

# **STAT 214 Spring 2025**

## **Week 7**

Austin Zane

# Outline

- GMM and EM
- Submitting jobs to Bridges2
- Train/validation/test splits

# **GMM & EM**

# Various resources

Intuition:

- Quick:

<https://stackoverflow.com/questions/11808074/what-is-an-intuitive-explanation-of-the-expectation-maximization-technique#answer-43561339>

Math:

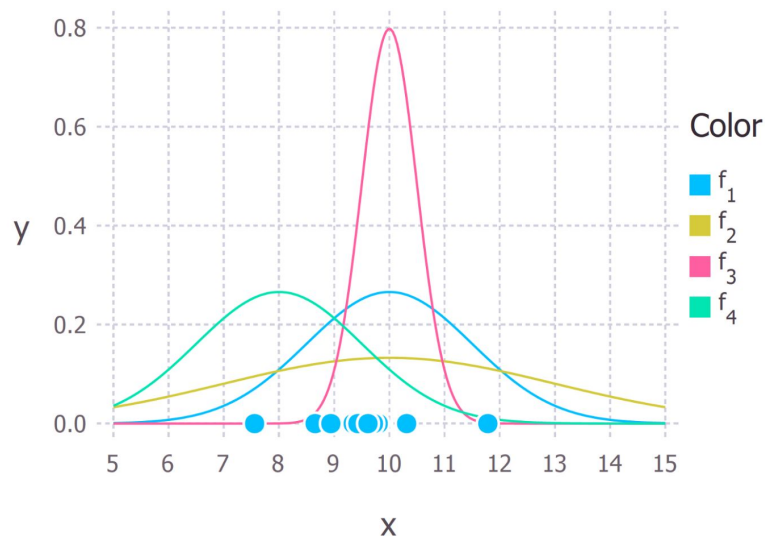
- Quick overview:

[http://www.seanborman.com/publications/EM\\_algorithm.pdf](http://www.seanborman.com/publications/EM_algorithm.pdf)

- More in-depth: <https://arxiv.org/pdf/1105.1476.pdf>

# Maximum Likelihood Estimation Review

- Data:  $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} p_\theta(x)$
- Likelihood:  $\mathcal{L}(\theta) = \prod_{i=1}^n p_\theta(x_i)$
- Log-likelihood:  $\ell(\theta) = \log \mathcal{L}(\theta)$
- MLE:  $\hat{\theta} = \arg \max_{\theta} \mathcal{L}(\theta)$   
$$= \arg \max_{\theta} \ell(\theta)$$



- **Intuition:** Find the value of  $\theta$  under which we would be least surprised to see a sample like the observed one.
- **Optimization problem:** take derivative, set equal to zero, solve for parameter.

# EM Overview

## Motivation 1: “Hard” Maximum Likelihood Estimation Problems

Say we have the following mixture of Gaussians problem:

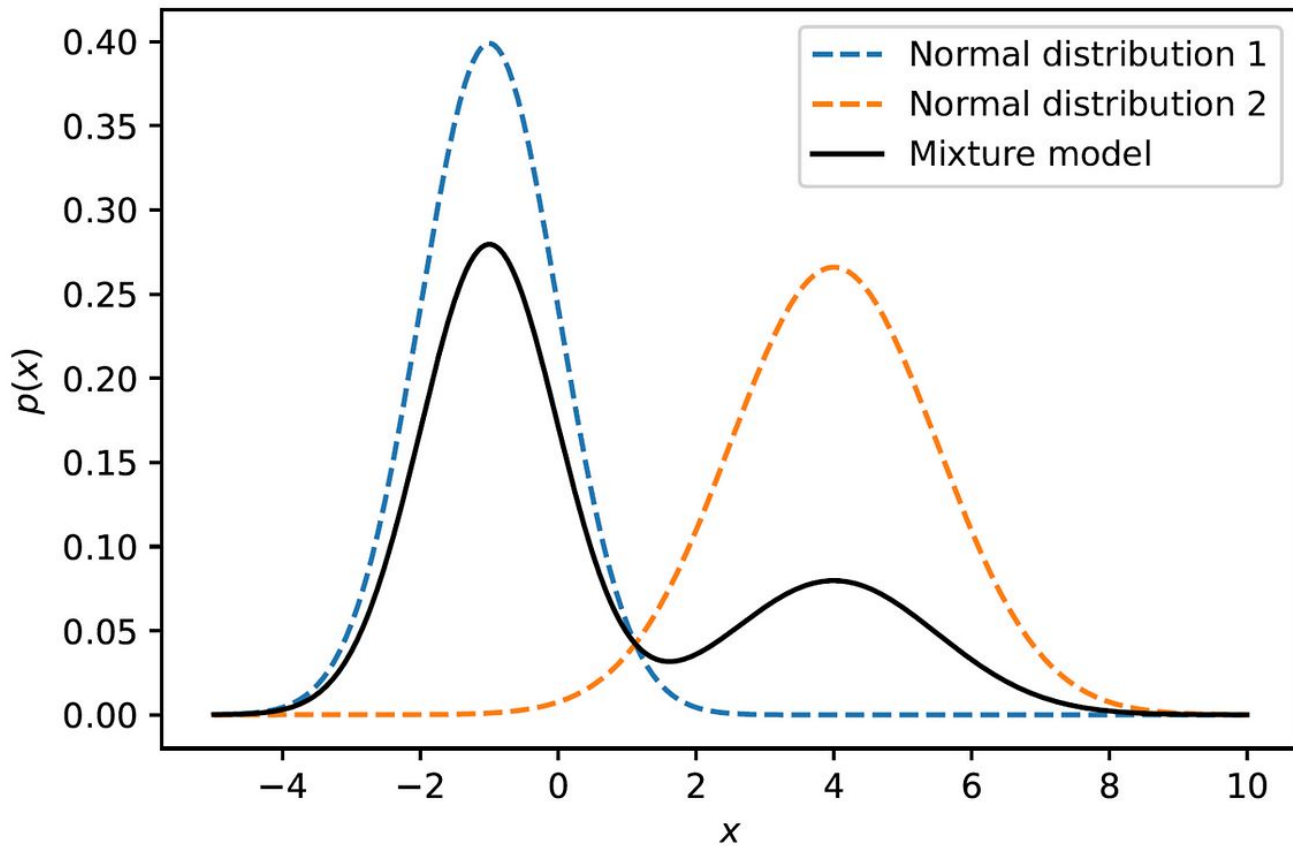
$$X_1, \dots, X_n \stackrel{i.i.d.}{\sim} \pi_1 N(\mu_1, \sigma_1^2) + \pi_2 N(\mu_2, \sigma_2^2), \pi_1 + \pi_2 = 1$$

- Because of the sum of normals, the log-likelihood isn't as helpful as before and taking derivatives w.r.t  $\theta = (\pi_1, \mu_1, \sigma_1, \mu_2, \sigma_2)$  and setting equal to zero, etc., doesn't lead to closed form solutions.
- Instead, we can introduce **latent variables** which tell us which of the two Gaussians each observation comes from:  $Z_i \stackrel{i.i.d.}{\sim} \text{Bernoulli}(\pi_1)$
- If we know the latent allocations then the problem simply becomes two easy MLE exercises.

$$X_i | Z_i = 1 \stackrel{i.i.d.}{\sim} N(\mu_1, \sigma_1^2)$$

$$X_i | Z_i = 2 \stackrel{i.i.d.}{\sim} N(\mu_2, \sigma_2^2)$$

# EM Overview



# EM Overview

## Motivation 2: Clustering

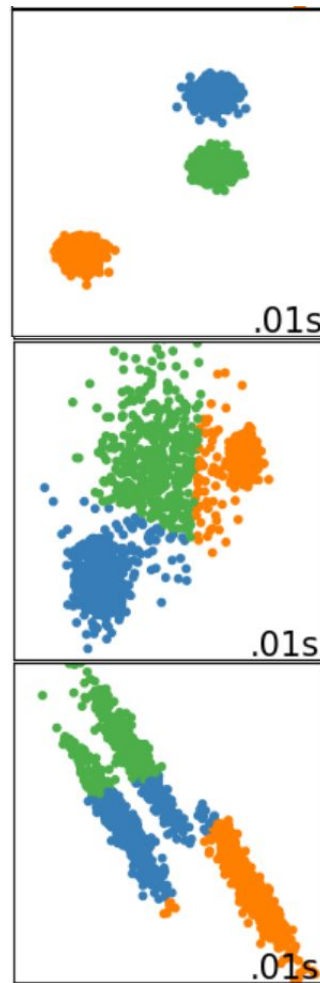
Since EM helps with finding solutions to mixture problems, it lends itself naturally to clustering.

### Recall:

- K-means performs well when clusters are homogeneous.
- But it fails miserably when faced with elongated or irregular shapes.

The EM generalizes K-means:

- Still performs great where K-means does.





# EM Overview

## Motivation 2: Clustering

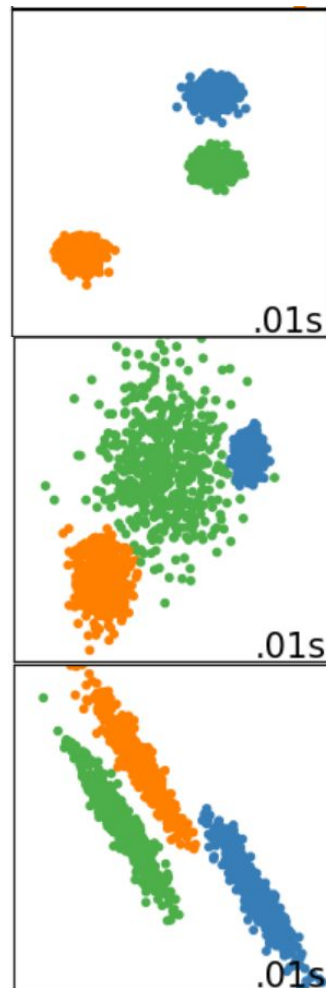
Since EM helps with finding solutions to mixture problems, it lends itself naturally to clustering.

### Recall:

- K-means performs well when clusters are homogeneous.
- But it fails miserably when faced with elongated or irregular shapes.

The EM generalizes K-means:

- Still performs great where K-means does.
- But it can also handle differences in spread and symmetry.
- EM is a “soft” clustering algorithm.



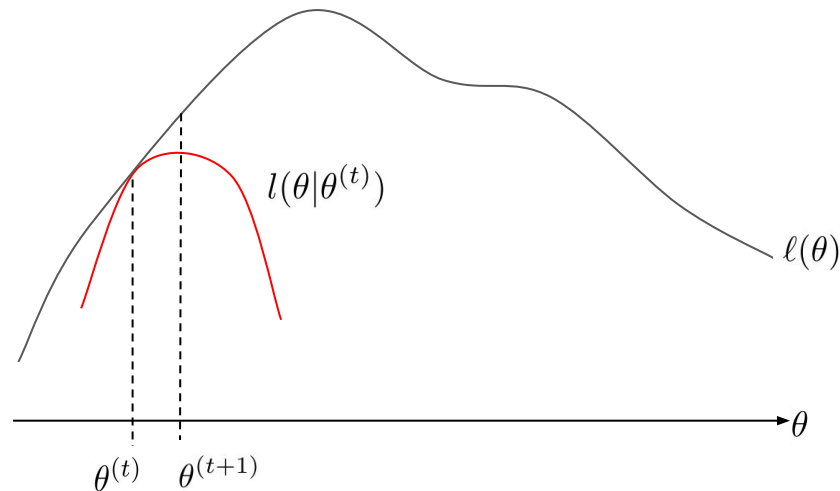
# EM algorithm intuition

Say we make a guess  $\theta^{(t)}$

The insight of the EM algorithm is that we can find a function  $l(\theta|\theta^{(t)})$  such that

- $l(\theta|\theta^{(t)}) \leq \ell(\theta)$
- $l(\theta^{(t)}|\theta^{(t)}) = \ell(\theta^{(t)})$

So any  $\theta$  that increases  $l(\theta|\theta^{(t)})$  also increases  $\ell(\theta)$



# EM algorithm steps

It turns out that maximizing  $l(\theta|\theta^{(t)})$  is equivalent\* to maximizing the expectation  $Q(\theta|\theta^{(t)}) = \mathbb{E}[\ell(\theta; X, Z)|\theta^{(t)}, X]$

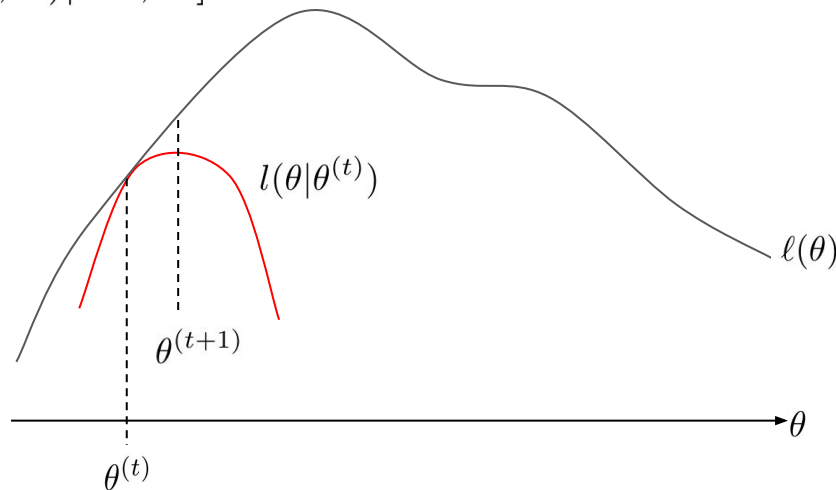
The EM algorithm involves two steps that are repeated until convergence:

1. **E:** Calculate the expectation

$$Q(\theta|\theta^{(t)}) = \mathbb{E}[\ell(\theta; X, Z)|\theta^{(t)}, X]$$

1. **M:** Maximize  $Q(\theta|\theta^{(t)})$  w.r.t.  $\theta$

Note: we can initialize with a random guess  $\theta^{(0)}$



\* see the Borman tutorial for the derivation

# EM: Gaussian Mixture Example

Same setup from before:

$$X_1, \dots, X_n \stackrel{i.i.d.}{\sim} \pi_1 N(\mu_1, \sigma_1^2) + \pi_2 N(\mu_2, \sigma_2^2), \pi_1 + \pi_2 = 1$$

$$X_i | Z_i = 1 \stackrel{i.i.d.}{\sim} N(\mu_1, \sigma_1^2) \quad X_i | Z_i = 2 \stackrel{i.i.d.}{\sim} N(\mu_2, \sigma_2^2)$$

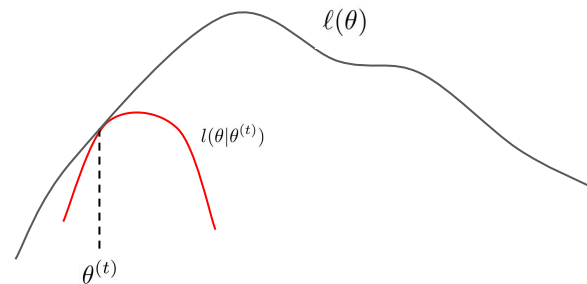
$$Z_i \stackrel{i.i.d.}{\sim} \text{Bernoulli}(\pi_1)$$

**Likelihood:**  $p_\theta(x_i, z_i) = p_\theta(x_i | z_i) p_\theta(z_i) = \begin{cases} \frac{1}{\sqrt{2\pi\sigma_1^2}} \exp\left\{-\frac{(x_i - \mu_1)^2}{2\sigma_1^2}\right\} \pi_1, & \text{if } Z_i = 1 \\ \frac{1}{\sqrt{2\pi\sigma_2^2}} \exp\left\{-\frac{(x_i - \mu_2)^2}{2\sigma_2^2}\right\} \pi_2, & \text{if } Z_i = 2 \end{cases}$

$$\log p_\theta(x_i, z_i) = \begin{cases} -\frac{1}{2} \log 2\pi - \log \sigma_1 - \frac{(x_i - \mu_1)^2}{2\sigma_1^2} + \log \pi_1, & \text{if } Z_i = 1 \\ -\frac{1}{2} \log 2\pi - \log \sigma_2 - \frac{(x_i - \mu_2)^2}{2\sigma_2^2} + \log \pi_2, & \text{if } Z_i = 2 \end{cases}$$

# E-Step: Compute $Q(\theta|\theta^{(t)}) = \mathbb{E}[\ell(\theta; X, Z)|\theta^{(t)}, X]$

$$\begin{aligned}
 Q(\theta|\theta^{(t)}) &= \sum_{i=1}^n \mathbb{E}[\ell(\theta; X_i, Z_i)|\theta^{(t)}, X] \\
 &= \sum_{i=1}^n \left\{ \mathbb{E}[\ell(\theta; X_i, Z_i)|\theta^{(t)}, X, Z_i = 1]\mathbb{P}(Z_i = 1|\theta^{(t)}, X) + \right. \\
 &\quad \left. \mathbb{E}[\ell(\theta; X_i, Z_i)|\theta^{(t)}, X, Z_i = 2]\mathbb{P}(Z_i = 2|\theta^{(t)}, X) \right\} \\
 &\quad \text{(law of total expectation)}
 \end{aligned}$$



$$\begin{aligned}
 &= \sum_{i=1}^n \left\{ \log \pi_1 - \frac{1}{2} \log 2\pi - \log \sigma_1 - \frac{(x_i - \mu_1)^2}{2\sigma_1^2} \right\} Z_{i,1}^{(t)} + \\
 &\quad \left( \log \pi_2 - \frac{1}{2} \log 2\pi - \log \sigma_2 - \frac{(x_i - \mu_2)^2}{2\sigma_2^2} \right) Z_{i,2}^{(t)} \Big\}
 \end{aligned}$$

Use Bayes rule

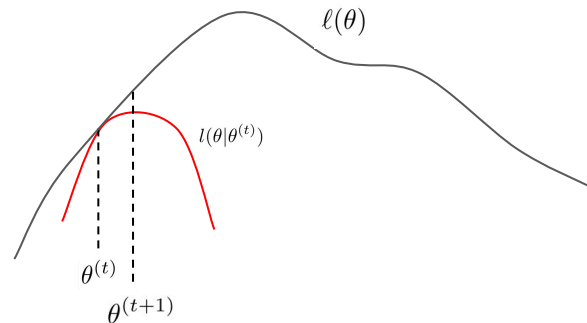
$$\begin{aligned}
 Z_{i,j}^{(t)} &= \mathbb{P}(Z_i = j|\theta^{(t)}, X) \\
 &= \frac{\pi_j^{(t)} \phi\left(\frac{x_i - \mu_j}{\sigma_j}\right)/\sigma_j}{\sum_{k=1}^2 \pi_k^{(t)} \phi\left(\frac{x_i - \mu_k}{\sigma_k}\right)/\sigma_k}
 \end{aligned}$$

Standard normal pdf

# M-Step: Maximize $Q(\theta|\theta^{(t)})$ w.r.t. $\theta = (\pi_1, \mu_1, \sigma_1, \mu_2, \sigma_2)$

- $\pi_1 : \frac{\partial Q}{\partial \pi_1} = \frac{\partial}{\partial \pi_1} \sum_{i=1}^n \left( \log \pi_1 Z_{i,1}^{(t)} + \log(1 - \pi_1) Z_{i,2}^{(t)} \right)$   
 $= \sum_{i=1}^n \frac{Z_{i,1}^{(t)}}{\pi_1} + \sum_{i=1}^n \frac{Z_{i,2}^{(t)}}{1 - \pi_1}$

$$\stackrel{\text{set}}{=} 0 \implies \pi_1^{(t+1)} = \frac{\sum_i Z_{i,1}^{(t)}}{\sum_i Z_{i,1}^{(t)} + Z_{i,2}^{(t)}} = \frac{1}{n} \sum_i Z_{i,1}^{(t)}$$

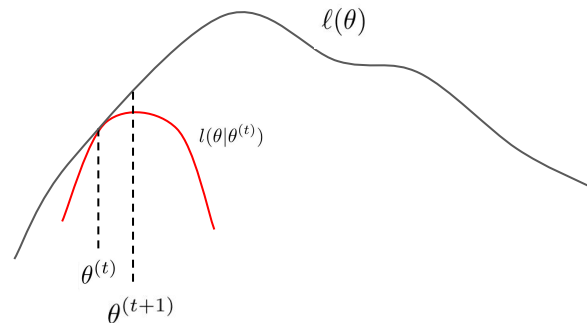


- $\mu_1 : \frac{\partial Q}{\partial \mu_1} = \frac{\partial}{\partial \mu_1} \sum_{i=1}^n \left( -\frac{(x_i - \mu_1)^2}{2\sigma_1^2} \right) Z_{i,1}^{(t)} \implies \mu_1^{(t+1)} = \frac{\sum_i Z_{i,1}^{(t)} X_i}{\sum_i Z_{i,1}^{(t)}}$   
 $\stackrel{\text{set}}{=} 0$

**M-Step: Maximize**  $Q(\theta|\theta^{(t)})$  **w.r.t.**  $\theta = (\pi_1, \mu_1, \sigma_1, \mu_2, \sigma_2)$

- $$\pi_1^{(t+1)} = \frac{\sum_i Z_{i,1}^{(t)}}{\sum_i Z_{i,1}^{(t)} + Z_{i,2}^{(t)}} = \frac{1}{n} \sum_i Z_{i,1}^{(t)}$$

- $$\mu_1^{(t+1)} = \frac{\sum_i Z_{i,1}^{(t)} X_i}{\sum_i Z_{i,1}^{(t)}}$$



- $$\sigma_1: \frac{\partial Q}{\partial \sigma_1} = \frac{\partial}{\partial \sigma_1} \sum_{i=1}^n \left( -\log \sigma_1 - \frac{(x_i - \mu_1)^2}{2\sigma_1^2} \right) Z_{i,1}^{(t)}$$

$$= \sum_{i=1}^n \left( -\frac{1}{\sigma_1} + \frac{(x_i - \mu_1)^2}{\sigma_1^3} \right) Z_{i,1}^{(t)}$$

$\stackrel{\text{set}}{=} 0$

$$\Rightarrow (\sigma_1^2)^{(t+1)} = \frac{\sum_i Z_{i,1}^{(t)} (X_i - \mu_1^{(t+1)})^2}{\sum_i Z_{i,1}^{(t)}}$$

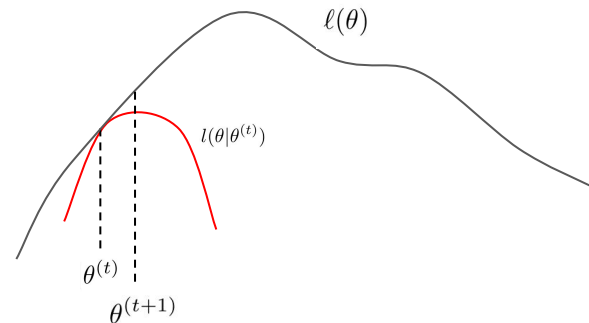
**M-Step: Maximize  $Q(\theta|\theta^{(t)})$  w.r.t.  $\theta = (\pi_1, \mu_1, \sigma_1, \mu_2, \sigma_2)$**

- $$\pi_1^{(t+1)} = \frac{\sum_i Z_{i,1}^{(t)}}{\sum_i Z_{i,1}^{(t)} + Z_{i,2}^{(t)}} = \frac{1}{n} \sum_i Z_{i,1}^{(t)}$$

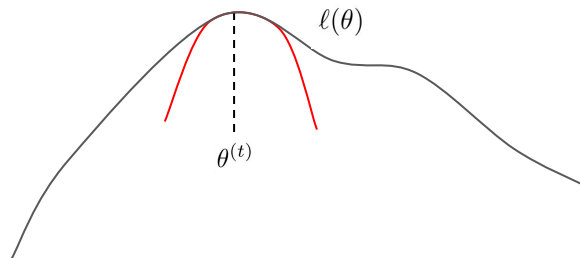
- $$\mu_1^{(t+1)} = \frac{\sum_i Z_{i,1}^{(t)} X_i}{\sum_i Z_{i,1}^{(t)}}$$

- $$(\sigma_1^2)^{(t+1)} = \frac{\sum_i Z_{i,1}^{(t)} (X_i - \mu_1^{(t+1)})^2}{\sum_i Z_{i,1}^{(t)}}$$

- Similar process for  $\mu_2$  &  $\sigma_2$



Repeat this process until we find a (local) maximum





**stat-214-gsi/discussion/week7/em.ipynb**

# **SLURM on Bridges2**

# Write a shell script to run your job

```
#!/bin/bash
#SBATCH -N 1
#SBATCH -p GPU-shared
#SBATCH -t 5:00:00
#SBATCH --gpus=v100-32:1

#type 'man sbatch' for more information and options
#this job will ask for 1 V100 GPU on a v100-32 node in GPU-shared for 5 hours
#this job would potentially charge 5 GPU SUS

#echo commands to stdout
set -x

# move to working directory
# this job assumes:
# - all input data is stored in this directory
# - all output should be stored in this directory
# - please note that groupname should be replaced by your groupname
# - PSC-username should be replaced by your PSC username
# - path-to-directory should be replaced by the path to your directory where the executable is
module load anaconda3
conda activate env_214
cd /ocean/projects/groupname/PSC-username/path-to-directory

#run Python script which is already in your project space

python train.py
```

# Submit your job

To submit:

```
sbatch shell_example.sh
```

To cancel:

```
scancel 12345 (replacing 12345 with the id of your job)
```

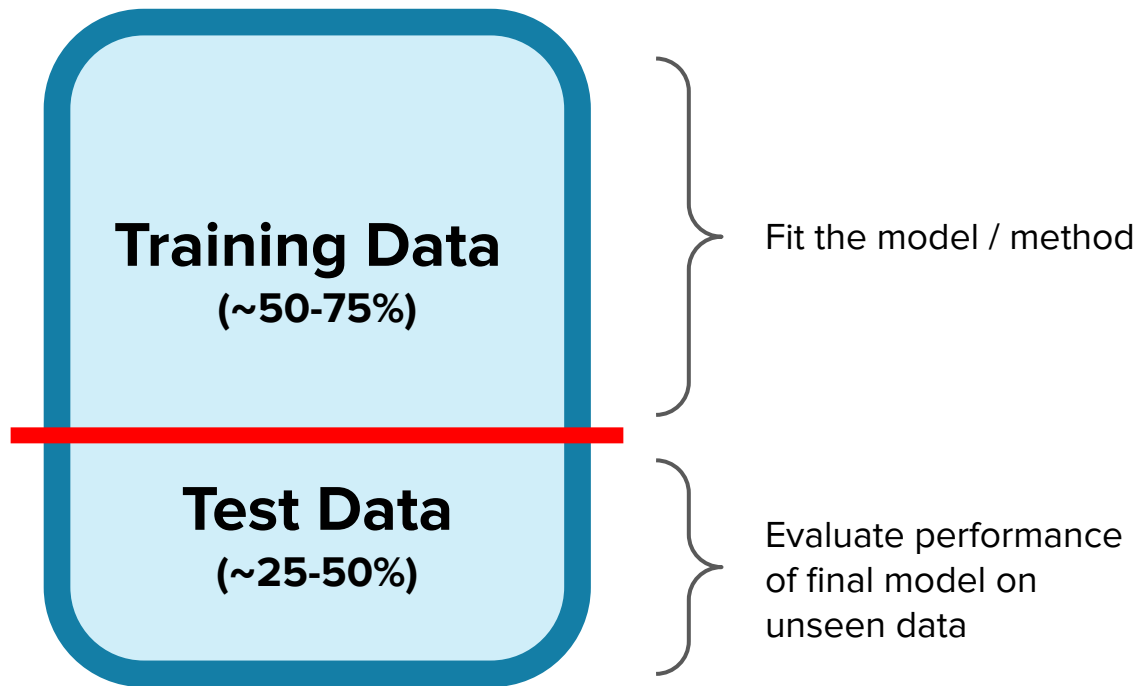
To see my running jobs:

```
squeue -u austin.zane
```

# **Splitting your data**

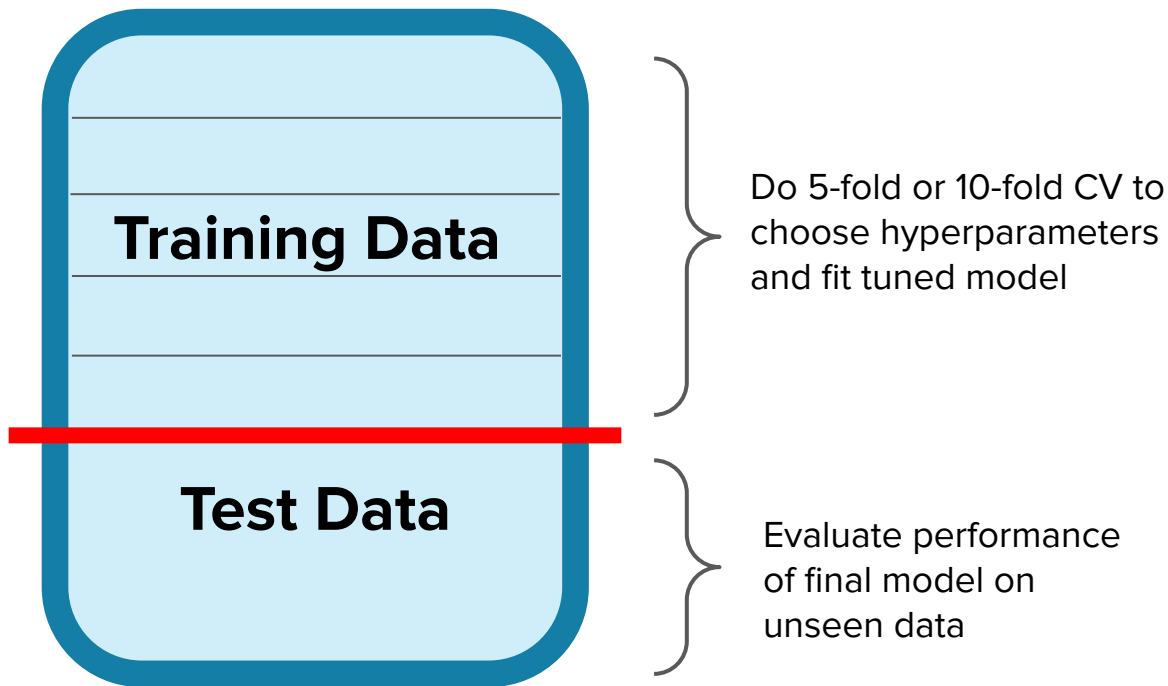
# Data splitting: a way to assess generalizability

Simplest case:

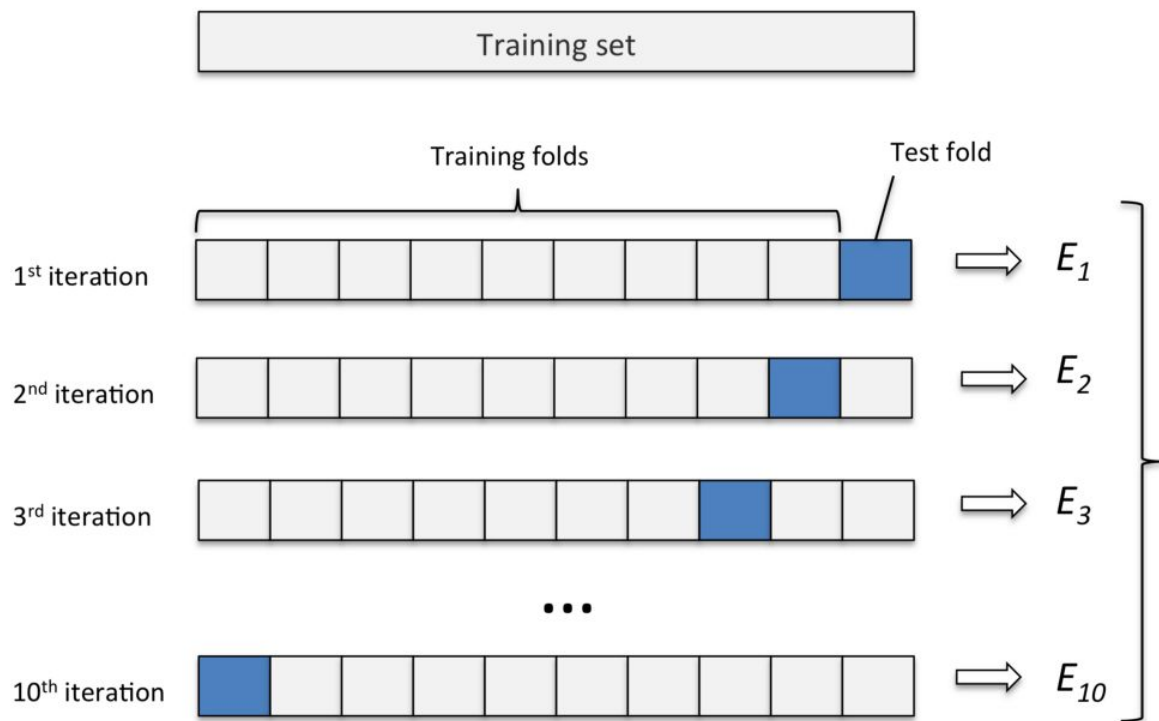


# Data splitting: a way to assess generalizability

With cross validation:

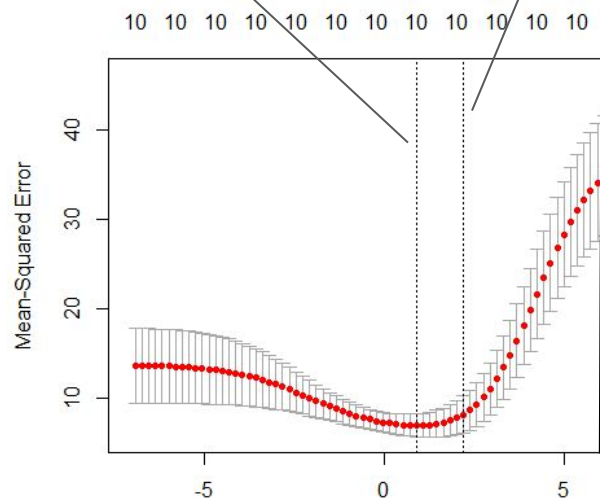


# K-fold cross validation for choosing hyperparameters



$\lambda$  with smallest mean CV error

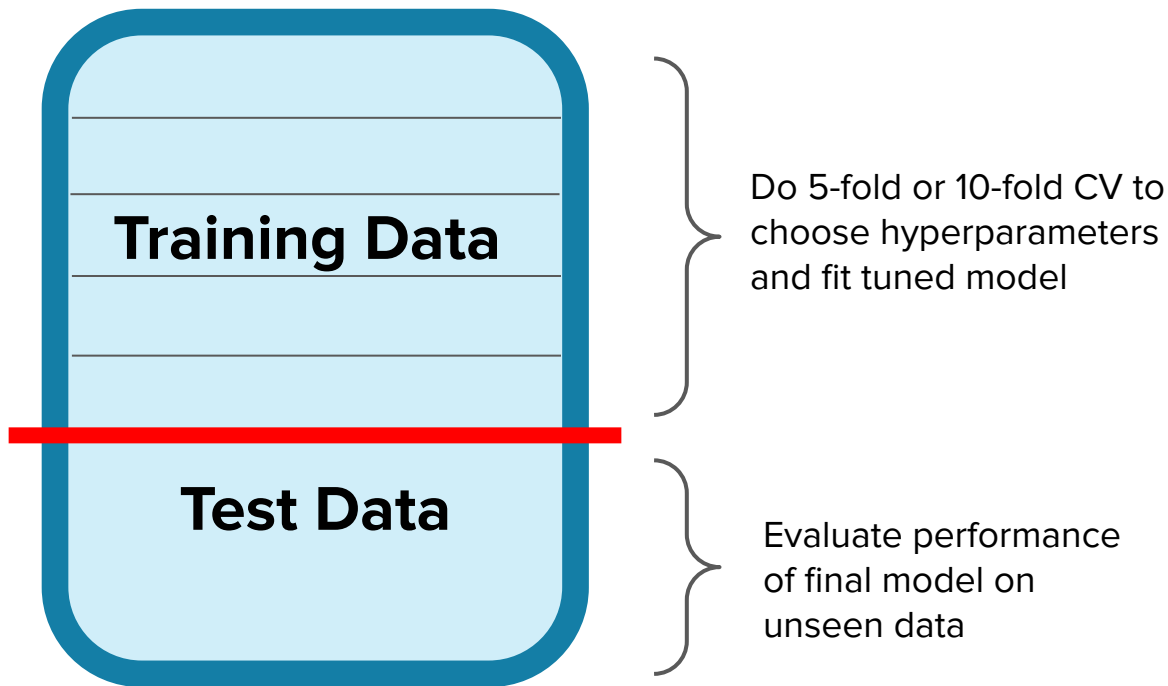
Optimal  $\lambda$  using the 1 SE rule



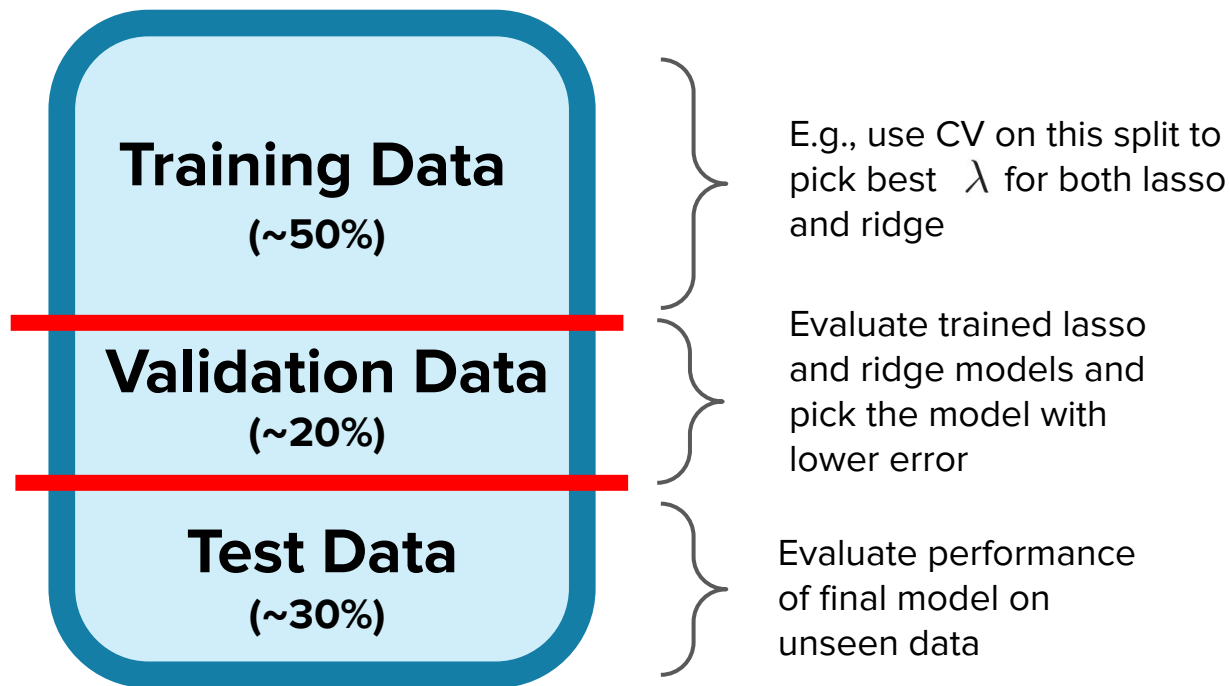


# Data splitting: a way to assess generalizability

With cross validation:



# Data splitting + tuning hyperparameters + multiple methods



# How much data to hold out for testing?

- Some say “30%” or give an arbitrary number.
- Realistically it depends on the problem.
  - If you have a billion observations, you might only need a several thousand in the test set to get a good idea of accuracy (assuming they are a simple random sample)
- Think about how low you need the standard error of your accuracy estimate to be.

