

# **Stochastic Gradient Descent**

## **Computational Statistics**

Adam Gorm Hoffmann

# Setup for stochastic gradient descent

- ▶ Estimate parameter in  $\Theta$  based on a large number of observations from  $\mathcal{X}$ .
- ▶ With some loss function  $L : \mathcal{X} \times \Theta \rightarrow \mathbb{R}$  the parameter can be estimated by minizing the risk (expected loss)  
 $H(\theta) = E(L(X, \theta))$  over  $\theta \in \Theta$ .

# Gradient Descent

Suppose that we have  $N$  iid. observations  $X_1, \dots, X_N$ .

A gradient descent using all the data, would be

$$\theta_{n+1} = \theta_n - \gamma_n \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} L(X_i, \theta_n)$$

with some decay schedule  $(\gamma_n)$ .

**Problem:** If the number of observations,  $N$ , is very large, or if the gradient  $\nabla_{\theta} L(X_i, \theta_n)$  is very expensive to compute, then this will be a very time consuming algorithm.

# Stochastic Gradient Descent (mini-batch)

**Idea:** Only use some of the data in each update.

Let  $M$  be the desired number of data points used for each update, and let  $I_n$  be a random sample of  $M$  indices from  $\{1, \dots, N\}$ . Then we compute the  $n$ 'th update using the  $M$  observations given by  $I_n$ .

$$\theta_{n+1} = \theta_n - \gamma_n \frac{1}{M} \sum_{i \in I_n} \nabla_{\theta} L(X_i, \theta_n).$$

# Epochs

We further split the computations into *epochs* where (almost) all data is seen once in each epoch. In epoch  $n$  we sample  $I_{n,j}$  such that

$$\bigcup_{j=1}^{\text{floor}(N/M)} I_{n,j} \approx \{1, \dots, N\}.$$

(If  $\text{floor}(N/M) < N/M$  then we leave out the remaining data points for simplicity.)

The learning rate  $\gamma_n$  stays the same within epoch  $n$ .

# Implementation: highest level

```
1  sgd <- function(par, n_obs, decay_schedule, epoch = batch,  
2                n_iter = 100, sampler = sample, cb = NULL, ...)  
3  {  
4    learning_rates <- decay_schedule(1:n_iter)  
5    for(k in 1:n_iter) {  
6      if(!is.null(cb))  
7        cb()  
8      epoch_sample <- sampler(n_obs)  
9      par <- epoch(par, epoch_sample, learning_rates[k], grad, ...)  
10   }  
11   if (!is.null(cb))  
12     cb()  
13   par  
14 }
```

# Implementation: mini-batches

```
1 batch <- function(par, epoch_sample, learning_rate, grad,  
2                   batch_size = 50, ...)  
3 {  
4   n_batches <- floor(length(epoch_sample) / batch_size)  
5   for(j in 0:(n_batches - 1)) {  
6     i <- epoch_sample[(j * batch_size + 1):  
7                       (j * batch_size + batch_size)]  
8     par <- par - learning_rate * grad(par, i, ...)  
9   }  
10  par  
11 }
```

Other update schemes:

- ▶ Momentum: remember last direction.
- ▶ Adam: rescale gradient components to avoid narrow valley.

## Specific problem: logistic regression model

$p_i(\beta) = P(Y_i = 1 \mid X_i = x_i)$  with log-odds

$$\log \frac{p_i(\beta)}{1 - p_i(\beta)} = f(x_i \mid \beta) = (\varphi_1(x_i), \dots, \varphi_p(x_i))^T \beta$$

where  $\beta \in \mathbb{R}^p$  and with fixed basis functions  $\varphi_1, \dots, \varphi_p : \mathbb{R} \rightarrow \mathbb{R}$ .

We want to optimize the penalized negative log likelihood

$$H(\beta) = -\frac{1}{N} \sum_{i=1}^N (y_i \log p_i(\beta) + (1 - y_i) \log (1 - p_i(\beta))) + \lambda \left\| f''_{\beta} \right\|_2^2$$

over  $\beta \in \mathbb{R}^p$ .



## Deriving $D_\beta H$

Let  $\Phi_{ij} = \varphi_j(x_i)$  and

$$\Omega_{ij} = \langle \varphi_i'', \varphi_j'' \rangle = \int \varphi_i''(z) \varphi_j''(z) dz$$

and  $h(z) = -\log(1 + \exp(z))$ .

## Deriving $D_{\beta}H$

Then

$$\log p(\beta) = \Phi\beta + h(\Phi\beta)$$

$$\log(\mathbf{1} - p(\beta)) = h(\Phi\beta).$$

so

$$\begin{aligned} H(\beta) &= -\frac{1}{N}\mathbf{y}^T(\Phi\beta + h(\Phi\beta)) - \frac{1}{N}(\mathbf{1} - \mathbf{y})^T h(\Phi\beta) + \lambda\beta^T\Omega\beta \\ &= \frac{1}{N}\left(-\mathbf{y}^T\Phi\beta - \mathbf{1}^T h(\Phi\beta)\right) + \lambda\beta^T\Omega\beta \end{aligned}$$

with derivatives

$$D_{\beta}H(\beta) = -\frac{1}{N}\Phi^T(\mathbf{y} - p(\beta)) + 2\lambda\Omega\beta$$

$$D_{\beta}^2H(\beta) = \frac{1}{N}\Phi^T W\Phi + 2\lambda\Omega$$

where  $W(\beta) = \text{diag}(p(\beta)) \text{diag}(1 - p(\beta))$ .

# Implementation of $H$

$$H(\beta) = -\frac{1}{N}(\mathbf{y}^T \log p(\beta) + (\mathbf{1} - \mathbf{y})^T \log(\mathbf{1} - p(\beta))) + \lambda \beta^T \Omega \beta$$

```
1 get_H <- function(lambda, x_vec, y_vec, knots, inner_knots)
2 {
3   force_all_args()
4   N <- length(x_vec)
5   Omega_mat <- Omega(inner_knots)
6   Phi_mat <- Phi(x_vec, knots)
7   function(beta) {
8     p_beta <- p(beta, Phi_mat[i, ])
9     -1/N * (crossprod(y_vec, log(p_beta)) +
10             crossprod(1 - y_vec, log(1 - p_beta))) +
11     lambda * crossprod(beta, Omega_mat %*% beta)
12   }
13 }
```

Forcing evaluation of arguments to avoid lazy evaluation problems

```
1 force_all_args <- function()
2   as.list(parent.frame())
```

# Implementation of $D_\beta H$

$$D_\beta H(\beta) = -\frac{1}{N} \Phi^T (\mathbf{y} - p(\beta)) + 2\lambda \Omega \beta$$

```
1  get_grad_H <- function(lambda, x_vec, y_vec, knots, inner_knots)
2  {
3      force_all_args()
4      Phi_mat <- Phi(x_vec, knots)
5      unpen_grad <- function(beta, i) {
6          n_i <- length(i)
7          Phi_i <- matrix(Phi_mat[i, ], nrow = n_i)
8          -crossprod(Phi_i, y_vec[i] - p(beta, Phi_i)) / n_i
9      }
10     if (lambda > 0) {
11         two_lambda_Omega <- 2 * lambda * Omega(inner_knots)
12         function(beta, i)
13             unpen_grad(beta, i) + two_lambda_Omega %*% beta
14     } else
15         unpen_grad
16 }
```

- Precompute as much as possible.
- Leave out penalization term if  $\lambda = 0$  to avoid wasting time on matrix multiplication for a 0-matrix.

# Implementation: components for $H$ and $D_{\beta\ell}$

```
1  inv_logit <- function(x)
2    exp(x) / (1 + exp(x))
3
4  get_knots <- function(inner_knots)
5    sort(c(rep(range(inner_knots), 3), inner_knots))
6
7  Phi <- function(x_vec, knots)
8    splineDesign(knots, x_vec)
9
10 p <- function(beta, phi_mat)
11   inv_logit(phi_mat %*% beta)
12
13 Omega <- function(inner_knots)
14 {
15   knots <- sort(c(rep(range(inner_knots), 3), inner_knots))
16   d <- diff(inner_knots) # The vector of knot differences; b - a
17   g_ab <- splineDesign(knots, inner_knots, derivs = 2)
18   knots_mid <- inner_knots[-length(inner_knots)] + d / 2
19   g_ab_mid <- splineDesign(knots, knots_mid, derivs = 2)
20   g_a <- g_ab[-nrow(g_ab), ]
21   g_b <- g_ab[-1, ]
22   (crossprod(d * g_a, g_a) +
23    4 * crossprod(d * g_ab_mid, g_ab_mid) +
24    crossprod(d * g_b, g_b)) / 6
25 }
26
27 d2f_two_norm_squared <- function(beta, inner_knots)
28   crossprod(beta, Omega(inner_knots) %*% beta)
```

# Decay schedule

```
1 decay_scheduler <- function(gamma0 = 1, a = 1, K = 1, gamma1, n1)
2 {
3   force_all_args()
4   if (!missing(gamma1) && !missing(n1))
5     K <- n1^a * gamma1 / (gamma0 - gamma1)
6   b <- gamma0 * K
7   function(n)
8     b / (K + n^a)
9 }
```

# Newton's method: highest level

```
1  newton <- function(init_guess, x, y, max_iter = 100, epsilon = 1e-5,  
2                        lambda = 0, cb = NULL)  
3  {  
4      n_param <- length(init_guess)  
5      par_new <- init_guess  
6      newton_update <- get_newton_update(n_param, x, y, lambda)  
7      for (i in 1:max_iter) {  
8          if (!is.null(cb))  
9              cb()  
10         par_old <- par_new  
11         par_new <- newton_update(par_old)  
12         if (sum((par_new - par_old)^2) <=  
13             epsilon * (sum(par_new^2) + epsilon))  
14             break  
15     }  
16     if (!is.null(cb))  
17         cb()  
18     par_new  
19 }
```

# Newton's method: update step

The solution to

$$D_{\beta}H(\beta^0) + (\beta - \beta^0)^T D_{\beta}^2H(\beta^0) = 0$$

is given by

$$\beta = (\Phi^T W \Phi + N \lambda \Omega)^{-1} \Phi^T W z_0$$

where

$$z_0 = \Phi \beta^0 + W(\beta^0)^{-1}(y - p(\beta^0))$$

and

$$W(\beta) = \text{diag}(p(\beta)) \text{diag}(1 - p(\beta)).$$



# Newton's method: update step

```
1  get_newton_update <- function(n_param, x, y, lambda)
2  {
3      force_all_args()
4      n_sim <- length(x)
5      inner_knots <- seq(min(x), max(x), length.out = n_param - 2)
6      knots <- get_knots(inner_knots)
7      Phi_mat <- Phi(x, knots)
8      tPhi <- t(Phi_mat)
9      Omega_mat <- Omega(inner_knots)
10
11     function(beta)
12     {
13         p_vec <- as.vector(p(beta, Phi_mat))
14         p_prod <- p_vec * (1 - p_vec)
15         W <- diag(p_prod)
16         W_inv <- diag(1 / p_prod)
17         z <- Phi_mat %*% beta + W_inv %*% (y - p(beta, Phi_mat))
18         solve(crossprod(Phi_mat, W %*% Phi_mat) +
19              n_sim * lambda * Omega_mat) %*%
20             crossprod(Phi_mat, W %*% z)
21     }
22 }
```

# Simulating data for initial tests

```
1  set.seed(2021)
2  n_beta <- 5
3  n_sim <- 1000
4  x <- seq(0, 1, length.out = n_sim)
5  inner_knots <- seq(min(x), max(x), length.out = n_beta - 2)
6  knots <- get_knots(inner_knots)
7  beta <- c(-1, 1, 3, -7, 1)
8  y <- as.numeric(runif(n_sim) < p(beta, x, knots))
9  beta_init <- rep(1, n_beta)
```

# Getting $H$ and $D_\beta H$ with $\lambda = 0$

- We set  $\lambda = 0$ , corresponding to usual unpenalized logistic regression, so we can compare results with `glm`.

```
1 grad_H <- get_grad_H(lambda = 0, x, y, knots, inner_knots)
2 H <- get_H(lambda = 0, x, y, knots, inner_knots)
```

- We confirm that our gradient implementation works by comparing to a numerical approximation.

```
1 all.equal(numDeriv::grad(H, beta),
2           as.vector(grad_H(beta, 1:length(x))), tolerance = 1e-8)
3
4 ## TRUE
```

# Run methods

```
1  set.seed(2021)
2
3  ## Stochastic gradient descent
4  sgd_batch_tracer <- tracer(c("objective", "par"),
5                             expr = quote(objective <- H(par)),
6                             N = 0)
7  ds <- decay_scheduler(gamma0 = 0.8, gamma1 = 0.008, n1 = 90)
8  beta_fit_1 <- sgd(par = beta_init,
9                    grad = grad_H,
10                     n_obs = n_sim,
11                     decay_schedule = ds,
12                     epoch = batch,
13                     n_iter = 200,
14                     batch_size = 1,
15                     cb = sgd_batch_tracer$tracer)
16
17  ## Newton
18  newton_tracer <- tracer(c("objective", "par_new"),
19                          expr = quote(objective <- H(par_new)))
20  newton_fit <- newton(init_guess = beta_init,
21                      x = x,
22                      y = y,
23                      max_iter = 100,
24                      cb = newton_tracer$tracer)
```

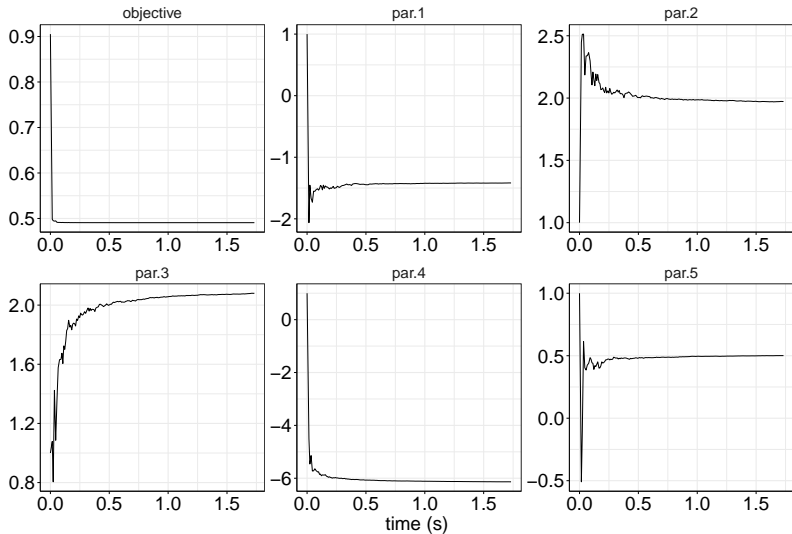
# Comparing results to glm (IWLS)

```
1 glm(y ~ splineDesign(knots, x) - 1, family = binomial)$coefficients
```

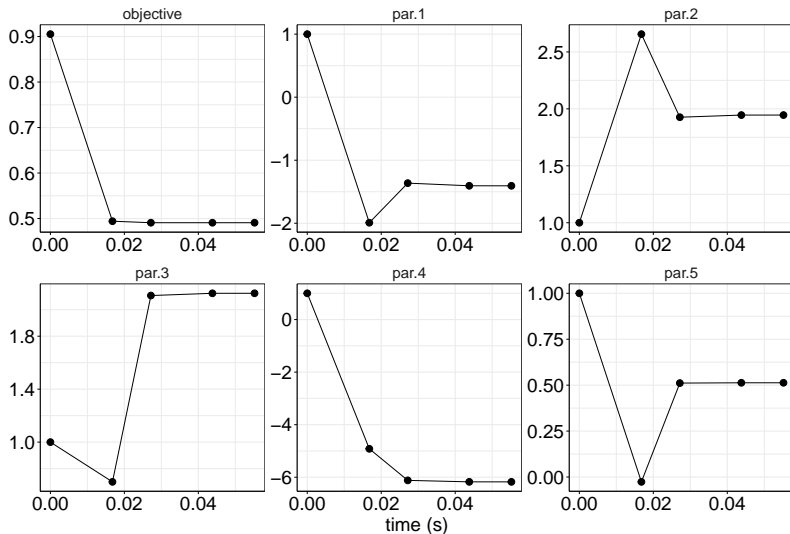
sgd	newton	glm
-1.4172427	-1.4058146	-1.4058146
1.9725305	1.9454677	1.9454677
2.0802703	2.1241283	2.1241283
-6.1350949	-6.1761810	-6.1761811
0.5009762	0.5129334	0.5129334

The methods yield approximately the same results, indicating that our implementations work with  $\lambda = 0$ .

# SGD convergence (200 iterations)

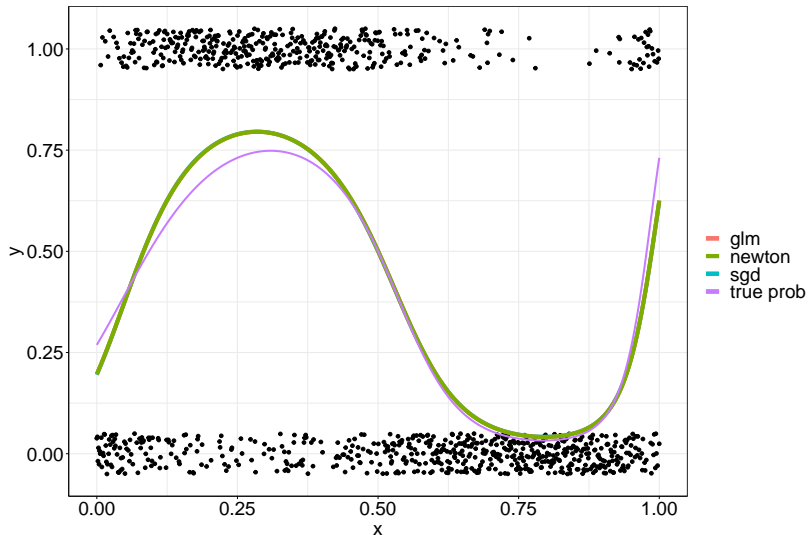


# Newton convergence (4 iterations)



► Okay convergence two steps from initial guess!

# Predictions (the same for all three methods)





## $\lambda > 0$ : Run methods

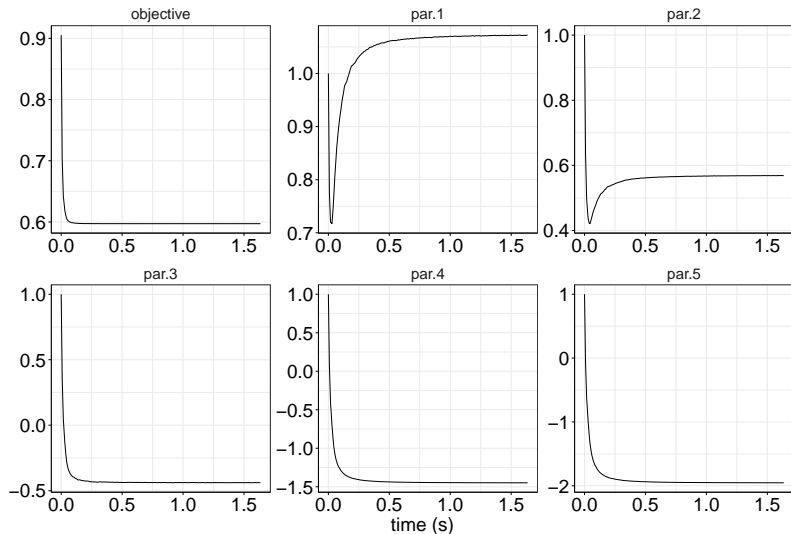
```
1 grad_H <- get_grad_H(lambda = 0.2, x, y, knots, inner_knots)
2 H <- get_H(lambda = 0.2, x, y, knots, inner_knots)
3 opt_fit <- optim(par = beta_init, fn = H, method = "BFGS",
4                 control = list(reltol = 1e-16))$par
5 ## Stochastic gradient descent
6 sgd_batch_tracer <- tracer(c("objective", "par"),
7                             expr = quote(objective <- H(par)),
8                             N = 0)
9 beta_fit_lam <- sgd(par = beta_init, grad = grad_H, n_obs = n_sim,
10                    decay_schedule = decay_scheduler(gamma0 = 0.015,
11                                                      gamma1 = 0.001,
12                                                      n1 = 90),
13                    epoch = batch, n_iter = 200, batch_size = 1,
14                    cb = sgd_batch_tracer$tracer)
15 ## Newton
16 newton_tracer <- tracer(c("objective", "par_new"),
17                         expr = quote(objective <- H(par_new)))
18 newton_fit_lam <- newton(init_guess = beta_init, x = x, y = y,
19                          max_iter = 100, lambda = 0.2,
20                          cb = newton_tracer$tracer)
```

## $\lambda > 0$ : Comparing $\beta$ estimates

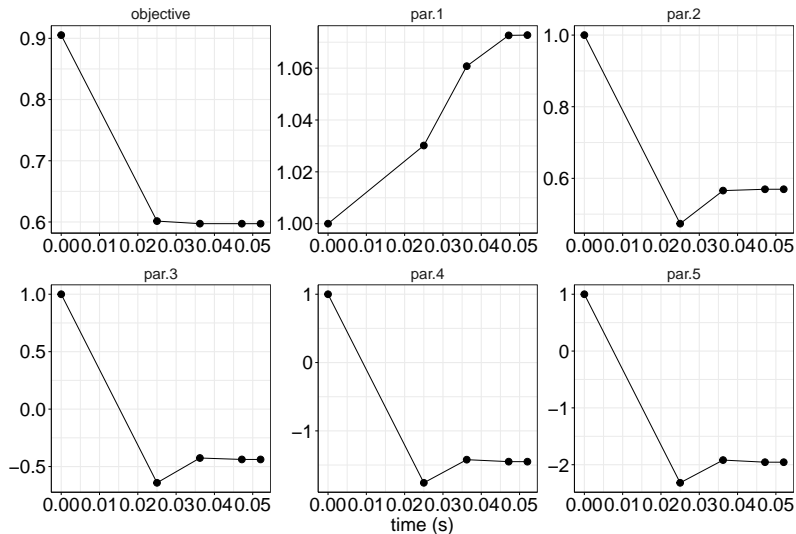
sgd	newton	optim
1.0719206	1.0727550	1.0735588
0.5686376	0.5695332	0.5695413
-0.4386992	-0.4381971	-0.4391370
-1.4484768	-1.4505377	-1.4501226
-1.9531086	-1.9558413	-1.9551820

- ▶ Approximately the same results!
- ▶ Strongly indicates that our implementations work.

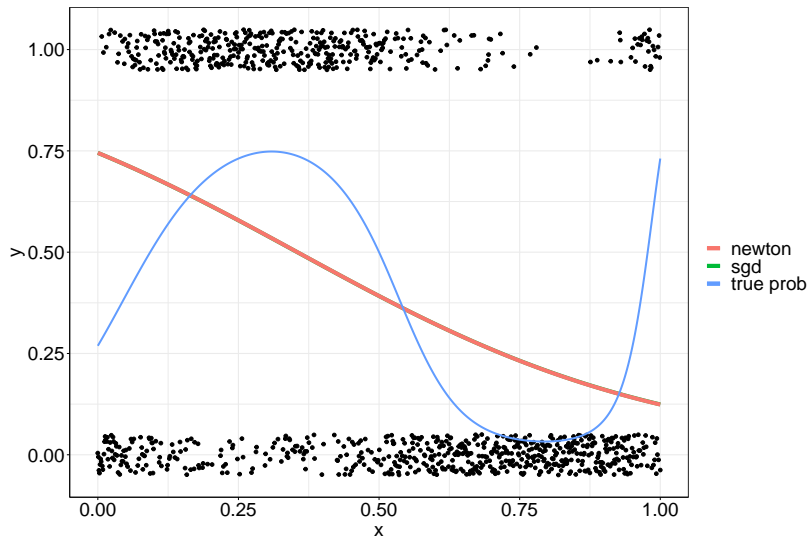
# $\lambda > 0$ : SGD convergence (200 iter.)



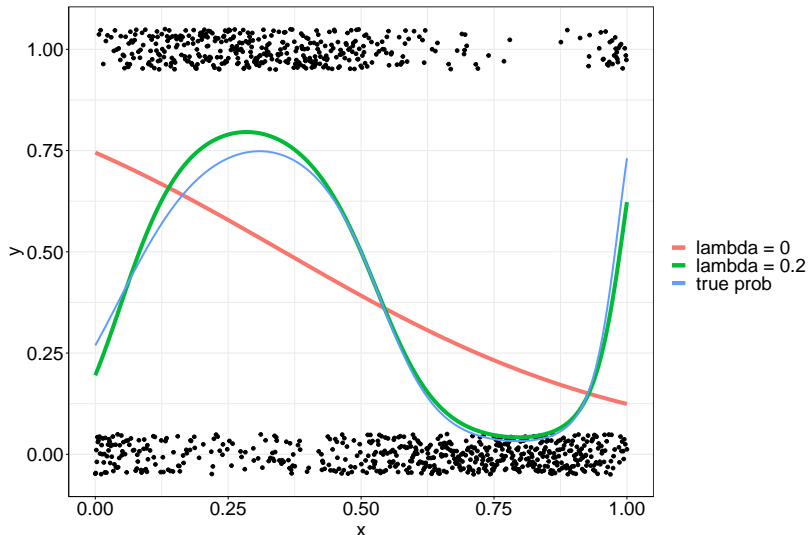
# $\lambda > 0$ : Newton convergence (4 iter.)



## $\lambda > 0$ : Fitted probabilities



## Comparing prob. fits with $\lambda = 0$ and $\lambda = 0.2$



- As expected,  $\lambda > 0$  gives a less volatile fit.
- But, in this case,  $\lambda = 0.2$  oversmooths.

**Conclusion:** Our `sgd` and `newton` implementations work for both  $\lambda = 0$  and  $\lambda > 0$ .

Can we improve the speed of the `sgd` implementation?

# Profiling: sgd and batch

```
sgd <- function(par, grad, n_obs, decay_schedule, epoch = batch,  
               n_iter = 100, sampler = sample, cb = NULL, ...)  
{  
  learning_rates <- decay_schedule(1:n_iter)  
  for(k in 1:n_iter) {  
    if(!is.null(cb))  
      cb()  
    epoch_sample <- sampler(n_obs)  
    par <- epoch(par, epoch_sample, learning_rates[k], grad, ...)  
  }  
  if (!is.null(cb))  
    cb()  
  par  
}
```

380

340

47830

```
batch <- function(par, epoch_sample, learning_rate, grad,  
                 batch_size = 50, ...)  
{  
  n_batches <- floor(length(epoch_sample) / batch_size)  
  for(j in 0:(n_batches - 1)) {  
    i <- epoch_sample[(j * batch_size + 1):  
                      (j * batch_size + batch_size)]  
    par <- par - learning_rate * grad(par, i, ...)  
  }  
  par  
}
```

20

100

1280

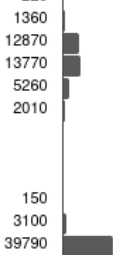
46380

► grad is the bottleneck.



# Profiling: grad\_H

```
get_grad_H <- function(lambda, x_vec, y_vec, knots, inner_knots)
{
  force_all_args()
  Phi_mat <- Phi(x_vec, knots)
  unpen_grad <- function(beta, i) {
    n_i <- length(i)
    Phi_i <- matrix(Phi_mat[i, ], nrow = n_i)
    y_p_diff <- y_vec[i] - p(beta, Phi_i)
    cp <- crossprod(Phi_i, y_p_diff)
    -cp / n_i
  }
  if (lambda > 0) {
    two_lambda_Omega <- 2 * lambda * Omega(inner_knots)
    function(beta, i) {
      two_lambda_Omega_beta <- two_lambda_Omega %*% beta
      unpen_grad(beta, i) + two_lambda_Omega_beta
    }
  } else
    unpen_grad
}
```



Line of Code	Time (ns)
force_all_args()	220
Phi_mat <- Phi(x_vec, knots)	1360
unpen_grad <- function(beta, i) {	12870
n_i <- length(i)	13770
Phi_i <- matrix(Phi_mat[i, ], nrow = n_i)	5260
y_p_diff <- y_vec[i] - p(beta, Phi_i)	2010
cp <- crossprod(Phi_i, y_p_diff)	
-cp / n_i	
}	
if (lambda > 0) {	150
two_lambda_Omega <- 2 * lambda * Omega(inner_knots)	3100
function(beta, i) {	
two_lambda_Omega_beta <- two_lambda_Omega %*% beta	
unpen_grad(beta, i) + two_lambda_Omega_beta	39790
}	
} else	
unpen_grad	
}	

- ▶ The indexing `Phi_mat[i, ]` is slow in R. Would be faster in C++.
- ▶ The rest are all vectorized operations in R and can't be improved much in C++.

## Profiling: p (not much to do about this one)

```
p <- function(beta, phi_mat)
{
  phi_beta_prod <- phi_mat %*% beta
  inv_logit(phi_beta_prod)
}
```

200

2980

7130

# What can we do?

- ▶ Reimplement the part of the gradient doing slicing in Rcpp, and as little else as possible, since everything is already vectorized.
- ▶ (Spoiler: It is actually a very good idea to implement everything in C++ so we can manage memory manually.)

## How to do it in an easy way?

- ▶ Use RcppArmadillo: Armadillo is a C++ interface to the BLAS and LAPACK linear algebra libraries also used by R.

```
1 #include <RcppArmadillo.h>  
2 // [[Rcpp::depends(RcppArmadillo)]]
```

```

1  arma::mat
2  crossprod(const arma::mat& A, const arma::mat& B)
3  {
4      return A.t() * B;
5  }
6
7  arma::colvec
8  inv_logit(const arma::colvec& x)
9  {
10     return exp(x) / (1 + exp(x));
11 }
12
13 arma::colvec
14 p(const arma::colvec& beta, const arma::mat& Phi_i)
15 {
16     return inv_logit(Phi_i * beta);
17 }
18
19 // [[Rcpp::export]]
20 arma::colvec
21 unpen_grad_cpp(const arma::colvec& beta, arma::uvec i,
22                double lambda, const arma::colvec& y,
23                const arma::mat& Phi)
24 {
25     i -= 1; // 0-indexing instead of 1-indexing.
26     return crossprod(Phi.rows(i), p(beta, Phi.rows(i)) - y.elem(i)) / i.n_elem;
27 }

```

# R wrapper for the C++ implementation

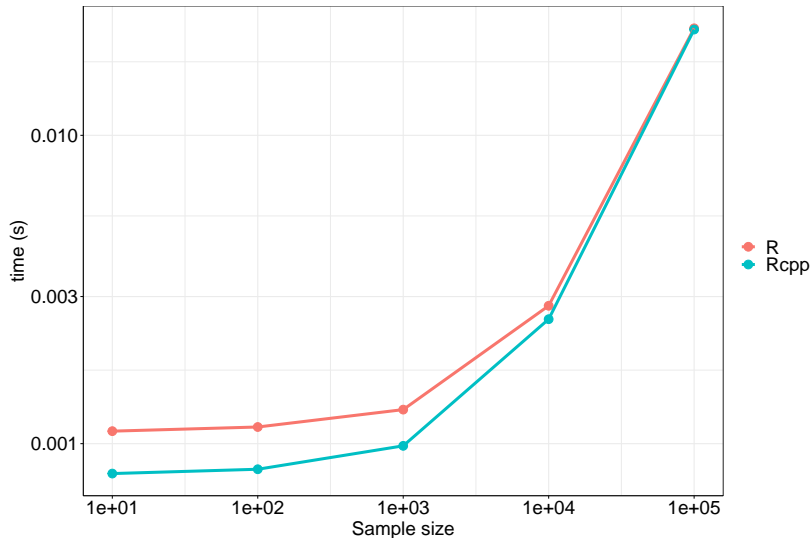
```
1  get_grad_H_cpp_wrap <- function(lambda, x_vec, y_vec, knots, inner_knots)
2  {
3      force_all_args()
4      Phi_mat <- Phi(x_vec, knots)
5      if (lambda > 0) {
6          two_lambda_Omega <- 2 * lambda * Omega(inner_knots)
7          function(beta, i)
8              unpen_grad_cpp(beta, i, lambda, y_vec, Phi_mat) +
9              two_lambda_Omega %*% beta
10     } else
11         function(beta, i)
12             unpen_grad_cpp(beta, i, lambda, y_vec, Phi_mat)
13 }
```

# Checking that the C++ gradient works

```
1  grad_H <- get_grad_H(lambda = 0.2, x, y, knots, inner_knots)
2  grad_H_cpp <- get_grad_H_cpp_wrap(lambda = 0.2, x, y, knots, inner_knots)
3  H <- get_H(lambda = 0.2, x, y, knots, inner_knots)
4
5  set.seed(2058)
6  beta_random <- rnorm(5, sd = 3)
7
8  all.equal(grad_H(beta_random, seq_along(x)),
9            grad_H_cpp(beta_random, seq_along(x)))
10
11  ## TRUE
12
13  max(abs(grad_H_cpp(beta_random, seq_along(x)) -
14          numDeriv::grad(H, beta_random)))
15
16  ## Approx. 4 * 10-9
```

numDeriv	R	C++
52.05423	52.05423	52.05423
-84.25675	-84.25675	-84.25675
-72.11850	-72.11850	-72.11850
324.56076	324.56076	324.56076
-220.28061	-220.28061	-220.28061

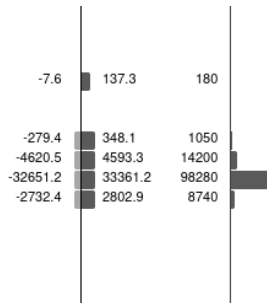
# Benchmarking the gradient implementations



- The Rcpp gradient is much faster for small sample sizes, but doesn't improve speed for sample sizes over  $10^6$ .

# Profiling `sgd` with Rcpp gradient for $10^6$ obs.

```
get_grad_H_cpp_wrap <- function(lambda, x_vec, y_vec, knots, innerknots)
{
  force_all_args()
  Phi_mat <- Phi(x_vec, knots)
  if (lambda > 0) {
    two_lambda_Omega <- 2 * lambda * Omega(inner_knots)
    function(beta, i) {
      two_lambda_Omega_beta <- two_lambda_Omega %*% beta
      unpen_grad <- unpen_grad_cpp(beta, i, lambda, y, Phi_mat)
      unpen_grad + two_lambda_Omega_beta
    }
  } else
    function(beta, i)
      unpen_grad_cpp(beta, i, lambda, y, Phi_mat)
}
```



- ▶ Gradient is still bottle neck.
- ▶ Very large memory allocation and deallocation for gradient!
- ▶ We must be able to do better with manual memory management in C++.



# C++: highest level

```
1 arma::colvec
2 sgd_cpp(arma::colvec par, const arma::colvec& learning_rates,
3         int n_iter, int batch_size, double lambda, const arma::colvec& y,
4         const arma::mat& Phi, const arma::mat& Omega)
5 {
6     int n_obs = y.n_elem;
7     arma::uvec epoch_order(n_obs);
8     for (int k = 0; k < n_iter; k++) {
9         epoch_order = arma::randperm(n_obs);
10        batch(par, learning_rates[k], batch_size, lambda,
11             n_obs, epoch_order, y, Phi, Omega);
12    }
13    return par;
14 }
```

- ▶ We only allocate memory for Phi, Omega, y, and the learning\_rates once, right when they are copied from R to C++.
- ▶ We pass them around by reference inside C++.

# C++: batch epoch

```
1 void
2 batch(arma::colvec& par, double learning_rate, int batch_size,
3       double lambda, int n_obs, arma::uvec& epoch_order,
4       const arma::colvec& y, const arma::mat& Phi, const arma::mat& Omega)
5 {
6     int lower_i = 0;
7     int upper_i = batch_size - 1;
8     int n_batches = floor(n_obs / batch_size);
9     for (int j = 0; j < n_batches; j++) {
10         par = par - learning_rate * grad_cpp(par,
11                                              epoch_order.subvec(lower_i, upper_i),
12                                              lambda, y, Phi, Omega);
13         lower_i += batch_size;
14         upper_i += batch_size;
15     }
16 }
```

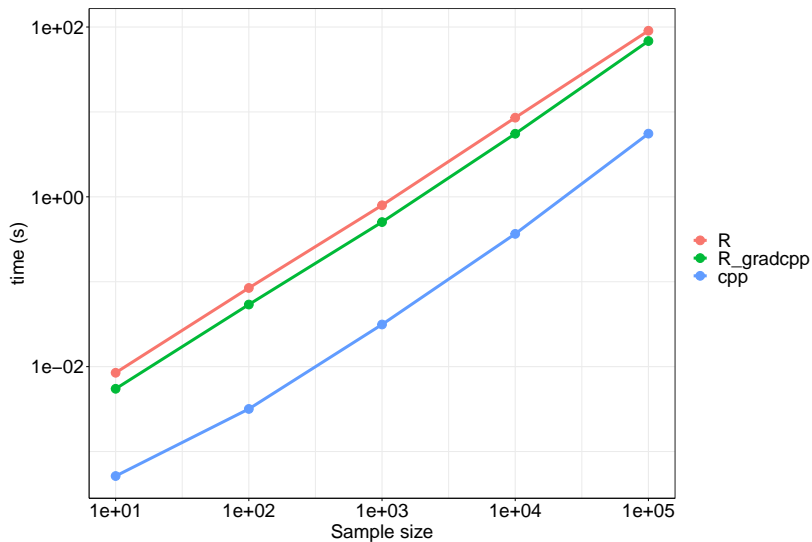
# C++: gradient

```
1  /* Version with 0-indexing for use in C++ */
2  arma::colvec
3  unpen_grad(const arma::colvec& beta, arma::uvec i,
4             double lambda, const arma::colvec& y,
5             const arma::mat& Phi)
6  {
7      return -crossprod(Phi.rows(i), p(beta, Phi.rows(i)) - y.elem(i)) / i.n_elem;
8  }
9
10 arma::colvec
11 grad_cpp(const arma::colvec& beta, arma::uvec i,
12          double lambda, const arma::colvec& y,
13          const arma::mat& Phi, const arma::mat& Omega)
14 {
15     if (lambda > 0)
16         return unpen_grad(beta, i, lambda, y, Phi) + 2 * lambda * Omega;
17     else
18         return unpen_grad(beta, i, lambda, y, Phi);
19 }
```

# C++: R wrapper

```
1  sgd_cpp_wrap <- function(par, x, y, decay_schedule, inner_knots,  
2                           lambda = 0, n_iter = 100, batch_size = 1)  
3  {  
4      learning_rates <- decay_schedule(1:n_iter)  
5      Om <- Omega(inner_knots)  
6      Phi_mat <- Phi(x, get_knots(inner_knots))  
7      sgd_cpp(par, learning_rates, n_iter, batch_size, lambda, y, Phi_mat, Om)  
8  }
```

# Benchmarking full implementations



► Staying in C++ pays off!

# Should you use stochastic gradient descent?

## Pros

- ▶ Can converge faster when the sample size is very large (by taking many small fast steps, instead of large expensive steps).
- ▶ Can be used in online learning settings, where the estimate is updated as more data becomes available.

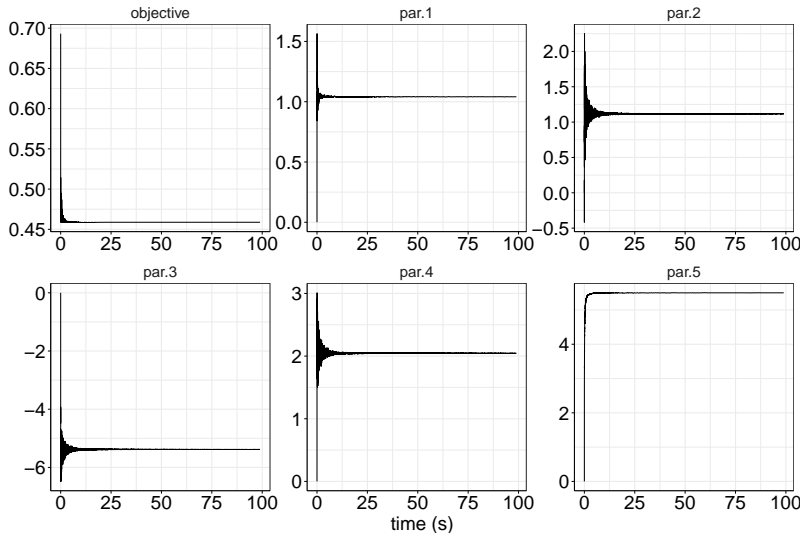
## Cons

- ▶ Time consuming to tune decay schedule parameters.
- ▶ Can still get trapped in local optima (alternative for functions with many local optima: simulated annealing, e.g., in the R package GenSA).

# Horse data: Fitting

```
1  ## Load data and scale and center x for faster convergence
2  x <- scale(horses$Temperature)
3  y <- horses$dead
4  n_beta <- 5
5  beta_init <- rep(0, n_beta)
6  inner_knots <- seq(min(x), max(x), length.out = n_beta - 2)
7  knots <- get_knots(inner_knots)
8
9  ## SGD
10 grad_H <- get_grad_H(lambda = 0, x, y, knots, inner_knots)
11 H <- get_H(lambda = 0, x, y, knots, inner_knots)
12 sgd_batch_tracer <- tracer(c("objective", "par"),
13                           expr = quote(objective <- H(par)),
14                           N = 0)
15 beta_fit_horse <- sgd(par = beta_init,
16                      grad = grad_H,
17                      n_obs = length(x),
18                      decay_schedule = decay_scheduler(gamma0 = 1,
19                                                       gamma1 = 0.3,
20                                                       n1 = 125),
21                      epoch = batch,
22                      n_iter = 20000,
23                      batch_size = 1,
24                      cb = sgd_batch_tracer$tracer)
25
26 ## Newton
27 newton_tracer <- tracer(c("objective", "par_new"),
28                       expr = quote(objective <- H(par_new)))
29
30 newton_fit_horse <- newton(init_guess = beta_init,
31                          x = x,
32                          y = y,
33                          max_iter = 100,
34                          cb = newton_tracer$tracer)
```

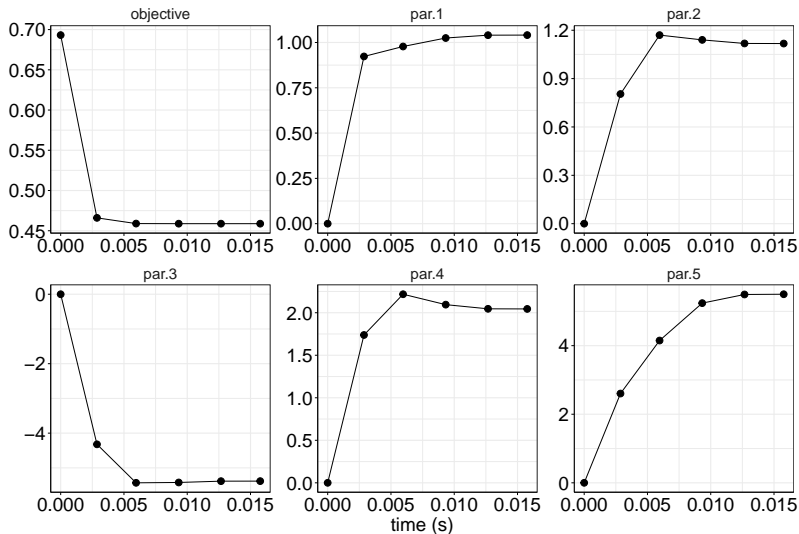
# Horse data: SGD convergence (20000 iter.)



- Hard to find a better decay schedule that takes large enough steps, but without fluctuating up and down.



# Horse data: Newton convergence (5 iter.)

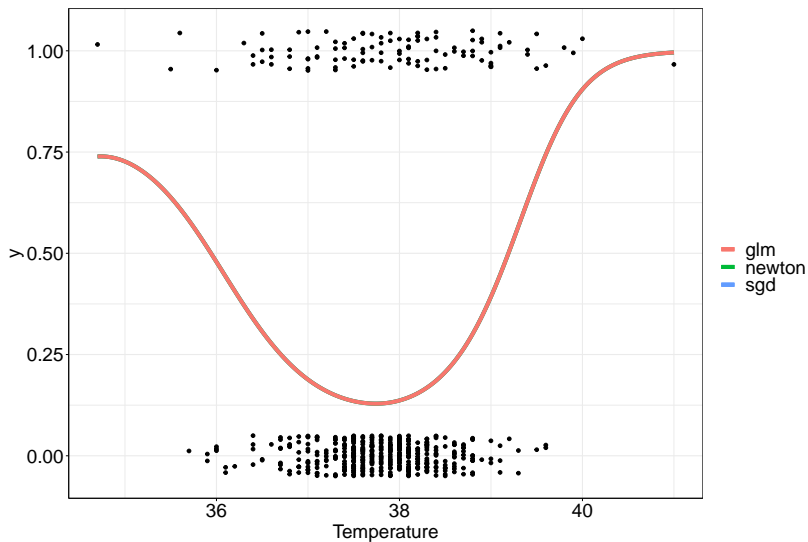


## Horse data: Estimates

sgd	newton	glm
1.040949	1.040815	1.040816
1.117849	1.117413	1.117412
-5.380647	-5.381459	-5.381458
2.045161	2.044704	2.044703
5.499648	5.499755	5.499760

- Approximately the same.

# Horse data: Fitted probabilities



► Approximately the same.

## Extra slides

- ▶ Gradient derivation.
- ▶ Comparing profiling of R and Rcpp.
- ▶ Momentum.
- ▶ Adam.
- ▶ Old slides for convergence and fits.
- ▶ `profvis` on `sgd_cpp`

## Rcpp implementation

```
batch <- function(par, epoch_sample, learning_rate, grad,  
                  batch_size = 50, ...)  
{  
  n_batches <- floor(length(epoch_sample) / batch_size)  
  for(j in 0:(n_batches - 1)) {  
    i <- epoch_sample[(j * batch_size + 1):  
                      (j * batch_size + batch_size)]  
    par <- par - learning_rate * grad(par, i, ...)  
  }  
  par  
}
```

30  
100  
1380  
27290

## R implementation

```
batch <- function(par, epoch_sample, learning_rate, grad,  
                  batch_size = 50, ...)  
{  
  n_batches <- floor(length(epoch_sample) / batch_size)  
  for(j in 0:(n_batches - 1)) {  
    i <- epoch_sample[(j * batch_size + 1):  
                      (j * batch_size + batch_size)]  
    par <- par - learning_rate * grad(par, i, ...)  
  }  
  par  
}
```

20  
100  
1280  
46380

# Rcpp implementation

```
get_grad_H.cpp.wrap <- function(lambda, x_vec, y_vec, knots, innerknots)
{
  force_all_args()
  Phi_mat <- Phi(x_vec, knots)
  if (lambda > 0) {
    two_lambda_Omega <- 2 * lambda * Omega(inner_knots)
    function(beta, i) {
      two_lambda_Omega_beta <- two_lambda_Omega %*% beta
      unpen_grad <- unpen_grad_cpp(beta, i, lambda, y, Phi_mat)
      unpen_grad + two_lambda_Omega_beta
    }
  } else
    function(beta, i)
      unpen_grad_cpp(beta, i, lambda, y, Phi_mat)
}
```

# R implementation

```
get_grad_H <- function(lambda, x_vec, y_vec, knots, inner_knots)
{
  force_all_args()
  Phi_mat <- Phi(x_vec, knots)
  unpen_grad <- function(beta, i) {
    n_i <- length(i)
    Phi_i <- matrix(Phi_mat[i, ], nrow = n_i)
    y_p.diff <- y_vec[i] - p(beta, Phi_i)
    cp <- crossprod(Phi_i, y_p.diff)
    -cp / n_i
  }
  if (lambda > 0) {
    two_lambda_Omega <- 2 * lambda * Omega(inner_knots)
    function(beta, i) {
      two_lambda_Omega_beta <- two_lambda_Omega %*% beta
      unpen_grad(beta, i) + two_lambda_Omega_beta
    }
  } else
    unpen_grad
}
```

# Momentum: keep memory of last direction

With memory  $\beta \in [0, 1)$ , use update direction

$$\rho_n = \beta \rho_{n-1} + (1 - \beta) \frac{1}{M} \sum_{i \in I_n} \nabla_{\theta} L(x_i, \theta_n)$$

to get

$$\theta_n = \theta_{n-1} - \gamma_n \rho_n.$$

```
1 momentum <- function() {  
2   rho <- 0  
3   function(par, epoch_sample, learning_rate, grad,  
4     batch_size = 50, mom_memory = 0.95, ...)  
5   {  
6     M <- floor(length(epoch_sample) / batch_size)  
7     for(j in 0:(M - 1)) {  
8       i <- epoch_sample[(j * batch_size + 1):  
9         (j * batch_size + batch_size)]  
10      rho <-<= mom_memory * rho + (1 - mom_memory) * grad(par, i, ...)  
11      par <-<= par - learning_rate * rho  
12    }  
13    par  
14  }  
15 }
```

# Adam: rescale gradient components to avoid narrow valley.

- ▶ Keeps memory of second moment for rescaling.
- ▶ “invariant to diagonal rescaling of the gradients” (Kingma and Ba, 2015).

```
1 adam <- function() {  
2   rho <- v <- 0  
3   function(par, epoch_sample, learning_rate, grad,  
4     batch_size = 50, mom_memory = 0.9, mom2_memory = 0.9, ...)  
5   {  
6     M <- floor(length(epoch_sample) / batch_size)  
7  
8     for(j in 0:(M - 1)) {  
9       i <- epoch_sample[(j * m + 1):(j * m + m)]  
10      gr <- grad(par, i, ...)  
11      rho <-< mom_memory * rho + (1 - mom_memory) * gr  
12      v <-< mom2_memory * v + (1 - mom2_memory) * gr^2  
13      par <- par - learning_rate * (rho / (sqrt(v) + 1e-8))  
14    }  
15    par  
16  }  
17 }
```



## Rcpp implementation

```
sgd <- function(par, grad, n_obs, decay_schedule, epoch = batch,  
               n_iter = 100, sampler = sample, cb = NULL, ...)  
{  
  learning_rates <- decay_schedule(1:n_iter)  
  for(k in 1:n_iter) {  
    if(!is.null(cb))  
      cb()  
    epoch_sample <- sampler(n_obs)  
    par <- epoch(par, epoch_sample, learning_rates[k], grad, ...)  
  }  
  if (!is.null(cb))  
    cb()  
  par  
}
```

10

290

28840

## R implementation

```
sgd <- function(par, grad, n_obs, decay_schedule, epoch = batch,  
               n_iter = 100, sampler = sample, cb = NULL, ...)  
{  
  learning_rates <- decay_schedule(1:n_iter)  
  for(k in 1:n_iter) {  
    if(!is.null(cb))  
      cb()  
    epoch_sample <- sampler(n_obs)  
    par <- epoch(par, epoch_sample, learning_rates[k], grad, ...)  
  }  
  if (!is.null(cb))  
    cb()  
  par  
}
```

380

340

47830

# Old slides

Slides using old version of gradient.

# Getting $H$ and $D_\beta H$ with $\lambda = 0$

- We set  $\lambda = 0$ , corresponding to usual unpenalized logistic regression, so we can compare results with `glm`.

```
1 grad_H <- get_grad_H(lambda = 0, x, y, knots, inner_knots)
2 H <- get_H(lambda = 0, x, y, knots, inner_knots)
```

- We confirm that our gradient implementation works by comparing to a numerical approximation.

```
1 all.equal(numDeriv::grad(H, beta),
2           as.vector(grad_H(beta, 1:length(x))), tolerance = 1e-8)
3
4 ## TRUE
```

# Run methods

```
1  set.seed(2021)
2
3  ## Stochastic gradient descent
4  sgd_batch_tracer <- tracer(c("objective", "par"),
5                             expr = quote(objective <- H(par)),
6                             N = 0)
7  ds <- decay_scheduler(gamma0 = 0.8, gamma1 = 0.008, n1 = 90)
8  beta_fit_1 <- sgd(par = beta_init,
9                    grad = grad_H,
10                     n_obs = n_sim,
11                     decay_schedule = ds,
12                     epoch = batch,
13                     n_iter = 200,
14                     batch_size = 1,
15                     cb = sgd_batch_tracer$tracer)
16
17  ## Newton
18  newton_tracer <- tracer(c("objective", "par_new"),
19                          expr = quote(objective <- H(par_new)))
20  newton_fit <- newton(init_guess = beta_init,
21                      x = x,
22                      y = y,
23                      max_iter = 100,
24                      cb = newton_tracer$tracer)
```

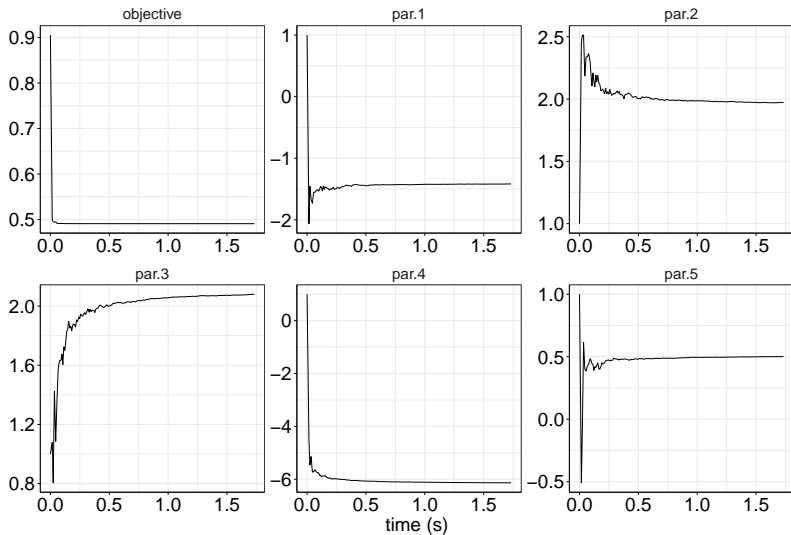
# Comparing results to glm (IWLS)

```
1 glm(y ~ splineDesign(knots, x) - 1, family = binomial)$coefficients
```

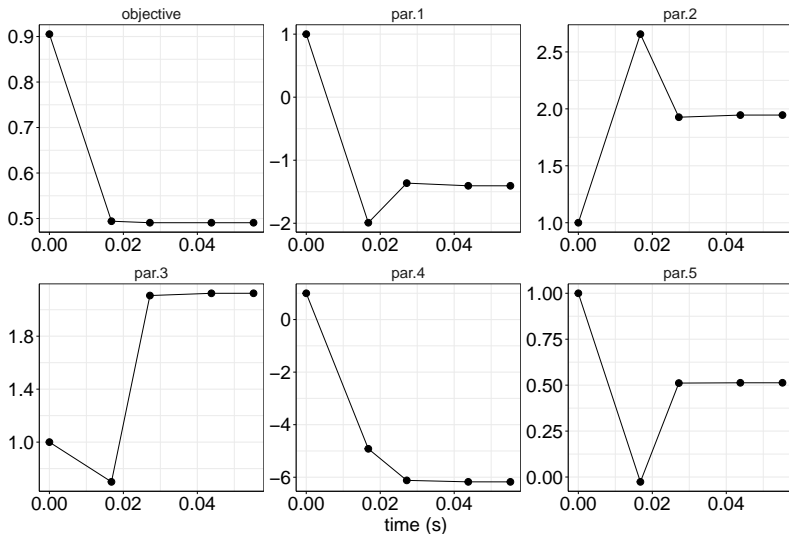
sgd	newton	glm
-1.4172427	-1.4058146	-1.4058146
1.9725305	1.9454677	1.9454677
2.0802703	2.1241283	2.1241283
-6.1350949	-6.1761810	-6.1761811
0.5009762	0.5129334	0.5129334

The methods yield approximately the same results, indicating that our implementations work with  $\lambda = 0$ .

# SGD convergence (200 iterations)

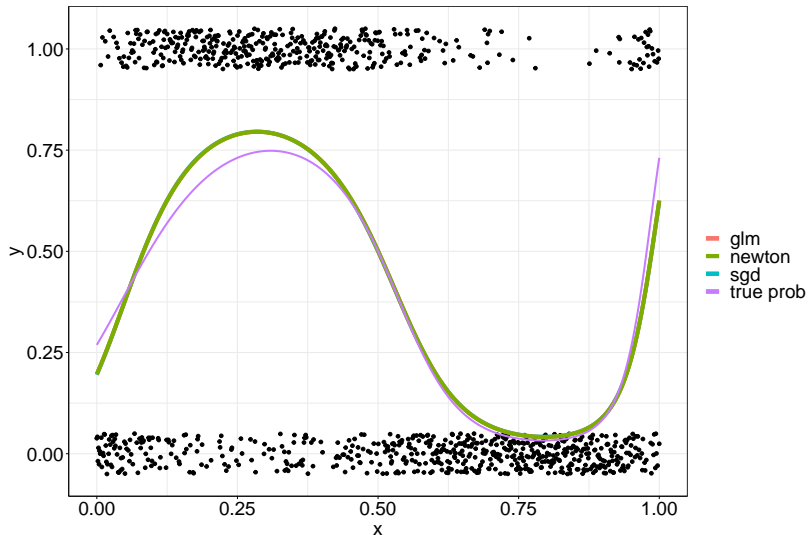


# Newton convergence (4 iterations)



► Okay convergence two steps from initial guess!

# Predictions (the same for all three methods)





# Profiling `sgd` and `batch`

```
sgd <- function(par, grad, n_obs, decay_schedule, epoch = batch,  
               n_iter = 100, sampler = sample, cb = NULL, ...)
```

```
  learning_rates <- decay_schedule(1:n_iter)
```

```
  for(k in 1:n_iter) {
```

```
    if(!is.null(cb))
```

```
      cb()
```

```
    epoch_sample <- sampler(n_obs)
```

```
    par <- epoch(par, epoch_sample, learning_rates[k], grad, ...)
```

```
  }
```

```
  if (!is.null(cb))
```

```
    cb()
```

```
  par
```

```
batch <- function(par, epoch_sample, learning_rate, grad,  
                 batch_size = 50, ...)
```

```
  n_batches <- floor(length(epoch_sample) / batch_size)
```

```
  for(j in 0:(n_batches - 1)) {
```

```
    i <- epoch_sample[(j * batch_size + 1):
```

```
                     (j * batch_size + batch_size)]
```

```
    par <- par - learning_rate * grad(par, i, ...)
```

```
  }
```

```
  par
```

330

330

43240

10

100

1420

41680

► `grad` is the bottleneck.

# Profiling grad\_H and p

```
p <- function(beta, phi_mat)                                220
{
  phi_beta_prod <- phi_mat %*% beta                        2720
  inv_logit(phi_beta_prod)                                6330
}

get_grad_H <- function(lambda, x_vec, y_vec, knots, inner_knots)
{
  force_all_args()
  Phi_mat <- Phi(x_vec, knots)
  Omega_mat <- Omega(inner_knots)
  function(beta, i) {                                       160
    Phi_i <- matrix(Phi_mat[i,], nrow = length(i))          13090
    Omega_beta <- Omega_mat %*% beta                          3220
    cp <- crossprod(Phi_i, y_vec[i] - p(beta, Phi_i))        16630
    -cp / length(x_vec[i]) + 2 * lambda * Omega_beta        6180
  }
}
```

- ▶ The indexing `Phi_mat[i, ]` is slow in R. Would be faster in C++.
- ▶ The rest are all vectorized operations in R, so they cannot be improved much in C++.

## $\lambda > 0$ : Run methods

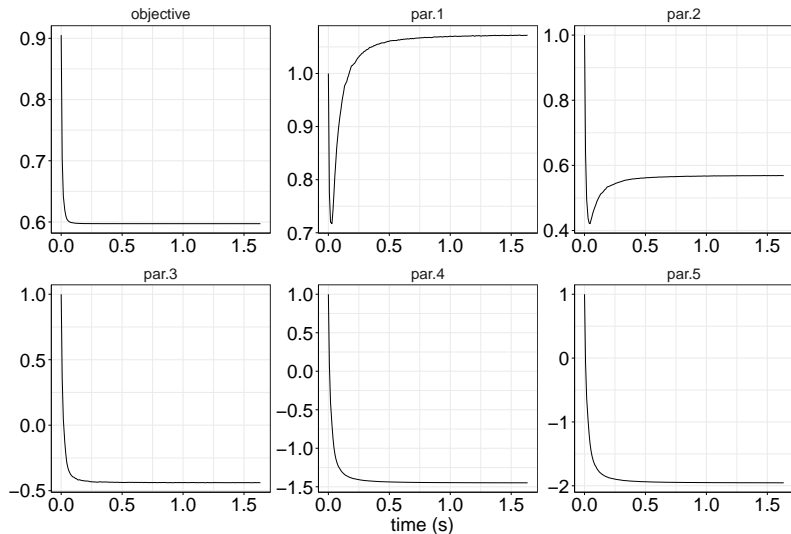
```
1  ## Stochastic gradient descent
2  grad_H <- get_grad_H(lambda = 0.2, x, y, knots, inner_knots)
3  H <- get_H(lambda = 0.2, x, y, knots, inner_knots)
4  sgd_batch_tracer <- tracer(c("objective", "par"),
5                             expr = quote(objective <- H(par)),
6                             N = 0)
7  beta_fit_lam <- sgd(par = beta_init,
8                     grad = grad_H,
9                     n_obs = n_sim,
10                    decay_schedule = decay_scheduler(gamma0 = 0.015,
11                                                    gamma1 = 0.001,
12                                                    n1 = 90),
13                    epoch = batch,
14                    n_iter = 200,
15                    batch_size = 1,
16                    cb = sgd_batch_tracer$tracer)
17
18  ## Newton
19  newton_tracer <- tracer(c("objective", "par_new"),
20                        expr = quote(objective <- H(par_new)))
21  newton_fit_lam <- newton(init_guess = beta_init, x = x, y = y,
22                        max_iter = 100,
23                        lambda = 0.2,
24                        cb = newton_tracer$tracer)
```

## $\lambda > 0$ : Comparing $\beta$ estimates

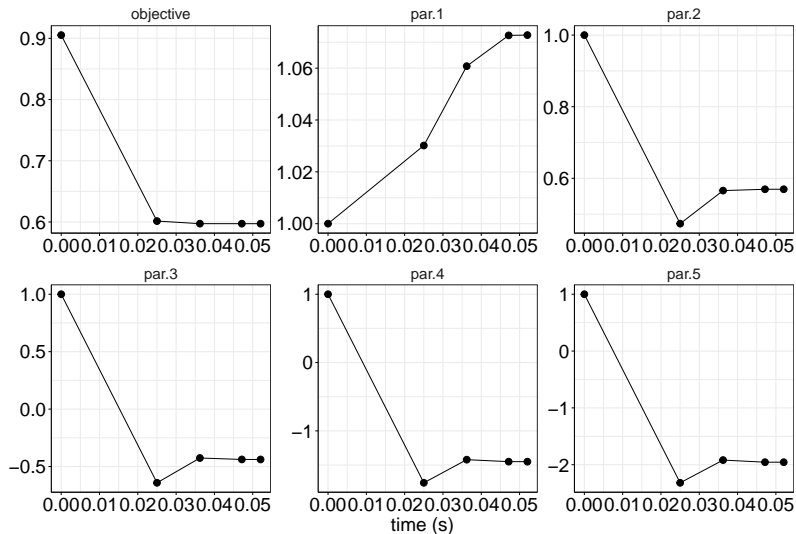
sgd	newton
1.0717191	1.0727550
0.5687178	0.5695332
-0.4380832	-0.4381971
-1.4485311	-1.4505377
-1.9531791	-1.9558413

- ▶ Approximately the same results!
- ▶ Indicates that there are no bugs in the implementations.

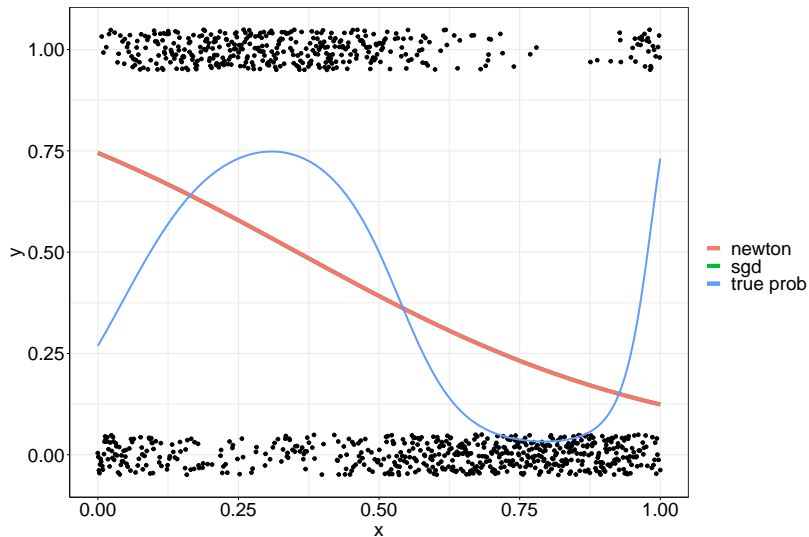
# $\lambda > 0$ : SGD convergence (200 iter.)



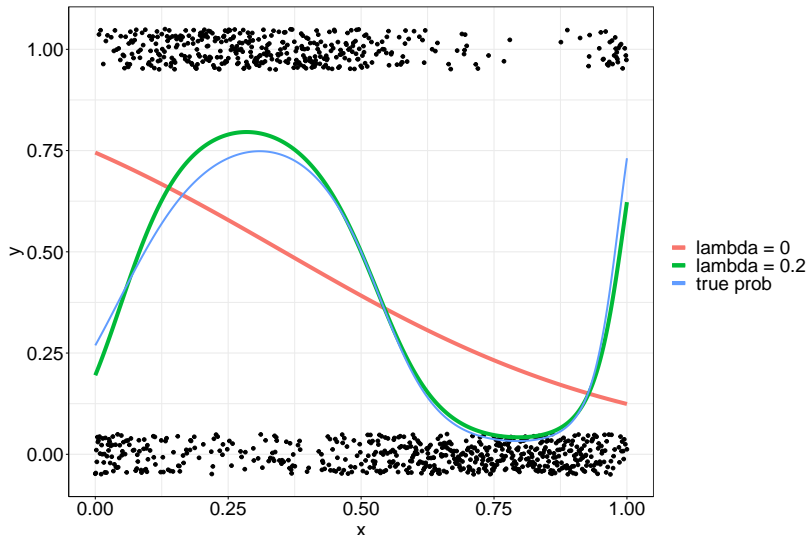
# $\lambda > 0$ : Newton convergence (4 iter.)



## $\lambda > 0$ : Fitted probabilities



## Comparing prob. fits with $\lambda = 0$ and $\lambda = 0.2$



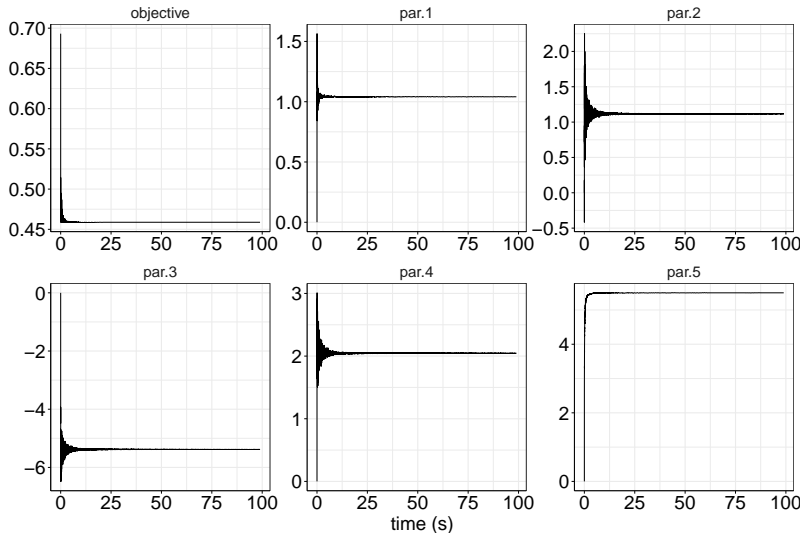
- As expected,  $\lambda > 0$  gives a less volatile fit.
- But, in this case,  $\lambda = 0.2$  oversmooths.



# Horse data: Fitting

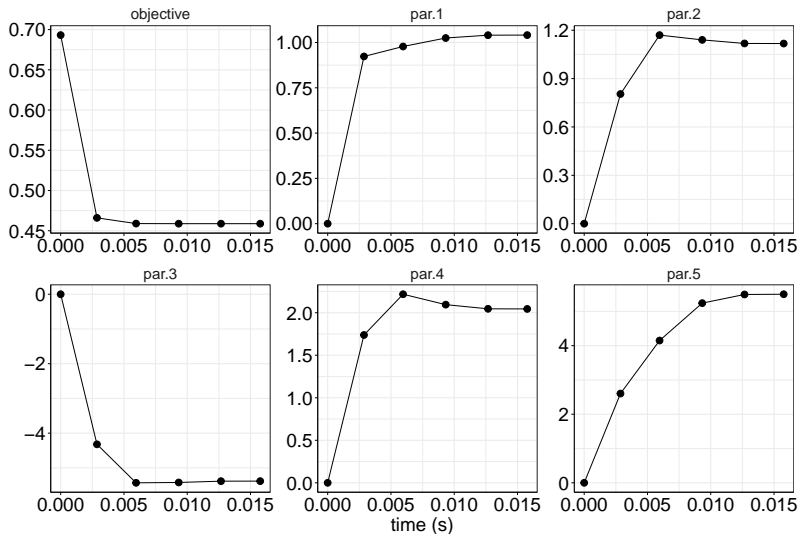
```
1  ## Load data and scale and center x for faster convergence
2  x <- scale(horses$Temperature)
3  y <- horses$dead
4  n_beta <- 5
5  beta_init <- rep(0, n_beta)
6  inner_knots <- seq(min(x), max(x), length.out = n_beta - 2)
7  knots <- get_knots(inner_knots)
8
9  ## SGD
10 grad_H <- get_grad_H(lambda = 0, x, y, knots, inner_knots)
11 H <- get_H(lambda = 0, x, y, knots, inner_knots)
12 sgd_batch_tracer <- tracer(c("objective", "par"),
13                             expr = quote(objective <- H(par)),
14                             N = 0)
15 beta_fit_horse <- sgd(par = beta_init,
16                       grad = grad_H,
17                       n_obs = length(x),
18                       decay_schedule = decay_scheduler(gamma0 = 1,
19                                                         gamma1 = 0.3,
20                                                         n1 = 125),
21                       epoch = batch,
22                       n_iter = 20000,
23                       batch_size = 1,
24                       cb = sgd_batch_tracer$tracer)
25
26 ## Newton
27 newton_tracer <- tracer(c("objective", "par_new"),
28                         expr = quote(objective <- H(par_new)))
29
30 newton_fit_horse <- newton(init_guess = beta_init,
31                           x = x,
32                           y = y,
33                           max_iter = 100,
34                           cb = newton_tracer$tracer)
```

# Horse data: SGD convergence (20000 iter.)



- Hard to find a better decay schedule that takes large enough steps, but without fluctuating up and down.

# Horse data: Newton convergence (5 iter.)

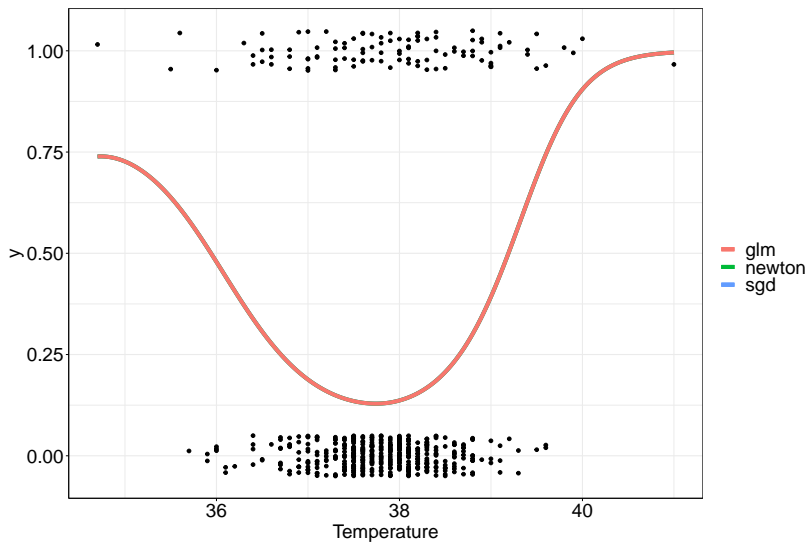


## Horse data: Estimates

sgd	newton	glm
1.040949	1.040815	1.040816
1.117849	1.117413	1.117412
-5.380647	-5.381459	-5.381458
2.045161	2.044704	2.044703
5.499648	5.499755	5.499760

- Approximately the same.

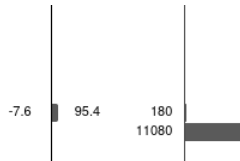
# Horse data: Fitted probabilities



► Approximately the same.

# C++ profiling

```
sgd_cpp_wrap <- function(par, x, y, decay_schedule, inner_knots,  
                          lambda = 0, n_iter = 100, batch_size = 1)  
{  
  learning_rates <- decay_schedule(1:n_iter)  
  Om <- Omega(inner_knots)  
  Phi_mat <- Phi(x, get_knots(inner_knots))  
  sgd_cpp(par, learning_rates, n_iter, batch_size, lambda, y, Phi_mat, Om)  
}
```



- ▶ profvis cannot check memory allocation and deallocation done in C++.
- ▶ But we know that we only allocate memory for all large objects once, so we must do better.