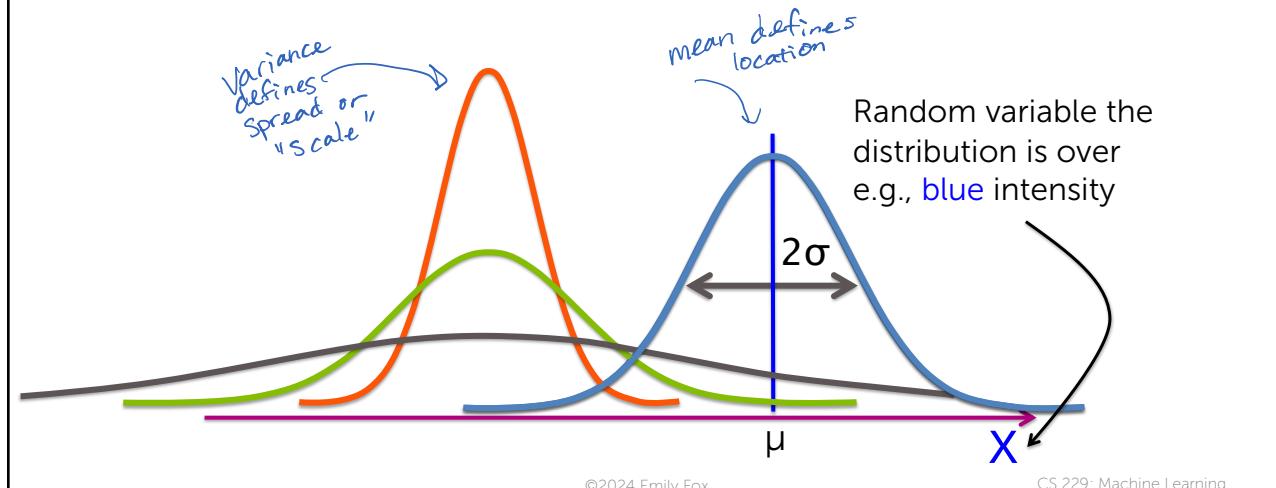
©2024 Emily Fox

CS 229: Machine Learning

70

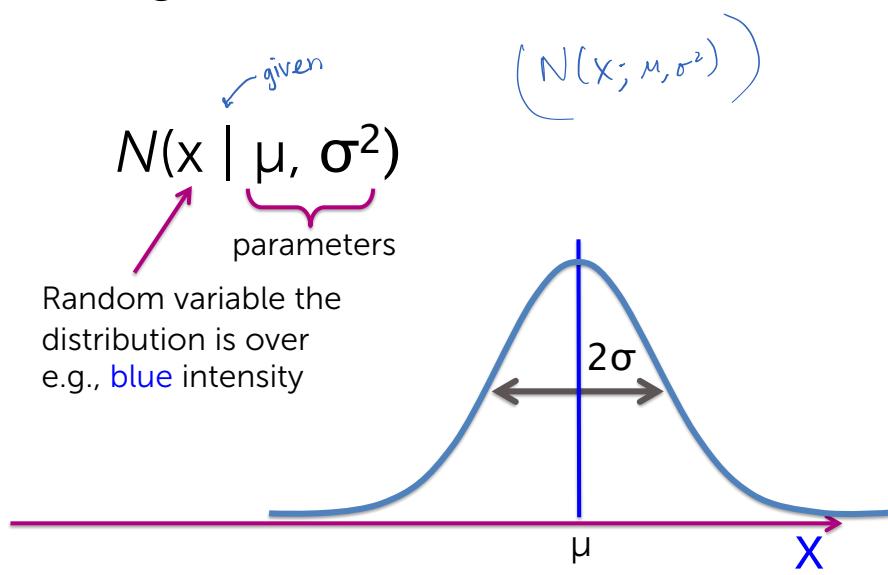
## 1D Gaussians

Fully specified by **mean**  $\mu$  and **variance**  $\sigma^2$  (or **st. dev.**  $\sigma$ )



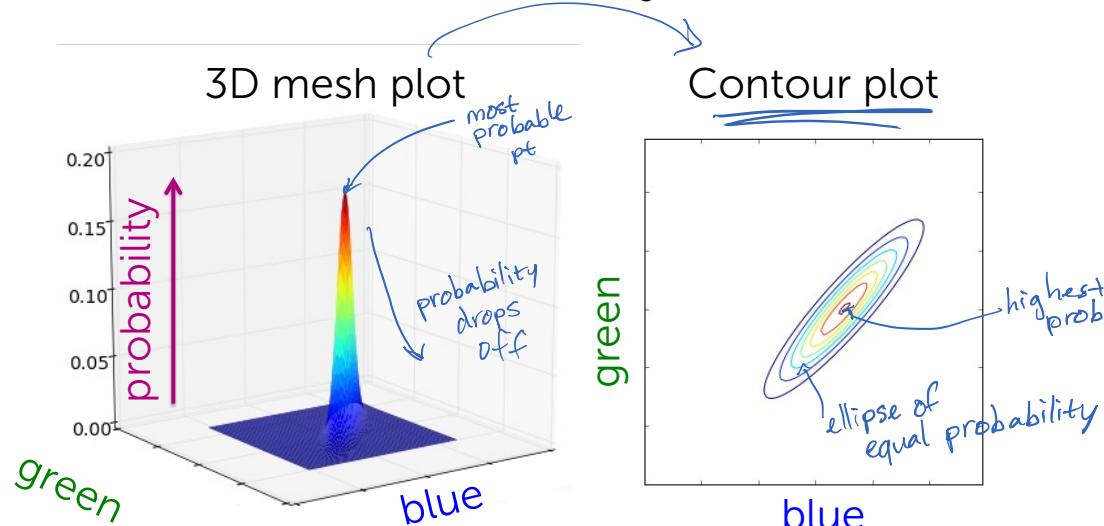
71

## Notating a 1D Gaussian distribution



72

## 2D Gaussians – Bird's eye view



©2024 Emily Fox

CS 229: Machine Learning

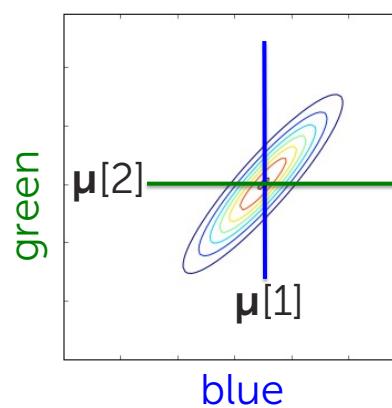
73

## 2D Gaussians – Parameters

Fully specified by **mean  $\mu$**  and **covariance  $\Sigma$**

$$\boldsymbol{\mu} = [\mu_{\text{blue}}, \mu_{\text{green}}]$$

mean centers the distribution in 2D



©2024 Emily Fox

CS 229: Machine Learning

74

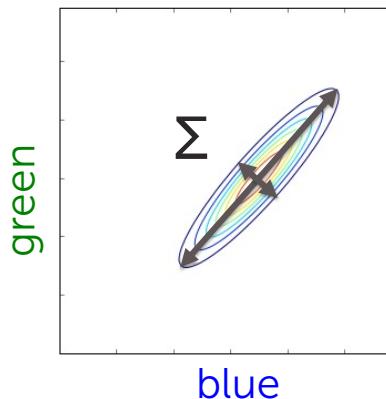
## 2D Gaussians – Parameters

Fully specified by **mean**  $\mu$  and **covariance**  $\Sigma$

$$\mu = [\mu_{\text{blue}}, \mu_{\text{green}}]$$

$$\Sigma = \begin{pmatrix} \sigma_{\text{blue}}^2 & \sigma_{\text{blue},\text{green}} \\ \sigma_{\text{green},\text{blue}} & \sigma_{\text{green}}^2 \end{pmatrix}$$

covariance determines orientation + spread



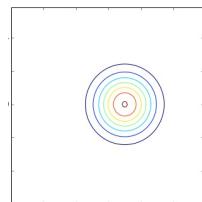
©2024 Emily Fox

CS 229: Machine Learning

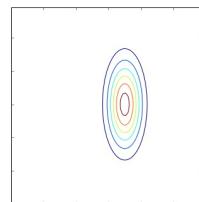
75

## Covariance structures

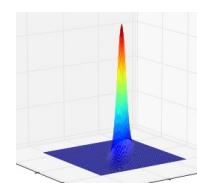
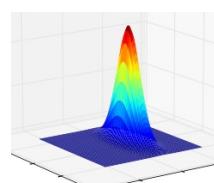
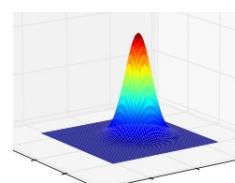
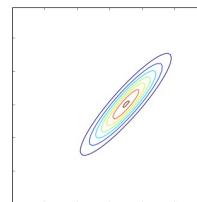
$$\Sigma = \begin{pmatrix} \sigma^2 & 0 \\ 0 & \sigma^2 \end{pmatrix}$$



$$\Sigma = \begin{pmatrix} \sigma_B^2 & 0 \\ 0 & \sigma_G^2 \end{pmatrix}$$



$$\Sigma = \begin{pmatrix} \sigma_B^2 & \sigma_{B,G} \\ \sigma_{G,B} & \sigma_G^2 \end{pmatrix}$$

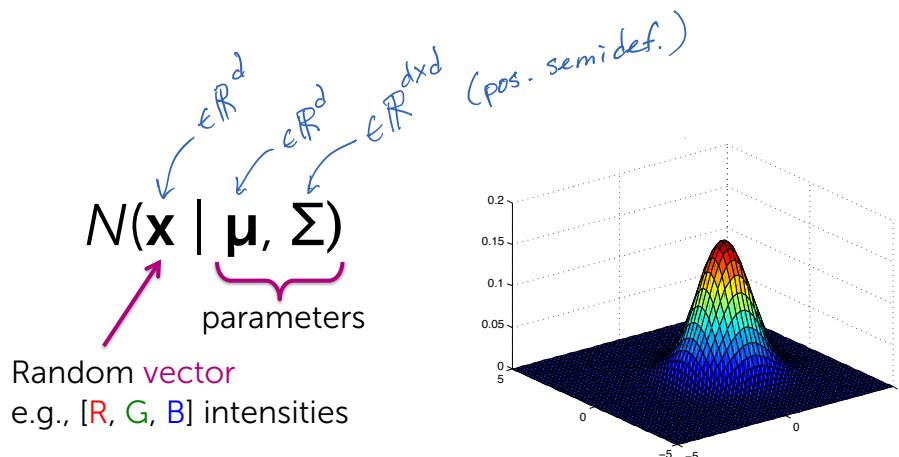


©2024 Emily Fox

CS 229: Machine Learning

76

## Notating a multivariate Gaussian



©2024 Emily Fox

CS 229: Machine Learning

77

Parameter estimation for 1D Gaussian distributions (and more background)

**OPTIONAL**

©2024 Emily Fox

CS 229: Machine Learning

78

## Some properties of Gaussians

- Affine transformation (multiplying by scalar and adding a constant)
  - $X \sim N(\mu, \sigma^2)$
  - $Y = aX + b \rightarrow Y \sim N(a\mu + b, a^2\sigma^2)$

same affine trans:  
 $E[aX+b] = aE[X]+b = a\mu+b$

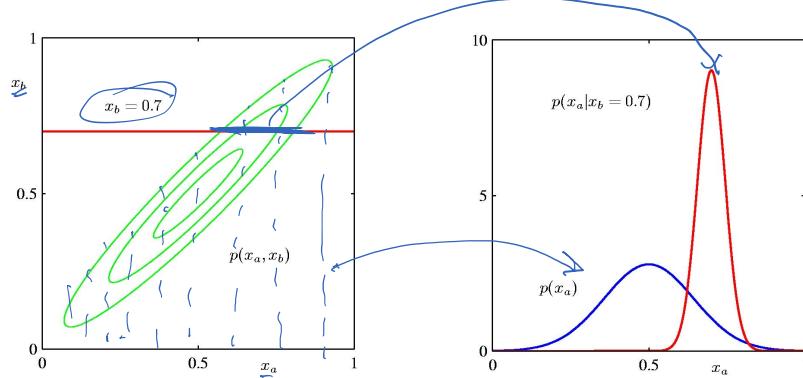
known (deterministic) quantities (scalars)
- Sum of independent Gaussian random variables
  - $X \sim N(\mu_X, \sigma^2_X)$
  - $Y \sim N(\mu_Y, \sigma^2_Y)$
  - $Z = X+Y \rightarrow Z \sim N(\mu_X+\mu_Y, \sigma^2_X+\sigma^2_Y)$

©2024 Emily Fox

CS 229: Machine Learning

79

## Conditional & marginal distributions



$$\begin{pmatrix} x_a \\ x_b \end{pmatrix} \sim N\left(\begin{pmatrix} \mu_a \\ \mu_b \end{pmatrix}, \begin{bmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{bmatrix}\right)$$

$$\text{Marg: } p(x_a) = \int p(x_a, x_b) dx_b \rightarrow x_a \sim N(\mu_a, \Sigma_{aa})$$

$$\text{Cond: } p(x_a | x_b) = \frac{p(x_a, x_b)}{\int p(x_a, x_b) dx_a} \rightarrow x_a | x_b \sim N\left(\mu_a + \Sigma_{ab} \Sigma_{bb}^{-1} (x_b - \mu_b), \Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba}\right)$$

©2024 Emily Fox

80

# Learning a 1D Gaussian

- Collect a bunch of data
  - Assume i.i.d. samples
- Learn parameters
  - Mean  $\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i$  ← Why?
  - Variance  $\hat{\sigma}^2$

$$p(x | \mu, \sigma) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

©2024 Emily Fox

CS 229: Machine Learning

81

## MLE for 1D Gaussian

- Prob. of i.i.d. samples  $D = \{x_1, \dots, x_N\}$ :

$$p(D | \mu, \sigma) = \left( \frac{1}{\sigma \sqrt{2\pi}} \right)^N \prod_{i=1}^N e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$

$$\hat{\mu}^{\text{MLE}}, \hat{\sigma}^{\text{MLE}} = \arg \max_{\mu, \sigma} p(D | \mu, \sigma) = \arg \max_{\mu, \sigma} \ln p(D | \mu, \sigma)$$

- Log-likelihood of data:

$$\begin{aligned} \ln p(D | \mu, \sigma) &= \ln \left[ \left( \frac{1}{\sigma \sqrt{2\pi}} \right)^N \prod_{i=1}^N e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} \right] \\ &= -N \ln \sigma \sqrt{2\pi} - \sum_{i=1}^N \frac{(x_i - \mu)^2}{2\sigma^2} \end{aligned}$$

©2024 Emily Fox

CS 229: Machine Learning

82

## MLE for mean of a 1D Gaussian

- What's MLE for the mean?

$$\frac{d}{d\mu} \ln p(D | \mu, \sigma) = \frac{d}{d\mu} \left[ -N \ln \sigma \sqrt{2\pi} - \sum_{i=1}^N \frac{(x_i - \mu)^2}{2\sigma^2} \right] = 0$$

*no depend.  
on  $\mu$*

$$\frac{d}{d\mu} \ln p(D | \mu, \sigma) = - \sum_{i=1}^N \frac{d}{d\mu} \frac{(x_i - \mu)^2}{2\sigma^2} = \sum_{i=1}^N \frac{x_i - \mu}{\sigma^2} = 0$$

$$N \hat{\mu}^{\text{MLE}} = \sum_{i=1}^N x_i$$

$$\hat{\mu}^{\text{MLE}} = \frac{1}{N} \sum_{i=1}^N x_i \quad \leftarrow \text{MLE doesn't depend  
on choice of } \sigma^2$$

©2024 Emily Fox

CS 229: Machine Learning

83

## MLE for variance

- Again, set derivative to zero:

$$\frac{d}{d\sigma} \ln p(D | \mu, \sigma) = \frac{d}{d\sigma} \left[ -N \ln \sigma \sqrt{2\pi} - \sum_{i=1}^N \frac{(x_i - \mu)^2}{2\sigma^2} \right]$$

*$\equiv$*

$$= \frac{d}{d\sigma} \left[ -N \ln \sigma \sqrt{2\pi} \right] - \sum_{i=1}^N \frac{d}{d\sigma} \left[ \frac{(x_i - \mu)^2}{2\sigma^2} \right] = 0$$

$$\Rightarrow -\frac{N}{\sigma} + \sum_{i=1}^N \frac{(x_i - \mu)^2}{\sigma^3} = 0$$

$$\Rightarrow \hat{\sigma}^2 \text{MLE} = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu}^{\text{MLE}})^2$$

use  $\mu = \hat{\mu}^{\text{MLE}}$  bc  
optimal choice of  
 $\mu$  doesn't depend  
on  $\sigma$

©2024 Emily Fox

CS 229: Machine Learning

84

## Learning 1D Gaussian parameters

- MLE:
$$\hat{\mu}_{MLE} = \frac{1}{N} \sum_{i=1}^N x_i$$
$$\hat{\sigma}_{MLE}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu}_{MLE})^2$$
- FYI, MLE for the variance of a Gaussian is **biased**
  - Expected value of estimator is **not** true parameter!
  - Unbiased variance estimator:

$$\hat{\sigma}_{unbiased}^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \hat{\mu}_{MLE})^2$$

©2024 Emily Fox

CS 229: Machine Learning