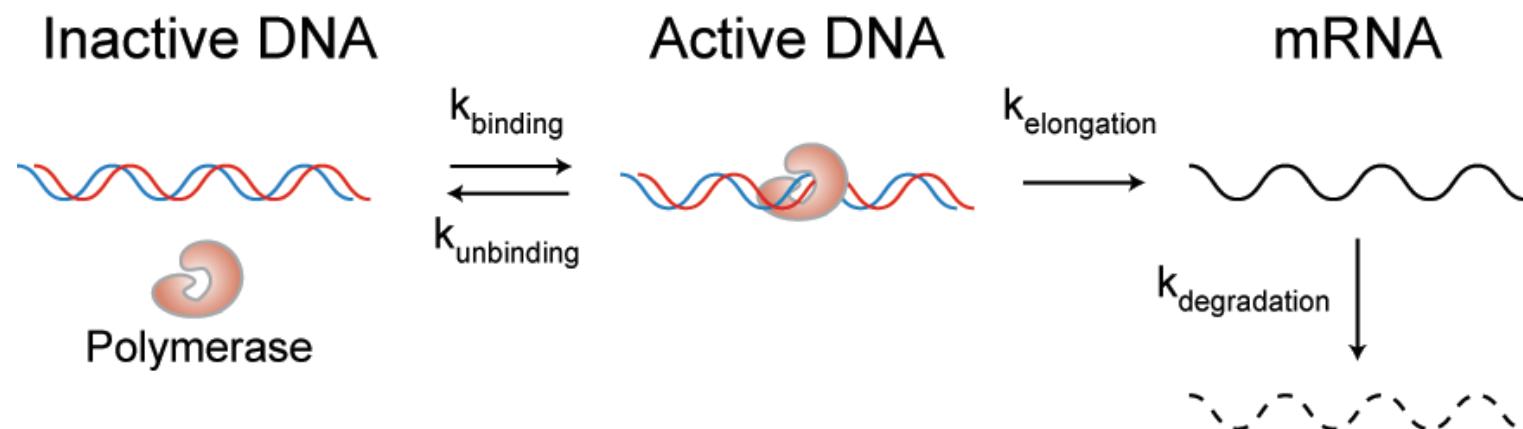


Parameterizing heterogeneity

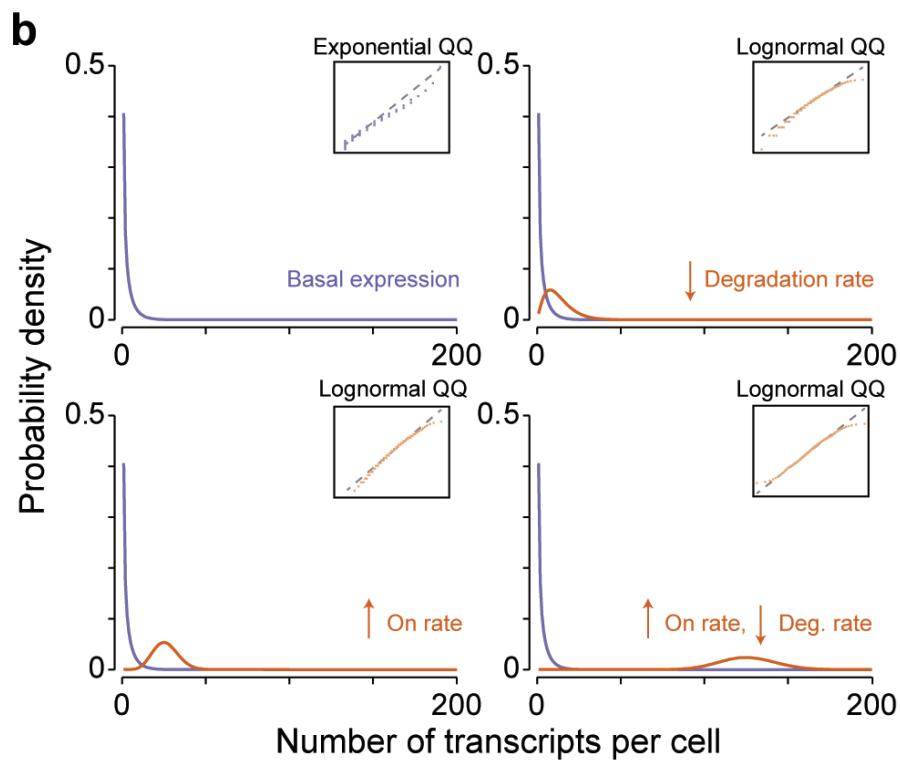
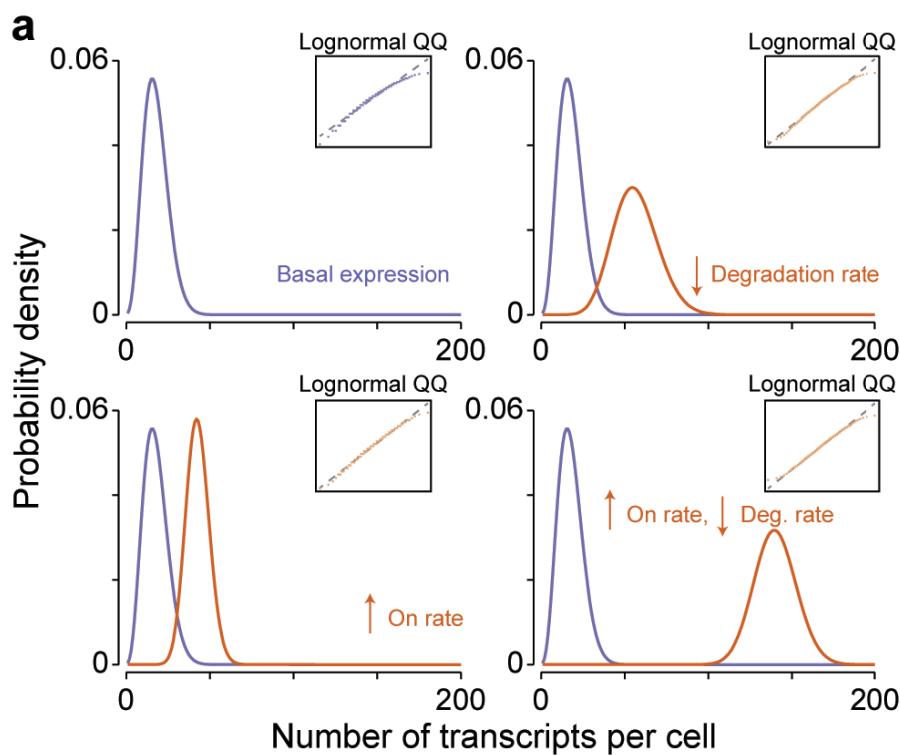
Limitations of stochastic profiles

- Still not single-cell resolution transcriptomics
- Only make a binary decision as to whether a gene is heterogeneously regulated, but no information of the extent of heterogeneity
- Solution: use mathematical modeling to determine the single-cell expression distribution

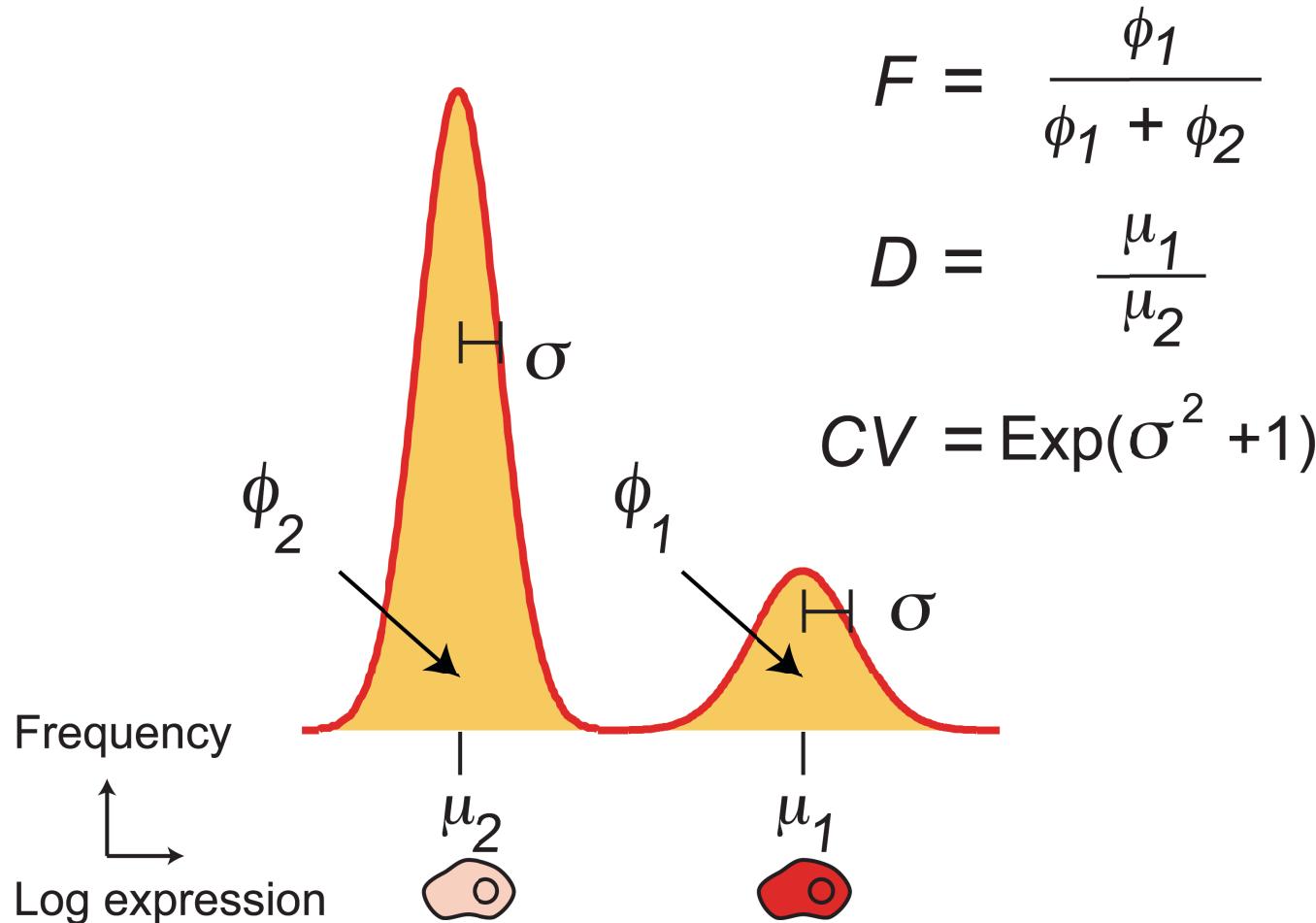
Single-cell gene expression and regulation



Probability models of gene regulation

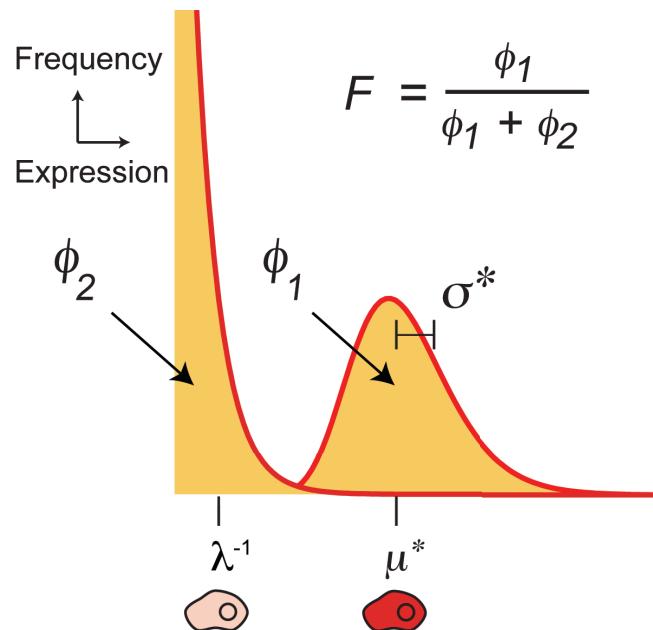


Dichotomous gene expression model



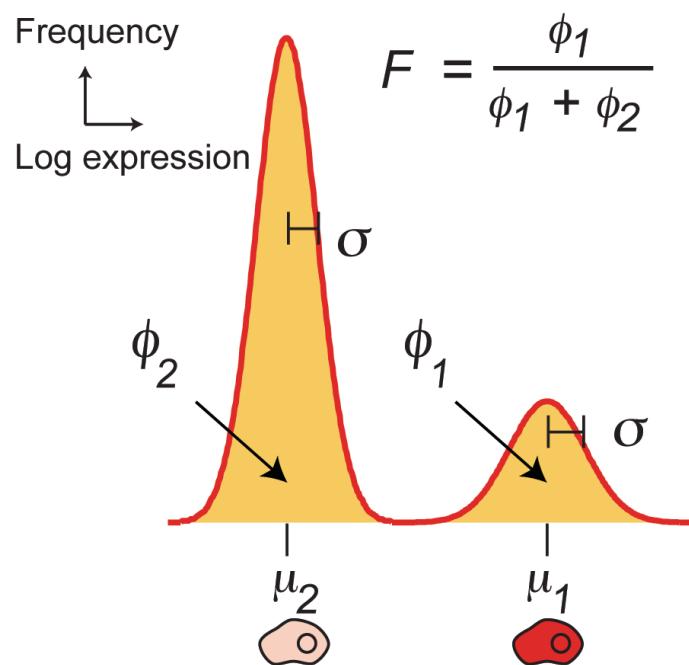
Dichotomous gene expression model

- Also have a model where the CVs are different for each lognormal population (rLN-LN)
- Also have a model where low-expression population is exponentially distribution

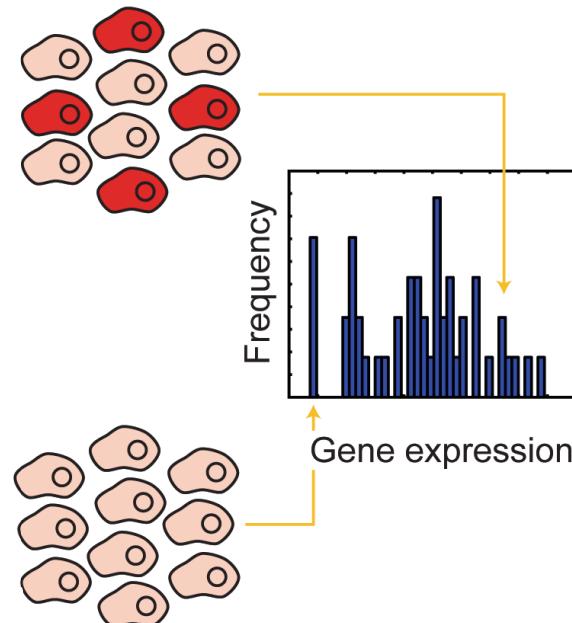


Parameterizing heterogeneity with Maximum-Likelihood Inference

1. Model of heterogeneity

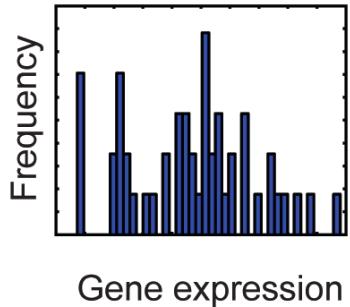


2. Random ten-cell profiling



Parameterizing heterogeneity with Maximum-Likelihood Inference

3. Maximize likelihood of pdf $f(x)$



measurements

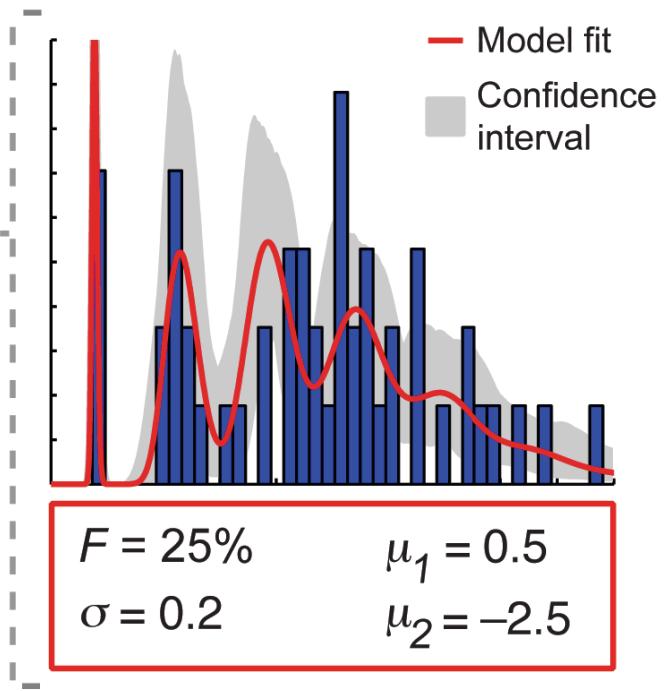
$$\rightarrow \text{Max} \sum_{i=1}^{10} \log f(x_i | \mu_1, \mu_2, F, \sigma)$$

$$f(x) = \sum_{j=0}^{10} \binom{10}{j} F^j (1-F)^{10-j} z_{j, 10-j}(x)$$

$z_{j, 10-j}(x)$ is pdf of $Z_1 + \dots + Z_{10}$, where

$$Z_c \stackrel{\text{ind}}{\sim} \begin{cases} \text{LN}(\mu_1, \sigma^2) & 1 \leq c \leq j \\ \text{LN}(\mu_2, \sigma^2) & j < c \leq 10 \end{cases}$$

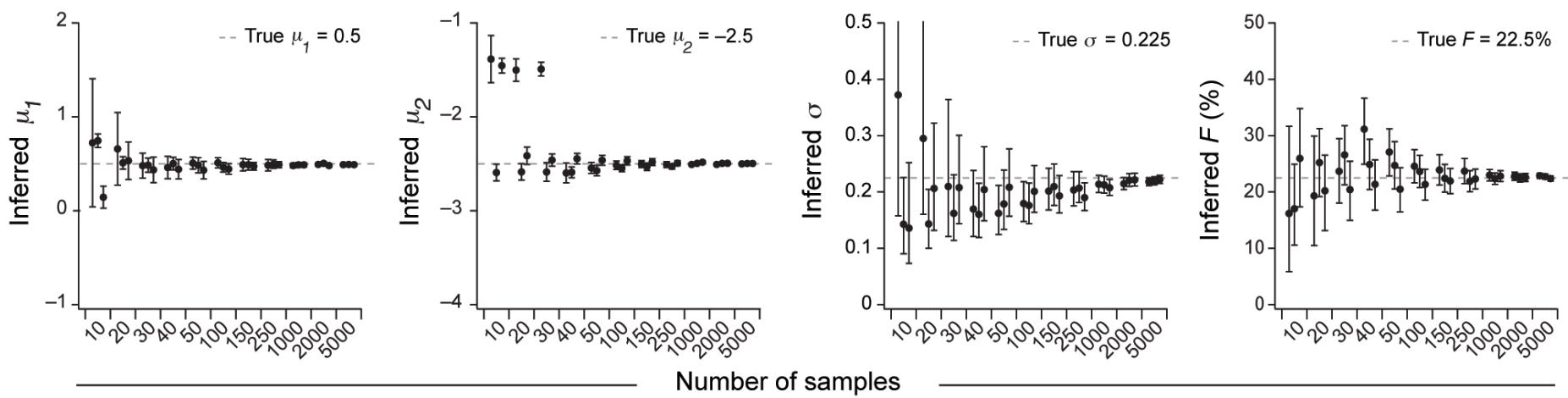
4. Estimate parameters



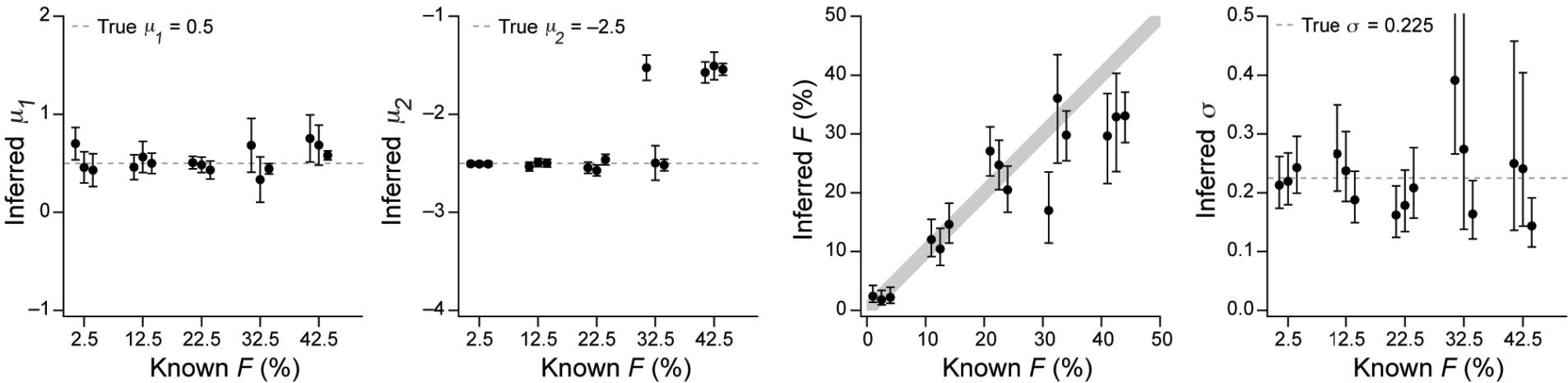
Maximum-likelihood estimation (MLE)

- Estimate parameters of a statistical model from a given data set
- Calculate the likelihood of the data
 - Probability of observing sample values given a parameter
- Maximizing the likelihood identifies the best parameter set of the model for the data

Sample-size requirements for inference

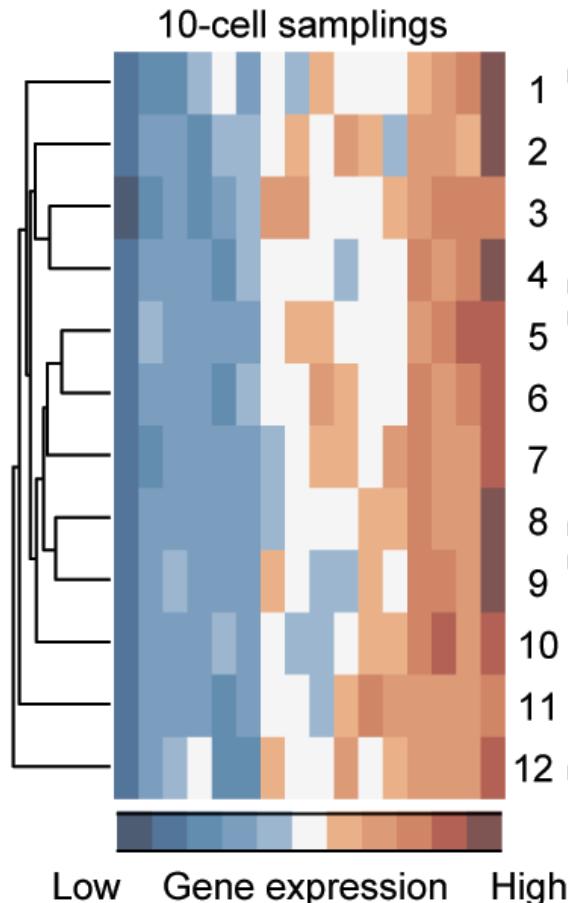


Parameter sensitivity in LN-LN estimation

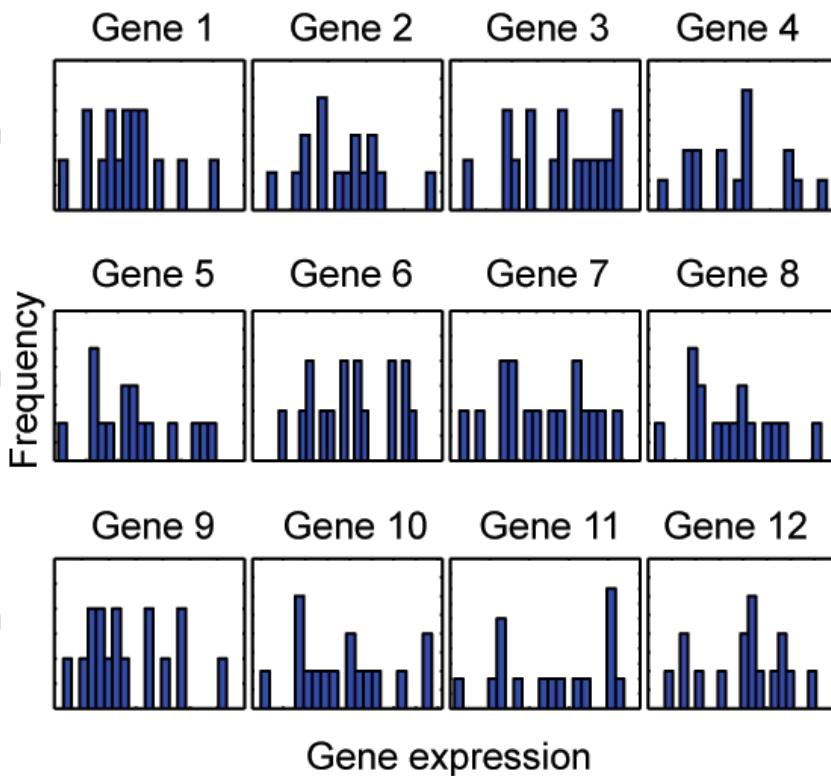


Globally parameterizing heterogeneities

1. Identify gene program

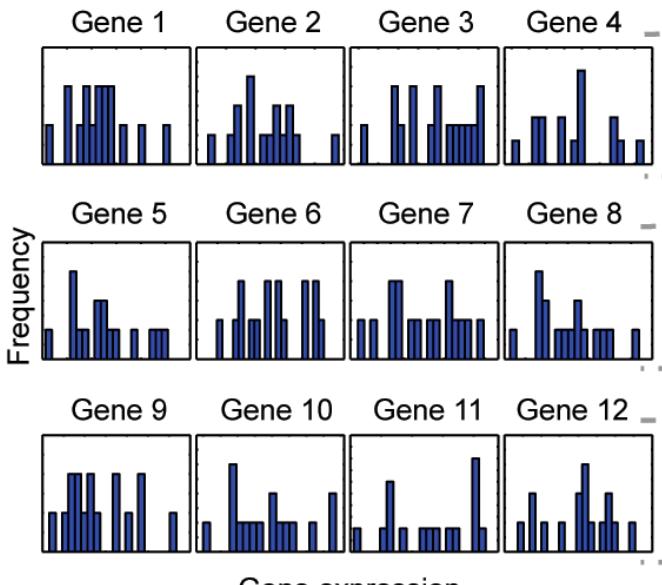


2. Estimate μ_1, μ_2 for each gene



Globally parameterizing heterogeneities

2. Estimate μ_1, μ_2 for each gene



3. Maximize likelihood of pdf $f(x)$

$$\text{Max } \sum_{k=1}^g \sum_{i=1}^m \log f^{(k)}(x_i^{(k)} | \hat{\mu}_1^{(k)}, \hat{\mu}_2^{(k)}, F, \sigma)$$

Fix μ_1, μ_2

$$f^{(k)}(x) = \sum_{j=0}^{10} \binom{10}{j} F^j (1-F)^{10-j} z_{j, 10-j}^{(k)}(x)$$

$z_{j, 10-j}^{(k)}(x)$ is pdf of $Z_1 + \dots + Z_{10}$, where

$$Z_c \stackrel{\text{ind}}{\sim} \begin{cases} \text{LN}(\mu_1^{(k)}, \sigma^2) & 1 \leq c \leq j \\ \text{LN}(\mu_2^{(k)}, \sigma^2) & j < c \leq 10 \end{cases}$$

$$\text{Max } \sum_{k=1}^g \sum_{i=1}^m \log f(x_i^{(k)} | \hat{\mu}_1^{(k)}, \hat{\mu}_2^{(k)}, F, \sigma)$$

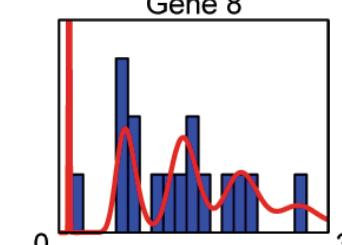
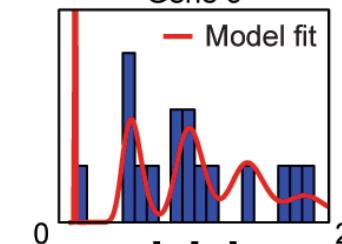
4. Estimate parameters

Cluster-wide parameters

$$F = 23\% \quad \sigma = 0.16$$

Gene 5

Gene 8



General work flow

- Input in data
 - Single-gene (e.g., qPCR) or cluster of genes
 - Not log transformed
- Perform inference with three mixture models
 - LN-LN, rLN-LN, EXP-LN
- Use Bayes' Information Criterion (BIC) to determine best model for data

Installing R package

- stochprofML package on CRAN
 - Install package from source file on thumbdrive or by downloading from CRAN

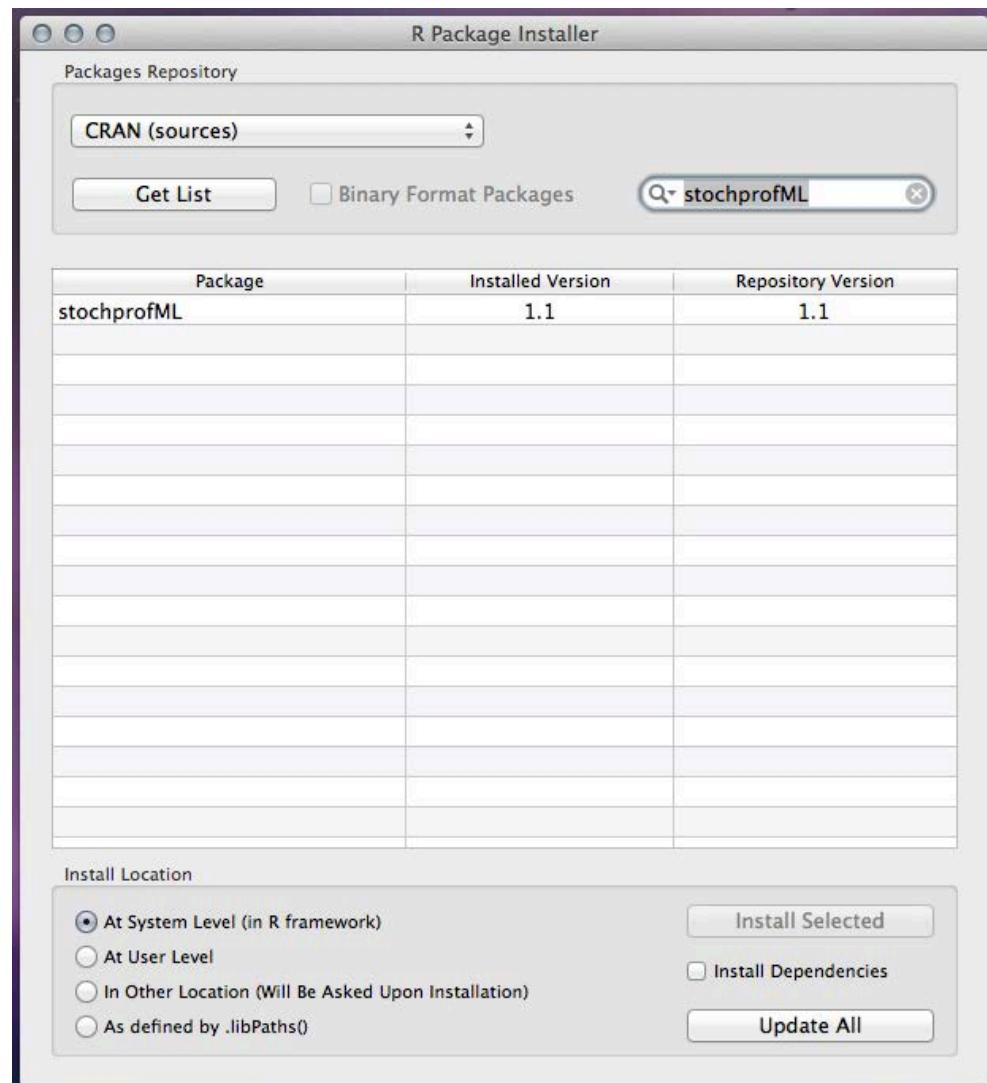


Image using R-Studio

Using R package

- Monte Carlo simulation variables to define
 - Model to use
 - LN-LN, rLN-LN (will take longer to run), EXP-LN (will take very long to run)
 - Number of cells sampled
 - Number of samples
 - F, μ_1, μ_2, σ

Using R package

- Use function **stochasticProfilingData()** to simulate data
- Use function **stochasticProfilingML()** to perform maximum likelihood inference for data (both simulated and real)

Exercises and discussion

- Generate simulated data and perform inference for LN-LN model
 - rLN-LN will take ~twice as long, EXP-LN will take ~1 hr
- Infer SOD2 expression frequency for LN-LN model
 - Data is stored in the package as the vector “sod2”
- During overflow or free time
 - Infer parameters for a cluster of 4 simulated genes

Key R functions

- `r.sum.of.mixtures.MODEL(k, n, F, mu, sigma)`
 - Simulate data
 - $k = \# \text{ of samples}$
 - $n = \# \text{ of cells}$
- `stochprof.loop`
 - Main function, takes in several arguments (see help file) and initiates inference algorithm

Using R package

```
samples <- c(NUMBER OF SAMPLES)
numberofcells <- c(NUMBER OF CELLS PER SAMPLE)
numberofpopulations <- c(NUMBER OF SUBPOPULATIONS)
p.vector <- c( $F$ ,  $1 - F$ )
mu.vector <- c( $\mu_1$ ,  $\mu_2$ )
sigma.vector <- c( $\sigma$ ,  $\sigma$ ) #sigmas are different in rLN-LN model
data <- r.sum.of.mixtures.LNLN(samples, numberofcells, p.vector, mu.vector,
sigma.vector)
#data <- qPCR data OR matrix of profiles where columns are genes, rows are samples
data <- matrix(data, nrow = samples, ncol = 1) #columns are gene, rows are samples
stochprof.loop("LN-LN", data, numberofcells, numberofpopulations)
```