

The EM algorithm

Computational Statistics

Adam Gorm Hoffmann

Setup where the EM algorithm is useful

- ▶ We observe incomplete data $X = M(Y)$.
- ▶ Log-likelihood $\log g(x | \theta)$ is impossible to find analytically, and expensive to compute numerically.
- ▶ Complete data log-likelihood $\log f(y | \theta)$ may be cheap to compute, but we don't observe y , only the incomplete data $x = M(y)$.
- ▶ It is cheap to compute $Q(\theta | \theta') = E_{\theta'}(\log f(Y | \theta) | X = x)$, the expected complete data log-likelihood in θ conditional on our incomplete data computed under the distribution $f_{\theta'}$.

The EM algorithm

Idea: Start with initial guess θ_0 and iteratively compute

$$\theta_{n+1} = \arg \max_{\theta} Q(\theta \mid \theta_n).$$

Stop when some convergence criterion is reached (e.g., after a specific number of iterations, or when $\|\theta_{n+1} - \theta_n\|$ is sufficiently small).

Why “EM”? For a given θ_n :

- ▶ **E**xpectation step: Compute the function $\theta \mapsto Q(\theta \mid \theta_n)$.
- ▶ **M**aximization step: Maximize $\theta \mapsto Q(\theta \mid \theta_n)$.

Each EM iteration increases the log-likelihood

Denote the cross-entropy $H(\theta \mid \theta') = E_{\theta'}(-\log h(Y \mid x, \theta) \mid X = x)$.

If $Q(\theta \mid \theta') > Q(\theta' \mid \theta')$, then

$$\begin{aligned}\ell(\theta) &= Q(\theta \mid \theta') + H(\theta \mid \theta') \\ &\geq Q(\theta \mid \theta') + H(\theta' \mid \theta') && \text{Gibbs' inequality} \\ &> Q(\theta' \mid \theta') + H(\theta' \mid \theta') && \text{Assumption} \\ &= \ell(\theta').\end{aligned}$$

General R implementation

```
1  ## Either Q or update_estimates must be supplied. If Q is supplied,
2  ## numerical optimization is used to update the estimates.
3  EM <- function(epsilon, par_init, update_estimates = NULL, Q = NULL,
4                callback = NULL)
5  {
6      if (is.null(update_estimates))
7          update_estimates <- get_update_estimates_general(Q)
8      epsilon_squared <- epsilon^2
9      new_est <- par_init
10     not_converged <- TRUE
11     while (not_converged) {
12         old_est <- new_est
13         new_est <- update_estimates(old_est)
14         not_converged <- crossprod(new_est - old_est) >
15             epsilon_squared * (crossprod(old_est) + epsilon)^2
16         if (!is.null(callback))
17             callback()
18     }
19     new_est
20 }
```

M step for general Q with optim

```
1  get_update_estimates_general <- function(Q)
2  {
3      force_all_args()
4      function(old_est)
5          optim(old_est, function(par) -Q(par, old_est),
6                control = list(reltol = 1e-16))$par
7  }
8
9  ## Forces evaluation of all arguments
10 force_all_args <- function()
11     as.list(parent.frame())
```

Specific problem setup

Let X follow the generalized t -distribution

$$f(x) = \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\sqrt{\nu\pi\sigma^2}\Gamma\left(\frac{\nu}{2}\right)} \left(1 + \frac{(x - \mu)^2}{\nu\sigma^2}\right)^{-\frac{\nu+1}{2}}.$$

There is no analytical solution for the MLE $(\hat{\mu}, \hat{\sigma}^2)$.

If $W \sim \chi_\nu^2$ and $X \mid W = w \sim N(\mu, \nu\sigma^2/w)$ then X follows the desired generalized t -distribution.

We can use the EM algorithm in this setup!

Specific problem setup

$Y = (X, W) \in \mathbb{R} \times (0, \infty)$ has joint density

$$f(x, w) = f(x \mid w)f(w) = \frac{1}{\sqrt{\pi\nu\sigma^2}2^{(\nu+1)/2}\Gamma(\nu/2)} w^{\frac{\nu-1}{2}} e^{-\frac{w}{2}\left(1+\frac{(x-\mu)^2}{\nu\sigma^2}\right)}$$

with fixed $\nu \in (0, \infty)$ and unknown parameters $\mu \in \mathbb{R}, \sigma^2 \in (0, \infty)$.

Goal: use the EM algorithm to estimate (μ, σ^2) when we only observe iid. X_1, \dots, X_n .

Simulating from the model

```
1  sim_y <- function(n, mu, sigma2, nu)
2  {
3      w <- rchisq(n, df = nu)
4      x <- rnorm(n, mean = mu, sd = sqrt(nu * sigma2 / w))
5      list(
6          x = x,
7          w = w
8      )
9  }
```

Numerical optimization in Stan

- ▶ Stan has a fast optimizer using LBFGS, BFGS or Newton.
- ▶ The marginal distribution as a Stan model:

```
data {  
  int<lower = 1> N;  
  real x[N];  
  int<lower = 1> nu;  
}  
parameters {  
  real mu;  
  real<lower = 0> sigma;  
}  
model {  
  x ~ student_t(nu, mu, sigma);  
}
```

Running the numerical optimization

```
1 stan_marginal <- stan_model(file = "marginal.stan")
2 stan_opt <- optimizing(stan_marginal,
3                       data = list(N = n,
4                                   x = y$x,
5                                   nu = nu),
6                       init = list(mu = 0, sigma = 1))
```

- ▶ Done instantly (3 ms when $n = 1000$, 900 ms when $n = 10^6$).
- ▶ We can use this to check that our implementation works.

The E step (deriving Q)

Up to an additive constant

$$Q(\mu, \sigma^2 \mid \mu_m, \sigma_m^2) = -\frac{n}{2} \log \sigma^2 - \frac{1}{2\nu\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \frac{1 + \nu}{1 + \frac{(x_i - \mu_m)^2}{\nu\sigma_m^2}}.$$

The M Step (maximizing Q)

Solving $\frac{d}{d\mu} Q(\mu, \sigma^2 \mid \mu_m, \sigma_m^2) = 0$ and $\frac{d}{d\sigma^2} Q(\mu, \sigma^2 \mid \mu_m, \sigma_m^2) = 0$ yields

$$\mu = \frac{\sum_{i=1}^n \frac{x_i}{1 + \frac{(x_i - \mu_m)^2}{\nu \sigma_m^2}}}{\sum_{i=1}^n \frac{1}{1 + \frac{(x_i - \mu_m)^2}{\nu \sigma_m^2}}}$$
$$\sigma^2 = \frac{1}{n\nu} \sum_{i=1}^n (x_i - \mu)^2 \frac{1 + n\nu}{1 + \frac{(x_i - \mu_m)^2}{\nu \sigma_m^2}}$$

Since Q is concave, this is a global maximum.

R implementation of Q

```
1  ## Only implemented up to additive constant,  
2  ## since we only need the derivatives.  
3  getQ <- function(x, nu)  
4  {  
5      force_all_args()  
6      function(par, par_prime)  
7          - length(x) * log(par[2]) / 2 -  
8            sum((x - par[1])^2 /  
9              (1 + (x - par_prime[1])^2 / (nu * par_prime[2]))) *  
10           (1 + nu) / (2 * nu * par[2])  
11  }
```

Profiling: General implementation

- ▶ Evaluating Q in `optim` is the only part that takes significant time.
- ▶ **Idea:** do specific implementation of M step.

R implementation of EM step

```
1  get_update_estimates_specific <- function(x, nu)
2  {
3      force_all_args()
4      n <- length(x)
5      function(old_est)
6      {
7          beta_inverse <- (1 + (x - old_est[1])^2 / (nu * old_est[2]))^(-1)
8          mu_new <- sum(x * beta_inverse) / sum(beta_inverse)
9          sigma2_new <- sum((x - mu_new)^2 *
10                          (1 + nu) * beta_inverse) / (n * nu)
11          c(mu = mu_new, sigma2 = sigma2_new)
12      }
13 }
```


Profiling: Specific implementation


- Updating the estimates is the only part that takes significant time.

```
EM <- function(epsilon, par_init, update_estimates = NULL, Q = NULL,
               callback = NULL)
{
  if (is.null(update_estimates))
    update_estimates <- get_update_estimates_general(Q)
  epsilon_squared <- epsilon^2
  new_est <- par_init
  not_converged <- TRUE
  while (not_converged) {
    old_est <- new_est
    new_est <- update_estimates(old_est)
    not_converged <- crossprod(new_est - old_est) >
      epsilon_squared * (crossprod(old_est) + epsilon)^2
    if (!is.null(callback))
      callback()
  }
  new_est
}
```

14320

Profiling: Specific implementation

```
get_update_estimates_specific <- function(x, nu)
{
  force(x)
  force(nu)
  function(old_est)
  {
    beta_inverse <- (1 + (x - old_est[1])^2 / (nu * old_est[2]))^(-1)
    mu_new <- sum(x * beta_inverse) / sum(beta_inverse)
    sigma2_new <- sum((x - mu_new)^2 * (1 + nu) * beta_inverse) / (n * nu)
    c(mu = mu_new, sigma2 = sigma2_new)
  }
}
```



Line	Time (ns)
beta_inverse <- (1 + (x - old_est[1])^2 / (nu * old_est[2]))^(-1)	11540
mu_new <- sum(x * beta_inverse) / sum(beta_inverse)	2400
sigma2_new <- sum((x - mu_new)^2 * (1 + nu) * beta_inverse) / (n * nu)	2660

- ▶ Everything is already vectorized.
- ▶ **Idea:** implement in Rcpp to avoid wasting time jumping between R and low-level libraries.

Rcpp implementation: estimates struct

```
1  struct estimates {  
2      double mu;  
3      double sigma2;  
4  };  
5  
6  estimates  
7  make_estimates(double mu, double sigma2)  
8  {  
9      estimates est;  
10     est.mu = mu;  
11     est.sigma2 = sigma2;  
12     return est;  
13 }
```

Rcpp implementation: highest level

```
1  estimates
2  EM(double epsilon, const NumericVector& x, double nu,
3     double mu_guess, double sigma2_guess)
4  {
5     double epsilon_squared = square(epsilon);
6     estimates old_est;
7     estimates new_est = make_estimates(mu_guess, sigma2_guess);
8     do {
9         old_est = new_est;
10        new_est = update_estimates(old_est, x, nu);
11    } while (diff_two_norm(new_est, old_est)
12            > epsilon_squared * square(two_norm(old_est) + epsilon));
13    return new_est;
14 }
15
16 // [[Rcpp::export]]
17 NumericVector
18 EM_cpp(double epsilon, const NumericVector& x, double nu,
19        double mu_guess, double sigma2_guess)
20 {
21     estimates est = EM(epsilon, x, nu, mu_guess, sigma2_guess);
22     NumericVector est_vec = {est.mu, est.sigma2};
23     return est_vec;
24 }
```

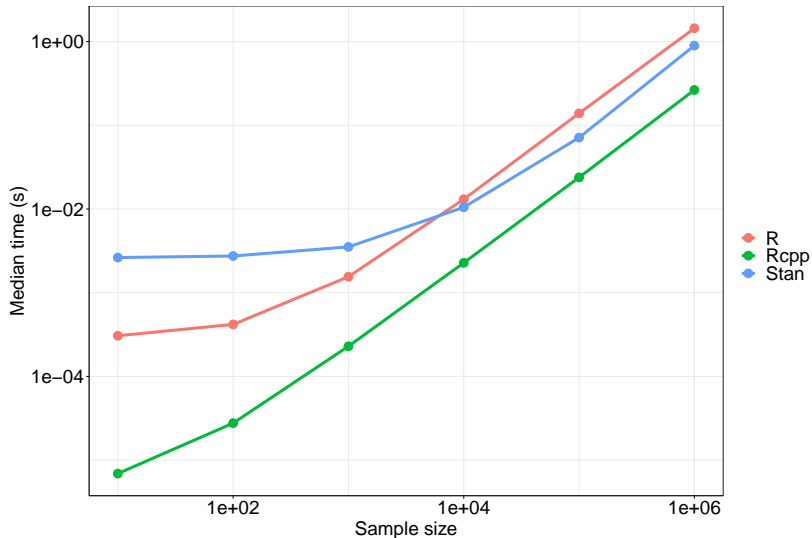
Rcpp implementation: EM step

```
1  estimates
2  update_estimates(estimates old_estimates, const NumericVector& x, double nu)
3  {
4      estimates new_estimates = make_estimates(0, 0);
5      double mu_divide_by = 0;
6      unsigned long long n_x = x.length();
7      double beta_inverse[n_x];
8
9      for (unsigned long long i = 0; i < n_x; i++) {
10         beta_inverse[i] = 1 / (1 + square((x[i] - old_estimates.mu))
11                                / (nu * old_estimates.sigma2));
12         new_estimates.mu += x[i] * beta_inverse[i];
13         mu_divide_by += beta_inverse[i];
14     }
15     new_estimates.mu /= mu_divide_by;
16     for (unsigned long long i = 0; i < n_x; i++)
17         new_estimates.sigma2 += square((x[i] - new_estimates.mu))
18                                * (1 + nu) * beta_inverse[i];
19     new_estimates.sigma2 /= (n_x * nu);
20     return new_estimates;
21 }
```

Rcpp implementation: small functions

```
1  double
2  square(double x)
3  {
4      return x * x;
5  }
6
7  double
8  two_norm(estimates est)
9  {
10     return square(est.mu) + square(est.sigma2);
11 }
12
13 double
14 diff_two_norm(estimates new_est, estimates old_est)
15 {
16     return square(new_est.mu - old_est.mu)
17         + square(new_est.sigma2 - old_est.sigma2);
18 }
```

Benchmarking EM implementations and Stan



- The Rcpp implementation is roughly 10 times faster than R, and also beats Stan's numerical optimization.

Gradient ascent: Marginal log-likelihood and gradient

```
1  get_marginal_loglik <- function(x, nu)
2  {
3      force_all_args()
4      function(par)
5          -n/2 * log(par[2]) -
6              (nu + 1)/2 * sum(log(1 + (x - par[1])^2 / (nu * par[2])))
7  }
8
9  get_marginal_grad <- function(x, nu)
10 {
11     force_all_args()
12     function(par)
13         c(
14             (nu + 1) * sum((x - par[1]) / ((x - par[1])^2 + nu * par[2])),
15             nu/2 * sum( ((x - par[1])^2 - par[2]) /
16                         (par[2] * (nu * par[2] + (x - par[1])^2)) )
17         )
18 }
```


Gradient ascent

```
1  grad_ascent <- function(init_guess, objective, grad,
2                             gamma = 0.01, epsilon = 1e-16,
3                             callback = NULL)
4  {
5      small_relative_ascent <- function(new_est, old_est)
6      {
7          objective_diff <- objective(new_est) - objective(old_est)
8          objective_diff >= 0 &&
9              objective_diff <= epsilon * (abs(objective(new_est)) + epsilon))
10     }
11
12     converged <- FALSE
13     new_est <- init_guess
14     while (!converged) {
15         old_est <- new_est
16         gr <- grad(old_est)
17         new_est <- old_est + gamma * gr
18         converged <- small_relative_ascent(new_est, old_est)
19         if (!is.null(callback))
20             callback()
21     }
22
23     new_est
24 }
```

MLE from complete data log-likelihood

If, in addition to X_1, \dots, X_n , we also observe W_1, \dots, W_n , then

$$\hat{\mu} = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i},$$
$$\hat{\sigma}^2 = \frac{1}{n\nu} \sum_{i=1}^n w_i (x_i - \hat{\mu})^2.$$

```
1 mle <- function(x, w, nu)
2 {
3   mu_hat <- sum(w * x) / sum(w)
4   c(mu_hat = mu_hat,
5     sigma2_hat = mean(w * (x - mu_hat)^2) / nu)
6 }
```

Testing that the implementations work

```
1  ## True parameters
2  mu <- 27
3  sigma2 <- 13
4  nu <- 5
5
6  ## Simulating data
7  n <- 1000
8  y <- sim_y(n, mu, sigma2, nu)
9
10 ## Running EM algorithm
11 update_estimates_specific <- get_update_estimates_specific(y$x, nu)
12 Q <- getQ(y$x, nu)
13 est_specific <- EM(1e-16, c(1,1),
14                   update_estimates = update_estimates_specific)
15 est_general <- EM(1e-16, c(1,1), Q = Q)
16 est_Rcpp <- EM_cpp(1e-16, y$x, nu, 1, 1)
17
18 ## Running gradient ascent
19 marginal_loglik <- get_marginal_loglik(y$x, nu)
20 marginal_grad <- get_marginal_grad(y$x, nu)
21 grad_ascent(c(1, 1), marginal_loglik, marginal_grad)
22
23 ## Calculate complete data MLE
24 est_mle <- mle(y$x, y$w, nu)
```

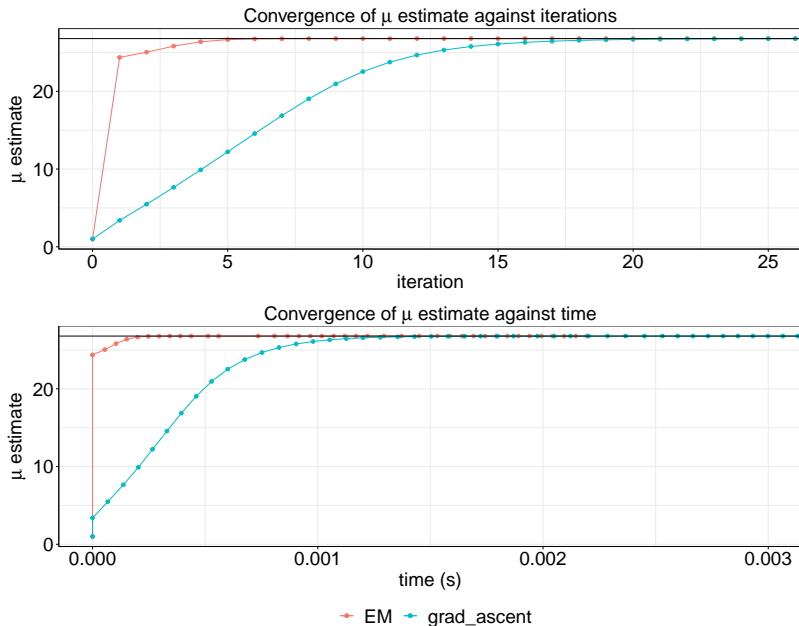
Results of test

- ▶ All EM implementations, the gradient ascent implementation, and Stan give the estimate $\tilde{\mu} = 26.78594$, $\tilde{\sigma}^2 = 13.10174$ up to 7 significant digits.
- ▶ The differences between the specific R and specific Rcpp estimates are of order 10^{-15} and the differences to the more general estimate and the gradient ascent estimate are of order 10^{-10} to 10^{-6} .
- ▶ The complete data MLE is $\hat{\mu} = 26.83363$ and $\hat{\sigma}^2 = 13.21392$, so the marginal data MLE (from the EM algorithm and gradient ascent) seems reasonable.
- ▶ Since Stan agrees with my EM implementations and gradient ascent implementation, they must be correct.

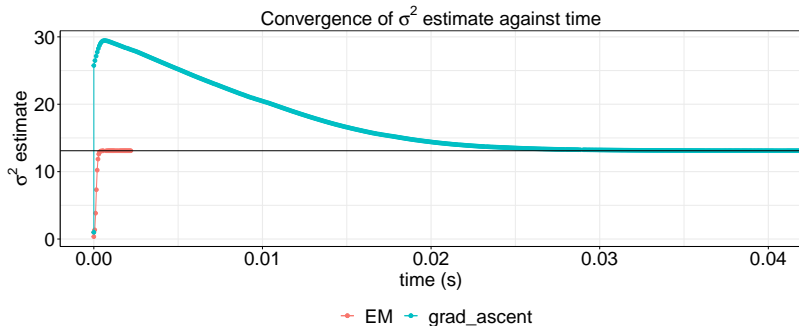
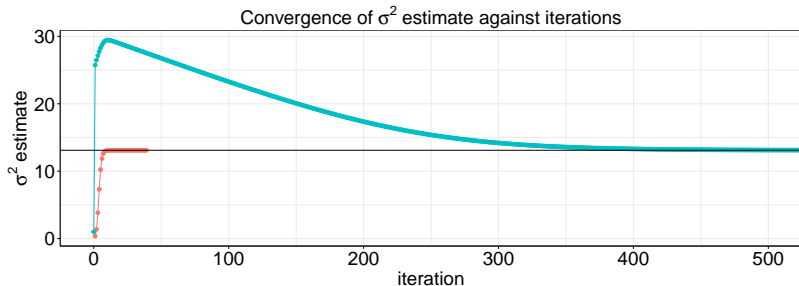
Comparing convergence of EM and gradient ascent

```
1  em_tracer <- tracer(c("old_est", "new_est"))
2  EM(1e-16, c(1,1), update_estimates = update_estimates_specific,
3     callback = em_tracer$tracer)
4
5  grad_tracer <- tracer(c("old_est", "new_est"))
6  grad_ascent(c(1, 1), marginal_loglik, marginal_grad,
7             callback = grad_tracer$tracer)
```

Convergence of μ estimate



Convergence of σ^2 estimate



Linear convergence

If

$$\limsup_{n \rightarrow \infty} \frac{\|\theta_n - \theta_\infty\|}{\|\theta_{n-1} - \theta_\infty\|} = r$$

for some $r \in (0, 1)$, then we say that the convergence is linear with asymptotic rate r .

r close to 0 means large improvements in each step, and r close to 1 means small improvements in each step.

Estimating the convergence rate

If the linear convergence is reached at n_0 , then (simplified)

$$\|\theta_n - \theta_\infty\| \approx r \|\theta_{n-1} - \theta_\infty\| \approx \dots \approx r^{n-n_0} \|\theta_{n-n_0} - \theta_\infty\|,$$

so

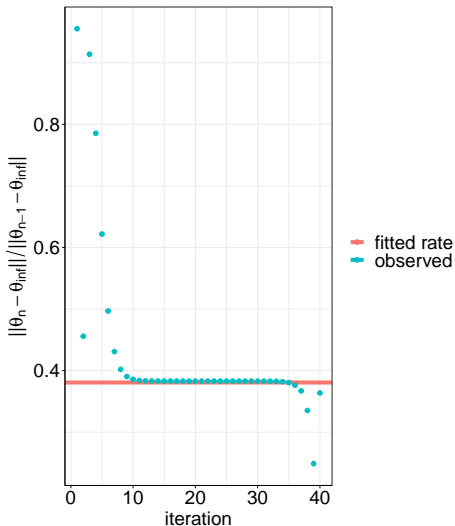
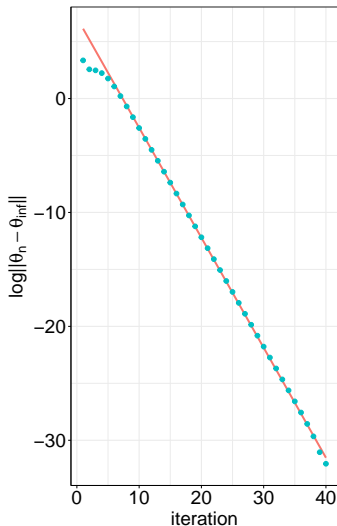
$$\log \|\theta_n - \theta_\infty\| \approx (n - n_0) \log(r) + \|\theta_{n-n_0} - \theta_\infty\|,$$

hence we can estimate $\log(r)$ as β when fitting the line

$$\log \|\theta_n - \theta_\infty\| = n\beta + \alpha.$$

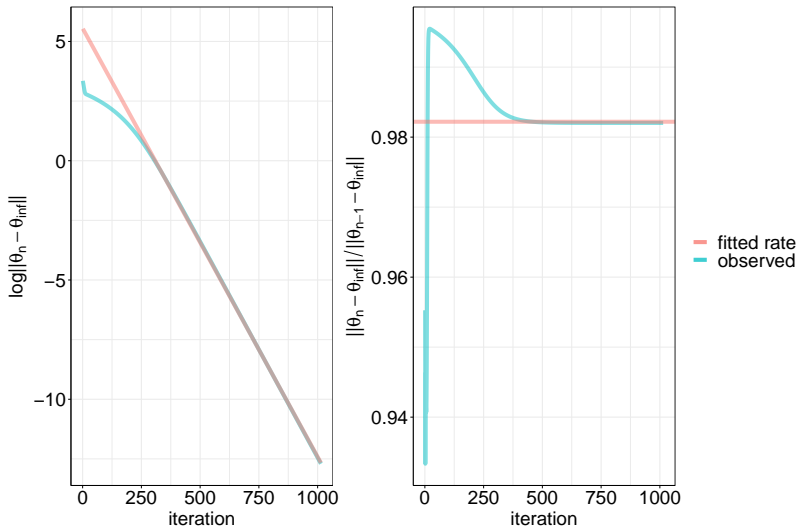
Convergence rate of EM

Linear convergence with $\hat{r} = \exp(\text{slope}) \approx 0.38$.



Convergence rate of gradient ascent

Linear convergence with $\hat{r} = \exp(\text{slope}) \approx 0.98$.



Conclusion from convergence investigation

- ▶ The EM algorithm converges in fewer steps and has better convergence rate than gradient ascent in our setup.
 - ▶ However, I have not spent much time tuning the parameters for gradient ascent, so it might be possible to do much better.
- ▶ My EM algorithm R implementation also converges in much less time than my gradient ascent R implementation.

Fisher information formulas

Let $\hat{i}_X = -D_{\hat{\theta}}^2 \ell(\hat{\theta})$ denote the observed Fisher information.

Then

$$\hat{i}_X = -D_{\bar{\theta}} \left(\nabla_{\theta} Q(\bar{\theta} \mid \bar{\theta}) \right) \Big|_{\bar{\theta}=\hat{\theta}}$$

$$\hat{i}_X = -D_{\hat{\theta}}^2 Q(\hat{\theta} \mid \hat{\theta}) - D_{\theta'} \nabla_{\theta} Q(\hat{\theta} \mid \hat{\theta})$$

$$\hat{i}_X = - \left(I - D_{\theta} \Phi(\hat{\theta})^T \right) D_{\hat{\theta}}^2 Q(\hat{\theta} \mid \hat{\theta})$$

where

$$\Phi(\theta') = \arg \max_{\theta} Q(\theta \mid \theta')$$

is the EM-map.

Implementing the formulas in R

Using the numDeriv package.

```
1  ## Calculate using definition
2  fisher_def <- function(est, loglik, x)
3    hessian(function(par) loglik(par, x = x), x = est)
4
5  ## Calculate with the three formulas
6  fisher1 <- function(est, Q)
7    -jacobian(function(par) grad(Q, par, par_prime = par), est)
8
9  D2Q <- function(est, Q)
10    hessian(Q, est, par_prime = est)
11
12  fisher2 <- function(est, Q)
13    -D2Q(est, Q) -
14      jacobian(function(par_prime) grad(Q, est, par_prime = par_prime),
15                est)
16
17  fisher3 <- function(est, Q, update_estimates = NULL)
18  {
19    if (is.null(update_estimates))
20      update_estimates <- get_update_estimates_general(Q)
21    -(diag(length(est)) - t(jacobian(update_estimates, est))) %*%
22      D2Q(est, Q)
23  }
```

They all give approximately the same result

$$\text{fisher_def: } \begin{pmatrix} 57.48859232 & -0.01589346 \\ -0.01589346 & 1.79719876 \end{pmatrix}$$

$$\text{fisher1: } \begin{pmatrix} 57.48859187 & -0.01589217 \\ -0.01589422 & 1.79720353 \end{pmatrix}$$

$$\text{fisher2: } \begin{pmatrix} 57.4885922 & -0.01589327 \\ -0.0158944 & 1.79719618 \end{pmatrix}$$

$$\text{fisher3: } \begin{pmatrix} 57.48859235 & -0.01589349 \\ -0.01589349 & 1.79719876 \end{pmatrix}$$

Empirical Fisher information

$$\hat{\mathcal{I}} = \sum_{i=1}^n \nabla_{\theta} Q_i(\hat{\theta} \mid \hat{\theta}) \nabla_{\theta} Q_i(\hat{\theta} \mid \hat{\theta})^T.$$

```
1 emp_fisher <- function(mu, sigma2, x, nu)
2 {
3   w_means <- (nu + 1) * 0.5 / (0.5 * (1 + ((x - mu)^2 / (sigma2 * nu))))
4   dQ_dmu <- w_means * ((x - mu) / (nu * sigma2))
5   dQ_dsigma2 <- -1/(2*sigma2) +
6     (w_means * (x - mu)^2) / (2 * nu * sigma2^2)
7   grad_list <- cbind(mu = dQ_dmu, sigma2 = dQ_dsigma2)
8   crossprod(grad_list)
9 }
```

$$\begin{pmatrix} 57.2396144 & -0.3093977 \\ -0.3093977 & 1.9327784 \end{pmatrix}$$

► Close, but not identical, to the observed Fisher information.