

# Optimisation Methods I

Matteo Fasiolo

# Optimisation Problems

Many problems in statistics can be characterised as solving

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} f(\theta), \quad (1)$$

where  $f$  is a real scalar valued function of vector  $\theta$ , usually called the **objective function**. For example,  $f$  might be:

- ▶ the negative log-likelihood of a statistical model;
- ▶ the (negative) Bayesian posterior density of  $\theta$ ;
- ▶ a dissimilarity measure for the alignment of two DNA sequences, where the alignment is controlled by  $\theta$ ;
- ▶ the total distance travelled, and  $\theta$  defines the order of delivery drop off points.

Note that

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} f(\theta) = \underset{\theta}{\operatorname{argmax}} -f(\theta),$$

so we are not losing anything by concentrating on minimisation.

In the optimization literature, the convention is to minimise the objective, which is interpreted as a **cost** or **penalty** function.

Our setup and assumptions:

- ▶  $\theta \in \Theta$ , the set of all possible values
- ▶ we can evaluate  $f(\theta)$  for all elements of  $\theta \in \Theta$
- ▶  $f$  is “well-behaved”, e.g.  $\theta_1 \approx \theta_2$  implies  $f(\theta_1) \approx f(\theta_2)$

To further restrict our focus, we will assume  $\Theta$  is continuous.

So excluding discrete optimisation problems, why?

There is a big divide between continuous and discrete problems.

We can not cover both in two lectures.

Continuous problem are:

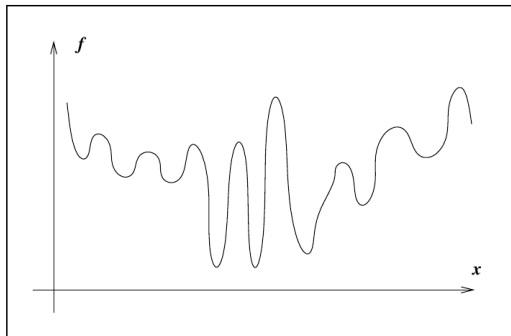
- ▶ generally easier, discrete problems can cost  $> O(n^k)$  for any  $k$
- ▶ more commonly encountered in Statistics
- ▶ often solvable via widely-used optimisation methods

We will consider unconstrained problems where  $\Theta = \mathbb{R}^p$ .

If  $\theta$  constrained:

- ▶ write  $\theta = g(\theta_u)$  (e.g.,  $\theta = \exp(\theta_u)$ );
- ▶  $\hat{\theta}_u = \operatorname{argmin}_{\theta_u} f(g(\theta_u))$

We will consider **local** optimisation methods.



Local = global if  $f(\theta)$  convex.

## A Simple Model Leading to Troubles

Consider the logistic map, that is

$$n_{t+1} = rn_t(1 - n_t/K), \quad t = 0, 1, 2, \dots$$

where

- ▶  $n_t$  is the population at time  $t$  (assume  $n_0$  is known)
- ▶  $r$  the growth rate parameter
- ▶  $K$  is the carrying capacity (if  $n_t = K$  then  $n_{t+1} = 0$ )

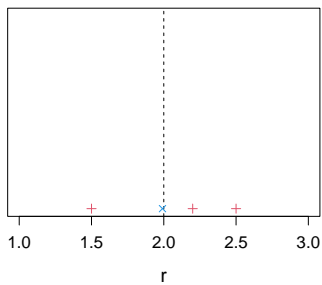
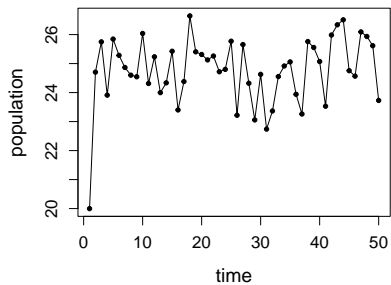
Suppose we observe  $y_t = n_t + \epsilon_t$  where  $\epsilon_t \sim N(0, \sigma^2)$ .

We want to estimate  $r$  and  $K$  by minimising

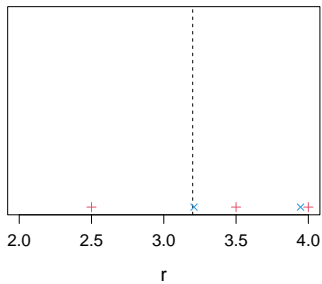
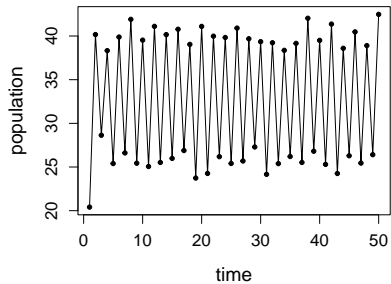
$$f(r, K) = \sum_{t=1}^T (y_t - n_t)^2$$

where  $n_t = n_t(r, K, n_{t-1})$ .

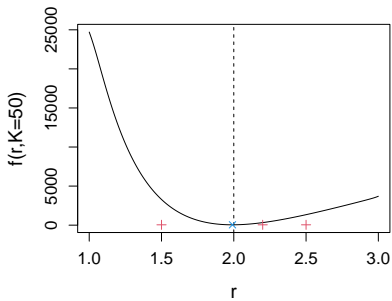
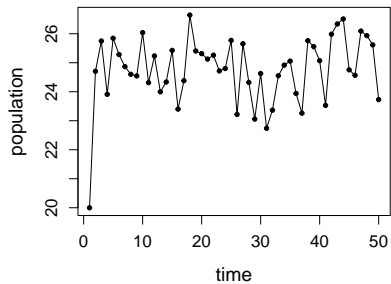
### Simulation with $r = 2$ and $k = 50$



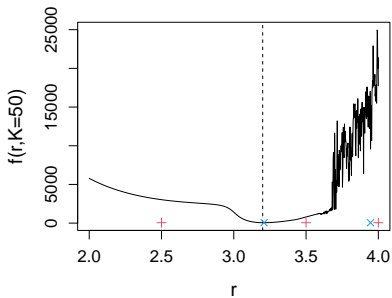
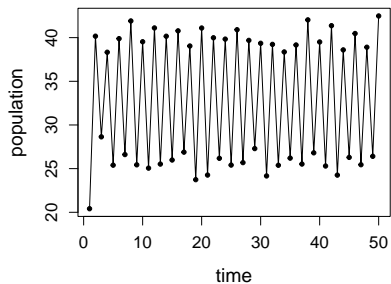
### Simulation with $r = 3.2$ and $k = 50$



## Simulation with $r = 2$ and $k = 50$

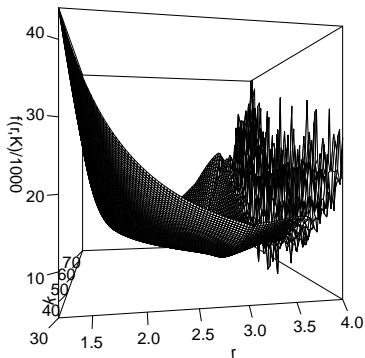
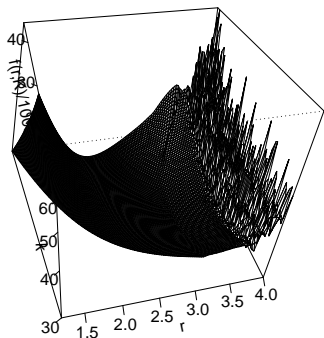


## Simulation with $r = 3.2$ and $k = 50$





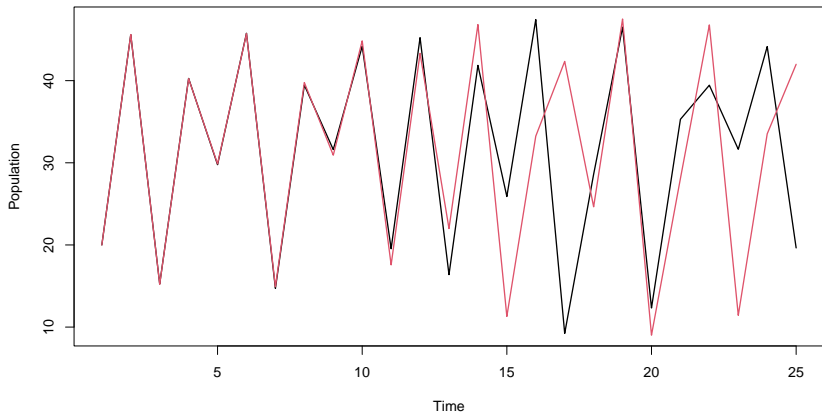
Likelihood with simulations from  $r = 3.2$  and  $k = 50$ :



So:

- ▶ simple models can lead to hard optimisation problems;
- ▶ when the optimiser does not converge, look at the objective;
- ▶ often possible to do this even when  $d > 2$  (see tomorrow).

## Digression: What's wrong with this model?



Population  $n_t$  when  $r = 3.8$  (black) and when  $r = 3.8001$  (red).

Likelihood might not be a good objective when fitting chaotic models, see Fasiolo et al (2016).

## Back to Local Optimisation

Given a guess  $\theta^{[0]}$  we consider methods that, at  $k$ -th iteration:

1. evaluate  $f(\theta^{[k]})$ , and possibly  $\nabla f(\theta)$  and  $\nabla^2 f(\theta)$  at  $\theta^{[k]}$ .
2. Use the information from step 1 to:
  - a. Find a **search direction**,  $\Delta$ , s.t.  $f(\theta^{[k]} + \alpha\Delta) < f(\theta^{[k]})$  for some  $\alpha > 0$ ;
  - b. Find  $\alpha$  s.t.  $f(\theta^{[k]} + \alpha\Delta)$  is sufficiently lower than  $f(\theta^{[k]})$  (**sufficient decrease condition**);
  - c. Set  $\theta^{[k+1]} = \theta^{[k]} + \alpha\Delta$ .
3. If a minimum has yet been reached, stop, otherwise back to 1.

Note: in 2.a. we are trying to find an “optimal” direction, not any descent direction.

How to find  $\mathbf{\Delta}$ ?

Assuming that  $f$  is twice differentiable, we can use Taylor expansion

$$f(\boldsymbol{\theta} + \mathbf{\Delta}) = f(\boldsymbol{\theta}) + \nabla f(\boldsymbol{\theta})^T \mathbf{\Delta} + \frac{1}{2} \mathbf{\Delta}^T \nabla^2 f(\boldsymbol{\theta}) \mathbf{\Delta} + R_{\boldsymbol{\theta}}^2(\mathbf{\Delta}),$$

where  $R_{\boldsymbol{\theta}}^2(\mathbf{\Delta})$  is the remainder of the expansion and

$$\nabla f(\boldsymbol{\theta}) = \begin{pmatrix} \frac{\partial f}{\partial \theta_1} \\ \frac{\partial f}{\partial \theta_2} \\ \vdots \\ \vdots \end{pmatrix} \quad \text{and} \quad \nabla^2 f(\boldsymbol{\theta}) = \begin{pmatrix} \frac{\partial^2 f}{\partial \theta_1^2} & \frac{\partial^2 f}{\partial \theta_1 \partial \theta_2} & \cdot & \cdot \\ \frac{\partial^2 f}{\partial \theta_2 \partial \theta_1} & \frac{\partial^2 f}{\partial \theta_2^2} & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \end{pmatrix}.$$

Note that

$$f(\boldsymbol{\theta} + \boldsymbol{\Delta}) = f(\boldsymbol{\theta}) + \nabla f(\boldsymbol{\theta})^T \boldsymbol{\Delta} + \frac{1}{2} \boldsymbol{\Delta}^T \nabla^2 f(\boldsymbol{\theta}) \boldsymbol{\Delta} + R_{\boldsymbol{\theta}}^2(\boldsymbol{\Delta}),$$

can help us identify a minimum when we find it.

Consider candidate  $\hat{\boldsymbol{\theta}}$ , then for  $\|\boldsymbol{\Delta}\| \ll 1$  we have

$$f(\hat{\boldsymbol{\theta}} + \boldsymbol{\Delta}) \simeq f(\hat{\boldsymbol{\theta}}) + \nabla f(\hat{\boldsymbol{\theta}})^T \boldsymbol{\Delta} + \frac{1}{2} \boldsymbol{\Delta}^T \nabla^2 f(\hat{\boldsymbol{\theta}}) \boldsymbol{\Delta},$$

and if  $\nabla f(\hat{\boldsymbol{\theta}}) = \mathbf{0}$  and  $\nabla^2 f(\hat{\boldsymbol{\theta}})$  is positive definite we have

$$f(\hat{\boldsymbol{\theta}} + \boldsymbol{\Delta}) \simeq f(\hat{\boldsymbol{\theta}}) + \frac{1}{2} \boldsymbol{\Delta}^T \nabla^2 f(\hat{\boldsymbol{\theta}}) \boldsymbol{\Delta} > f(\hat{\boldsymbol{\theta}}) \quad \text{for any small } \boldsymbol{\Delta}.$$

Now back to search for  $\boldsymbol{\Delta}$  based on Taylor expansion.

## Steepest Descent (aka Gradient Descent)

SD uses a first order Taylor expansion around current  $\theta$

$$f(\theta + \Delta) \simeq f(\theta) + \nabla f(\theta)^T \Delta.$$

Under local model, along which direction  $\Delta$  should we move?

Note that local model does not provide any info about  $\|\Delta\|$ .

Let  $p = \dim(\theta)$  then if:

- ▶  $p = 1$  if local model is a line
- ▶  $p = 2$  a plane
- ▶  $p \geq 3$  a hyper-plane

Let's separate the magnitude and the length of the step.

So, let's assume  $\|\mathbf{\Delta}\| = 1$  and that the local model is

$$f(\boldsymbol{\theta} + \alpha\mathbf{\Delta}) \simeq f(\boldsymbol{\theta}) + \alpha\nabla f(\boldsymbol{\theta})^T \mathbf{\Delta},$$

for some step-length  $\alpha > 0$ .

Then

$$f(\boldsymbol{\theta} + \alpha\mathbf{\Delta}) \simeq f(\boldsymbol{\theta}) + \alpha\|\nabla f(\boldsymbol{\theta})\|\cos(\phi),$$

where  $\phi$  is the angle between  $\nabla f(\boldsymbol{\theta})$  and  $\mathbf{\Delta}$ .

For any fixed  $\alpha$ ,  $f(\boldsymbol{\theta} + \alpha\mathbf{\Delta})$  is minimal when  $\phi = \pi$ , so

$$\mathbf{\Delta} = -\frac{\nabla f(\boldsymbol{\theta})}{\|\nabla f(\boldsymbol{\theta})\|}.$$

That's the **steepest descent** direction.

But, given  $\Delta$ , how to choose  $\alpha$ ?

This is a **line search problem**,  $\alpha$  moves us along  $\theta + \alpha\Delta$ .

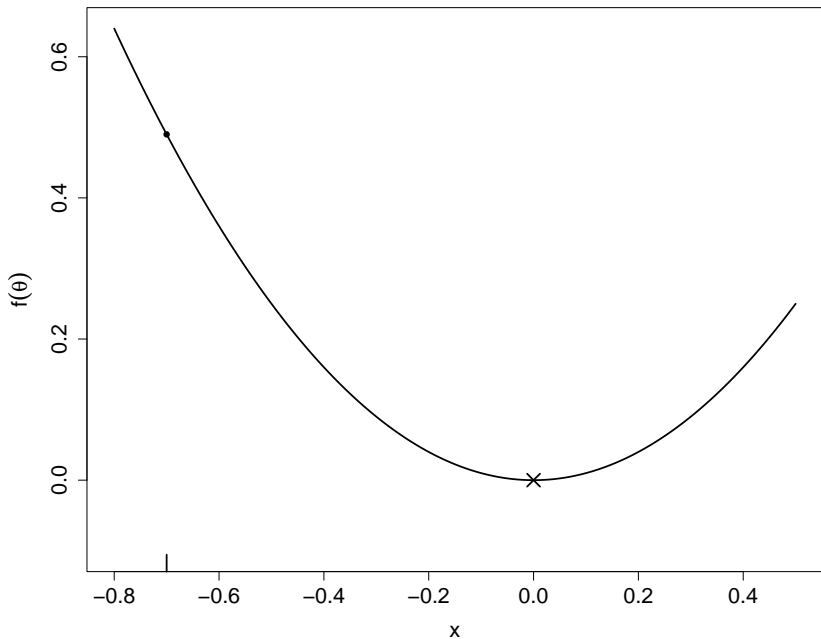
Local linear model does not give any information to choose  $\alpha$ .

Possible approaches (omitting iteration index  $k$ ):

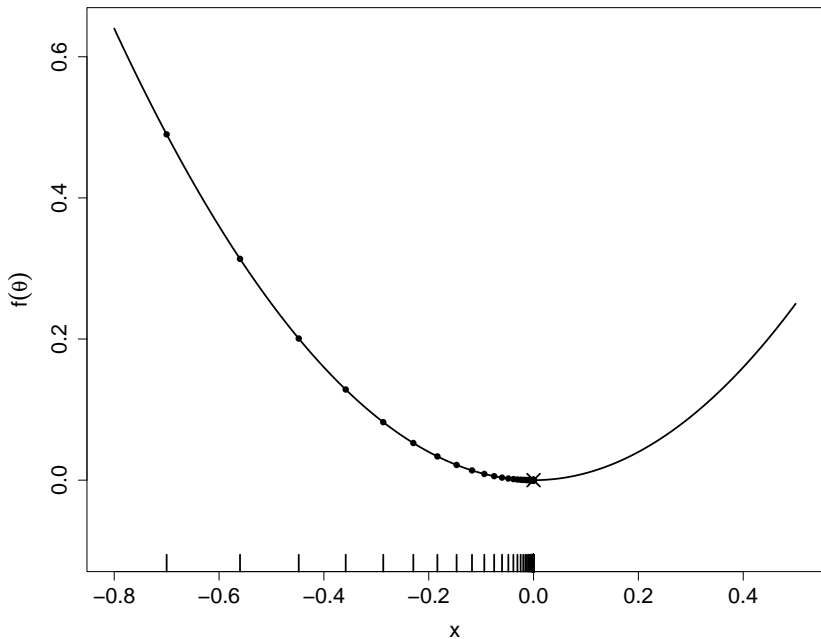
1. Find  $\alpha^* = \underset{\alpha}{\operatorname{argmin}} f(\theta + \alpha\Delta)$  (optimal but expensive);
2. Set  $\alpha = \alpha_0$  and then **backtrack**:
  - ▶ if  $f(\theta + \alpha\Delta) > f(\theta)$  set  $\alpha \leftarrow \alpha/2$  and retry;



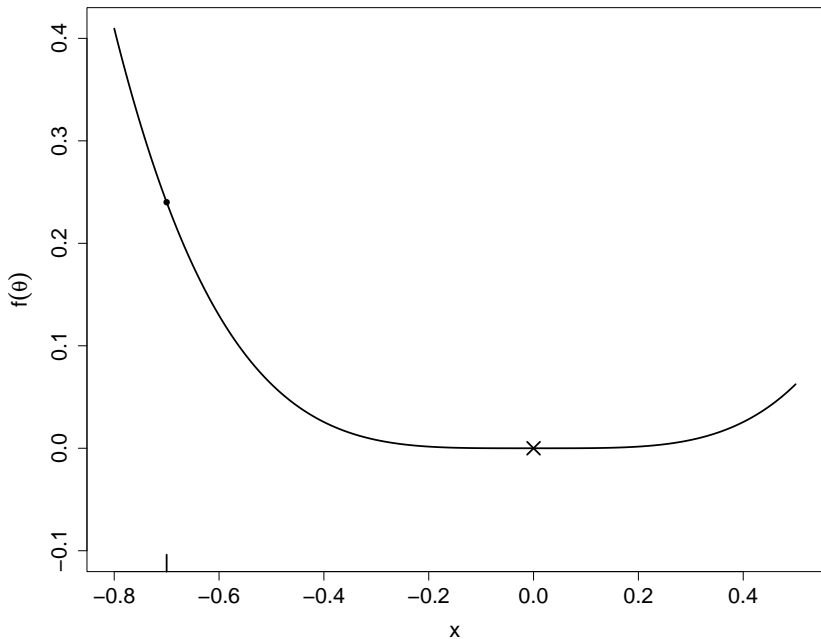
SD on  $f(\theta) = \theta^2$ ,  $f'(\theta) = 2\theta$ , with 100 steps and  $\alpha_0 = 0.1$ .



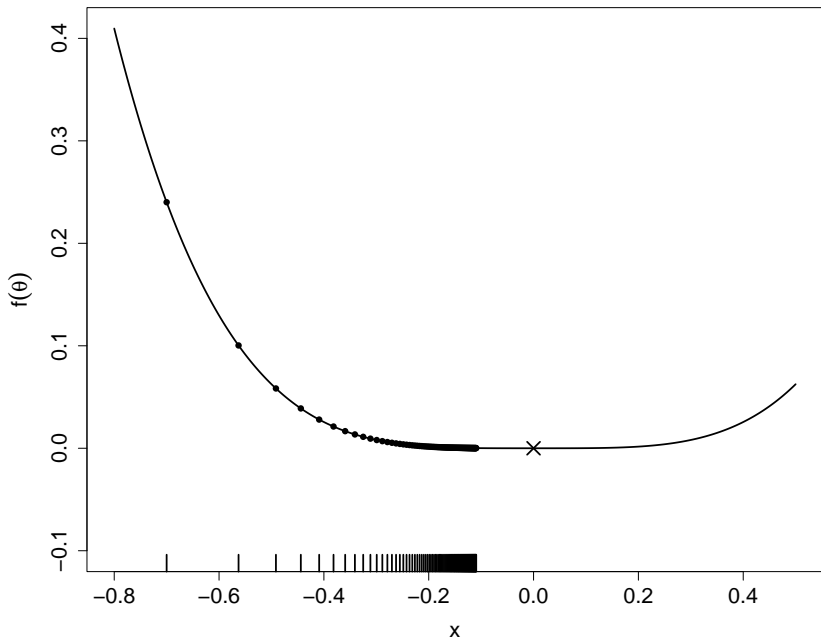
SD on  $f(\theta) = \theta^2$ ,  $f'(\theta) = 2\theta$ , with 100 steps and  $\alpha_0 = 0.1$ .



SD on  $f(\theta) = \theta^4$ ,  $f'(\theta) = 4\theta^3$ , with 100 steps and  $\alpha_0 = 0.1$ .



SD on  $f(\theta) = \theta^4$ ,  $f'(\theta) = 4\theta^3$ , with 100 steps and  $\alpha_0 = 0.1$ .



So one might look for more complicated line searches.

E.g. we might look for **sufficient** decrease of  $f(\theta + \alpha \Delta)$ .

But further problem is that  $\Delta$  from SD is itself not amazing.

It is based on a local linear approximation to the objective.

It can lead to zig-zagging behaviour.

See animated examples from Finn Lindgren:

```
library(devtools)
devtools::install_bitbucket("finnlindgren/FLtools")

library(FLtools)
FLtools::optimisation()
```

## Conclusions on SD

Proposed step  $k$ -th is:

$$\boldsymbol{\theta}^{[k+1]} \leftarrow \boldsymbol{\theta}^{[k]} - \alpha \nabla f(\boldsymbol{\theta}^{[k]}).$$

It's prudent to couple it with a line search for  $\alpha$  to avoid:

- ▶ very slow convergence (small steps)
- ▶ divergence (overly large steps)

See Nocedan and Wright (2006), Chapter 3.

It's simple, but not a great method when we actually want to reach the minimum.

In ML estimation, we want to reach the MLE exactly.

Then we use asymptotic theory to quantify uncertainty via Hessian of log-likelihood.

So SD not widely used in maximum likelihood estimation.

Exception is Big Data problems where **stochastic gradient descent** (SGD) is used.

At each iteration we compute gradient on a random subsample.

SGD is method of choice in Machine Learning, e.g. to fit NNs.

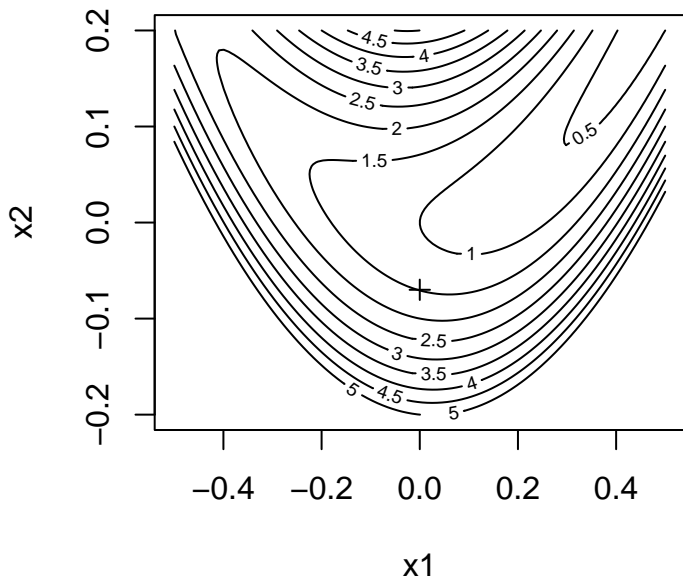
But aim is exactly minimising the objective on train data.

Noise of SGD and early stopping are used to **regularise** the fit.

So SGD is used not purely as an optimiser.

## How to improve SD

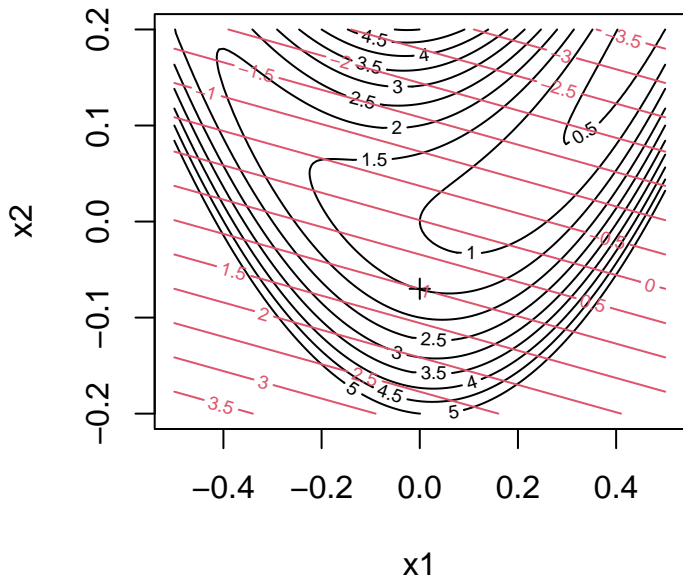
Recall the local linear model is  $f(\boldsymbol{\theta} + \boldsymbol{\Delta}) \simeq f(\boldsymbol{\theta}) + \nabla f(\boldsymbol{\theta})^T \boldsymbol{\Delta}$ .





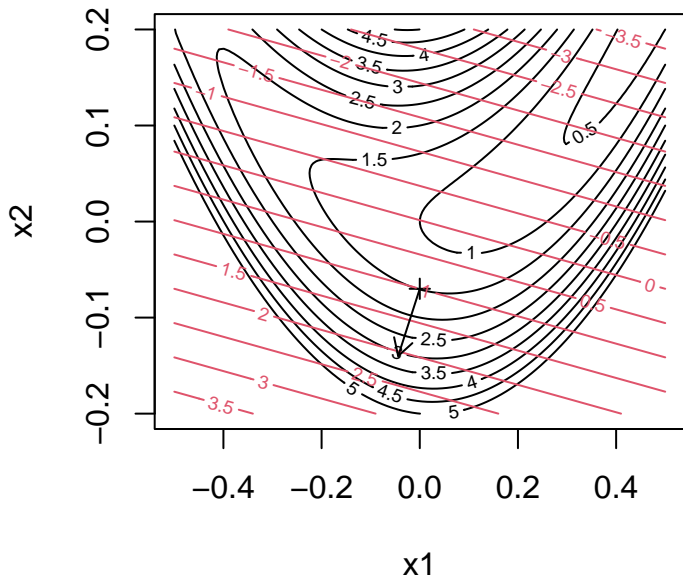
## How to improve SD

Recall the local linear model is  $f(\theta + \Delta) \simeq f(\theta) + \nabla f(\theta)^T \Delta$ .



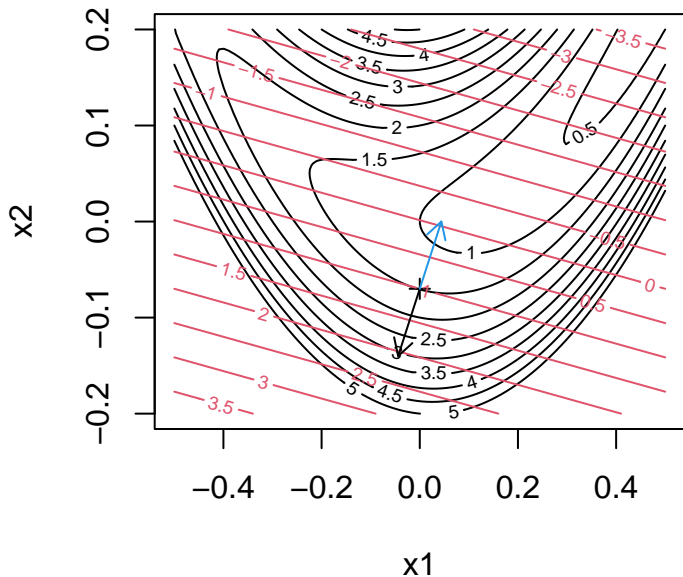
## How to improve SD

Recall the local linear model is  $f(\theta + \Delta) \simeq f(\theta) + \nabla f(\theta)^T \Delta$ .



## How to improve SD

Recall the local linear model is  $f(\theta + \Delta) \simeq f(\theta) + \nabla f(\theta)^T \Delta$ .



Note:

- ▶ Linear model is tangent to  $f(\boldsymbol{\theta})$  at  $\boldsymbol{\theta}$
- ▶  $\nabla f(\boldsymbol{\theta})$  is  $\perp$  to contour of  $f$  passing by  $\boldsymbol{\theta}$ .
- ▶ Any  $\boldsymbol{\Delta}$  s.t.  $\boldsymbol{\Delta}^T \nabla f(\boldsymbol{\theta}) < 0$  goes downhill on local model

Hence we can modify SD step via **preconditioning** matrix  $\mathbf{A}^{[k]}$ :

$$\boldsymbol{\theta}^{[k+1]} \leftarrow \boldsymbol{\theta}^{[k]} - \alpha \mathbf{A}^{[k]} \nabla f(\boldsymbol{\theta}^{[k]}).$$

Effect of  $\mathbf{A}^{[k]}$  is to rotate and rescale SD step.

What are the conditions on  $\mathbf{A}^{[k]}$  to ensure that we go downhill on local linear model?

$\mathbf{A}^{[k]}$  must be **positive definite**:  $\mathbf{x}^\top \mathbf{A}^{[k]} \mathbf{x} \succ \forall \mathbf{x}$ .

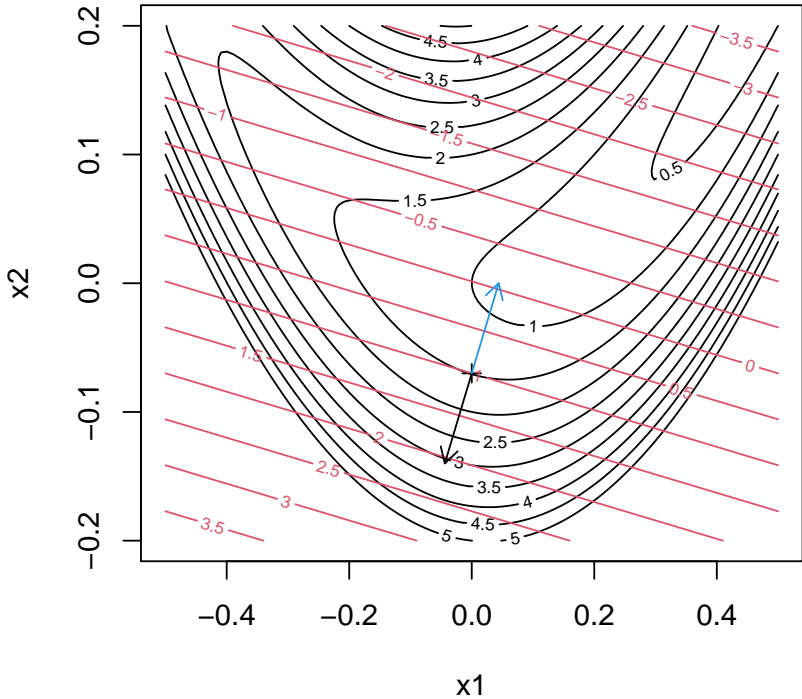
Why does it work? Drop index  $k$  and look at local linear model

$$f(\boldsymbol{\theta} + \boldsymbol{\Delta}) \simeq f(\boldsymbol{\theta}) + \nabla f(\boldsymbol{\theta})^\top \boldsymbol{\Delta}.$$

plug in  $\boldsymbol{\Delta} = -\mathbf{A} \nabla f(\boldsymbol{\theta})$

$$f(\boldsymbol{\theta} + \boldsymbol{\Delta}) \simeq f(\boldsymbol{\theta}) - \nabla f(\boldsymbol{\theta})^\top \mathbf{A} \nabla f(\boldsymbol{\theta}) < f(\boldsymbol{\theta}).$$

P.d.  $\mathbf{A}$  ensures that it does not rotate  $-\nabla f(\boldsymbol{\theta})$  more than  $\pi/2$ .



## How to find a good preconditioner $\mathbf{A}$ ?

**Newton's method** chooses  $\mathbf{A}$  via a better local model of  $f(\boldsymbol{\theta})$ .

Consider 2nd order Taylor approximation

$$f(\boldsymbol{\theta} + \boldsymbol{\Delta}) \simeq f(\boldsymbol{\theta}) + \nabla f(\boldsymbol{\theta})^\top \boldsymbol{\Delta} + \frac{1}{2} \boldsymbol{\Delta}^\top \nabla^2 f(\boldsymbol{\theta}) \boldsymbol{\Delta}.$$

For  $\|\boldsymbol{\Delta}\| < 1$  this model is better, with smaller remainder.

How do we use it to choose step?

We minimise the local quadratic model for  $f(\boldsymbol{\theta})$  w.r.t.  $\boldsymbol{\Delta}$ .

Differentiating the local model

$$f(\boldsymbol{\theta} + \boldsymbol{\Delta}) \simeq \tilde{f}(\boldsymbol{\theta} + \boldsymbol{\Delta}) = f(\boldsymbol{\theta}) + \nabla f(\boldsymbol{\theta})^\top \boldsymbol{\Delta} + \frac{1}{2} \boldsymbol{\Delta}^\top \nabla^2 f(\boldsymbol{\theta}) \boldsymbol{\Delta}.$$

w.r.t.  $\boldsymbol{\Delta}$  leads to

$$\nabla_{\boldsymbol{\Delta}} \tilde{f}(\boldsymbol{\theta} + \boldsymbol{\Delta}) = \nabla f(\boldsymbol{\theta}) + \nabla^2 f(\boldsymbol{\theta}) \boldsymbol{\Delta}$$

and setting this to zero leads to

$$\boldsymbol{\Delta} = - \left( \nabla^2 f(\boldsymbol{\theta}) \right)^{-1} \nabla f(\boldsymbol{\theta}).$$



Hence  $k$ -th iteration is

$$\boldsymbol{\theta}^{[k+1]} \leftarrow \boldsymbol{\theta}^{[k]} - \left( \nabla^2 f(\boldsymbol{\theta}^{[k]}) \right)^{-1} \nabla f(\boldsymbol{\theta}^{[k]}).$$

or

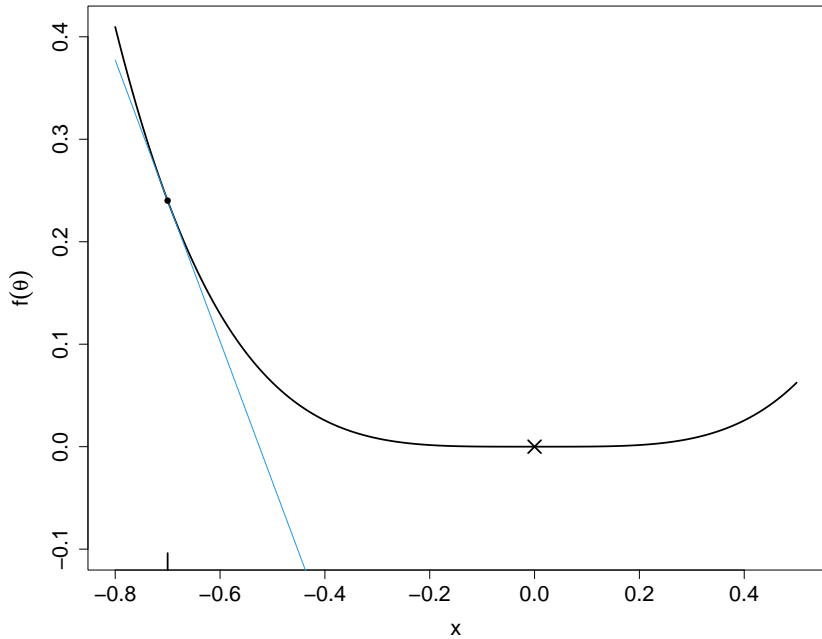
$$\boldsymbol{\theta}^{[k+1]} \leftarrow \boldsymbol{\theta}^{[k]} - \mathbf{A}^{[k]} \nabla f(\boldsymbol{\theta}^{[k]}).$$

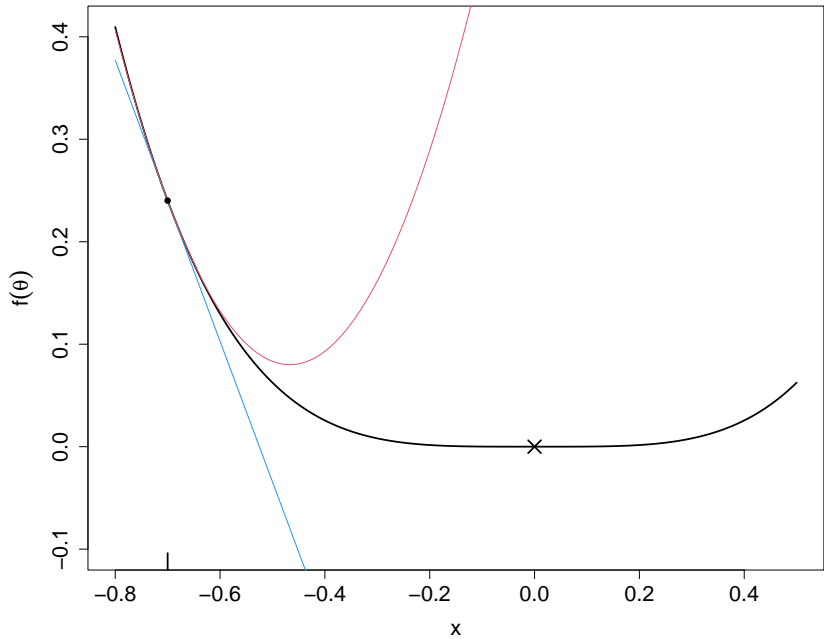
where  $\left( \nabla^2 f(\boldsymbol{\theta}^{[k]}) \right)^{-1}$  is the preconditioner  $\mathbf{A}^{[k]}$ .

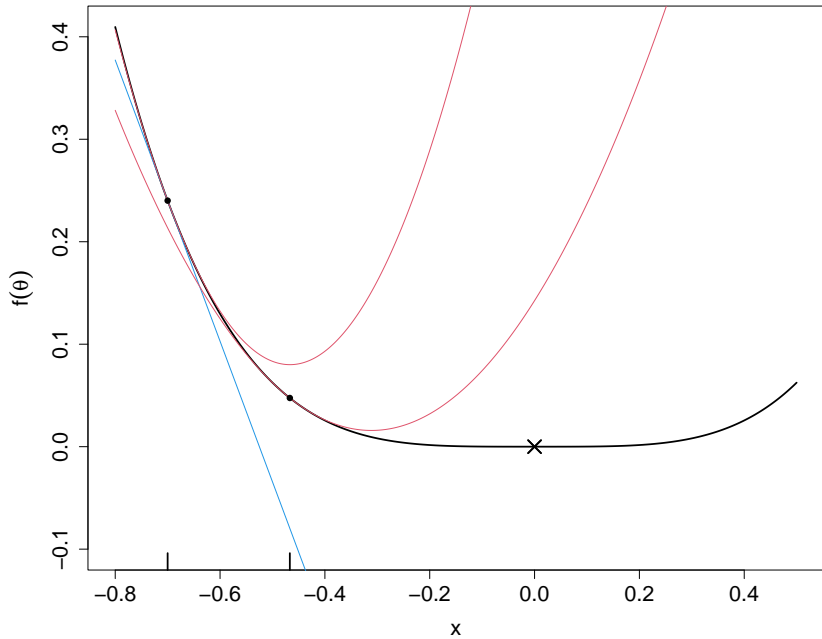
Questions:

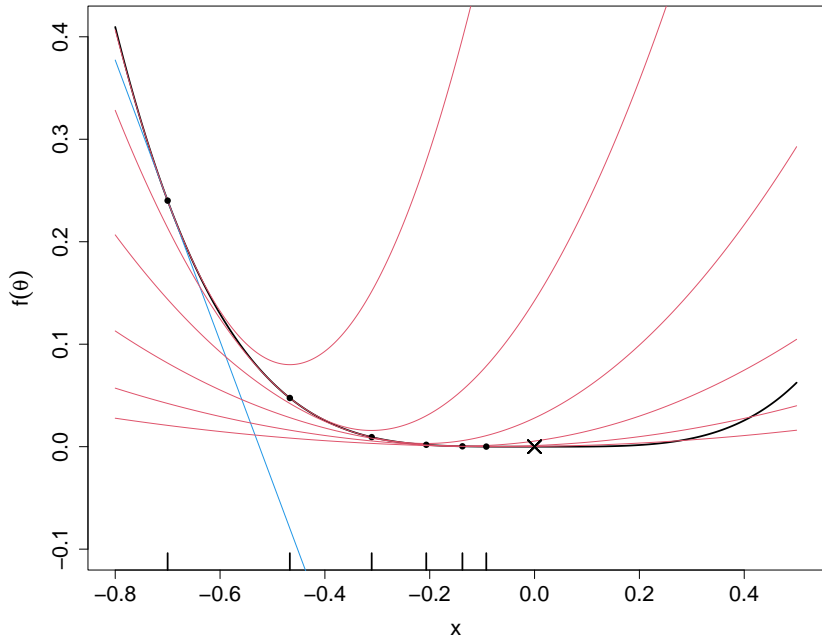
1. Where this the step size  $\alpha$  go?
2. What happens if  $\mathbf{H} = \nabla^2 f(\boldsymbol{\theta})$  is not positive definite?

Assume 2 is fine, consider 1.









So quadratic model suggests iterating using

$$\boldsymbol{\theta}^{[k+1]} \leftarrow \boldsymbol{\theta}^{[k]} - \alpha \boldsymbol{\Delta}^{[k]},$$

with  $\alpha = 1$ , that is taking full steps.

In practice, Newton should be coupled with a **backtracking** search.

Set  $\alpha = \alpha_0 = \mathbf{1}$  (!! ) and then **backtrack**:

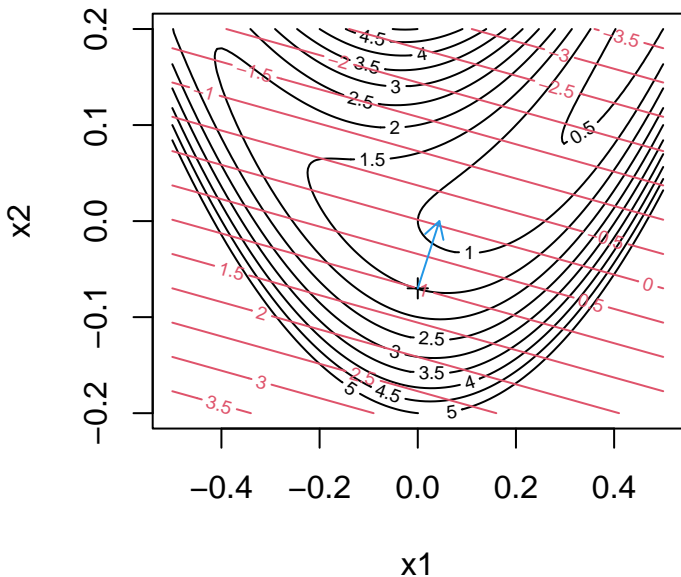
- ▶ if  $f(\boldsymbol{\theta} + \alpha \boldsymbol{\Delta})$  is not **sufficiently lower** than  $f(\boldsymbol{\theta})$  set  $\alpha \leftarrow \alpha/2$  and retry.

Sufficiently lower because step might be too long, see Nocedal and Wright 2006, Fig 3.2.

Note that for SD **we did have not a clue** about  $\alpha_0$ .

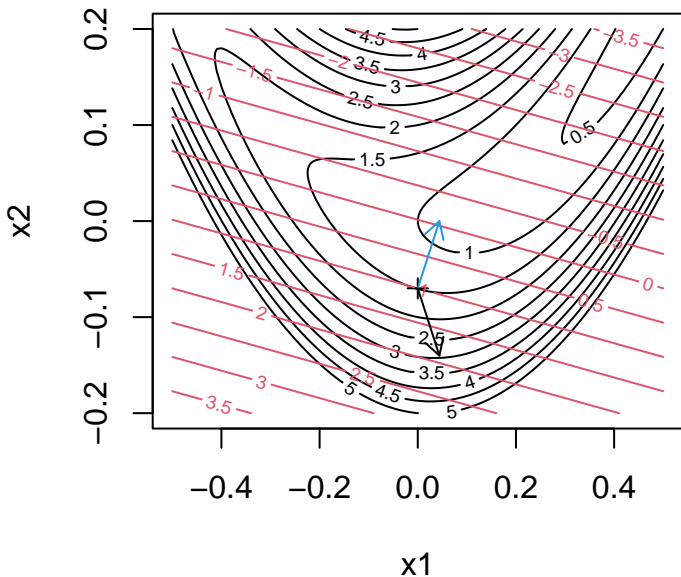
Back to what happens if  $\mathbf{H} = \nabla^2 f(\theta)$  is not positive definite?

Using  $\mathbf{H}$  directly might sent you uphill!



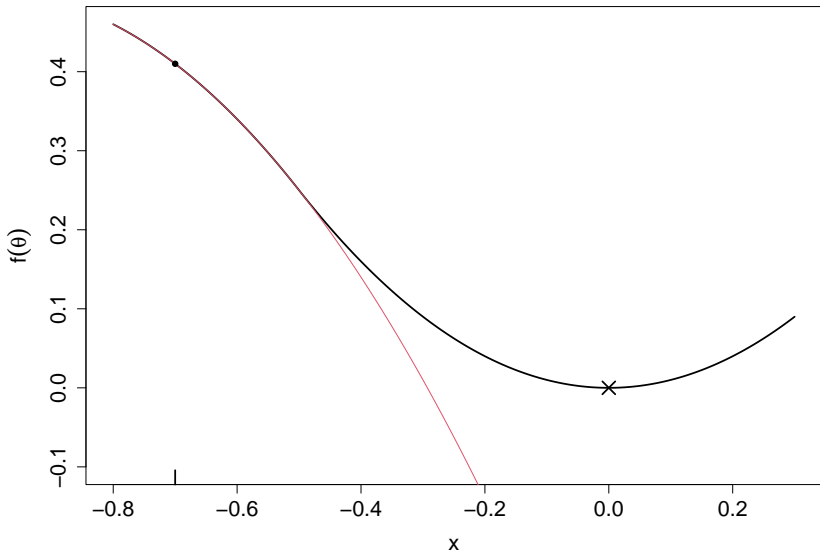
Back to what happens if  $\mathbf{H} = \nabla^2 f(\theta)$  is not positive definite?

Using  $\mathbf{H}$  directly might sent you uphill!





In 1D dimension, consider:



Here setting  $\nabla_{\Delta} \tilde{f}(\theta + \Delta) = \mathbf{0}$  sends us far to the left.

If this happens, don't panic.

Hessian might still be useful, but needs some fixing.

$\mathbf{H}$  might be problematic in certain directions, e.g.

$$\mathbf{H} = \begin{bmatrix} H_{11} & 0 & \dots & 0 \\ 0 & H_{22} & & \\ \vdots & & \ddots & \\ 0 & & & H_{pp} \end{bmatrix}.$$

Here Newton is a scaled SD, where  $\mathbf{H}^{-1}$  scales each dimension.

But if  $H_{jj} < \epsilon$  (where  $0 < \epsilon \ll 1$  is a small tolerance) then either

- ▶  $|H_{jj}| < \epsilon$  (massive step)  $\rightarrow$  set  $H_{jj} \leftarrow \epsilon$
- ▶ or  $H_{jj} < -\epsilon$  (wrong direction)  $\rightarrow$  set  $H_{jj} \leftarrow |H_{jj}|$

In second case we assume that  $|H_{jj}|$  is informative.

“Targeted intervention” idea extends to full  $\mathbf{H}$  (correlated  $\theta$ ):

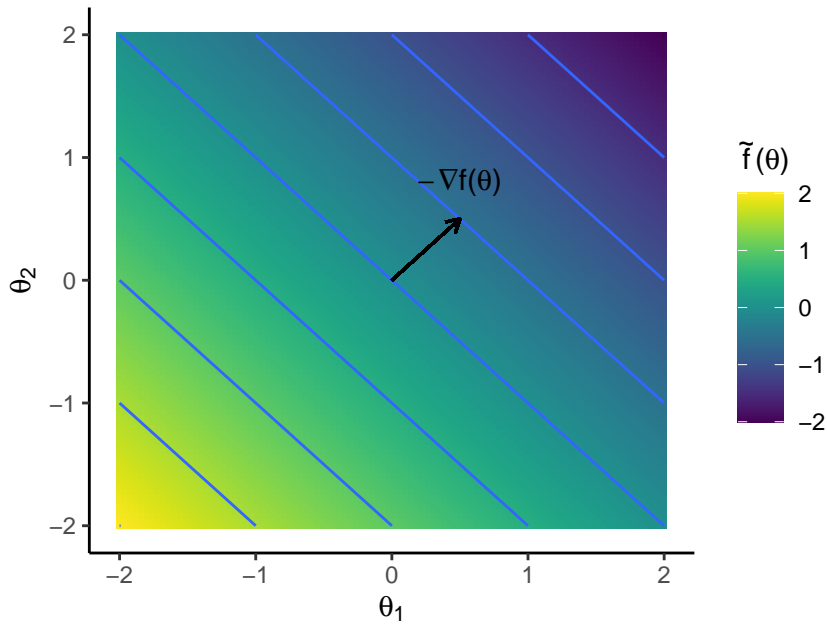
1. Eigen-decompose  $\mathbf{H}$ .
2. For any eigenvalue  $\lambda_j < \epsilon$  then either:
  - ▶  $|\lambda_j| < \epsilon$  (indefinite)  $\rightarrow$  set  $\lambda_j \leftarrow \epsilon$
  - ▶ or  $\lambda_j < -\epsilon$  (neg. definite)  $\rightarrow$  set  $\lambda_j \leftarrow |\lambda_j|$

Other solutions:

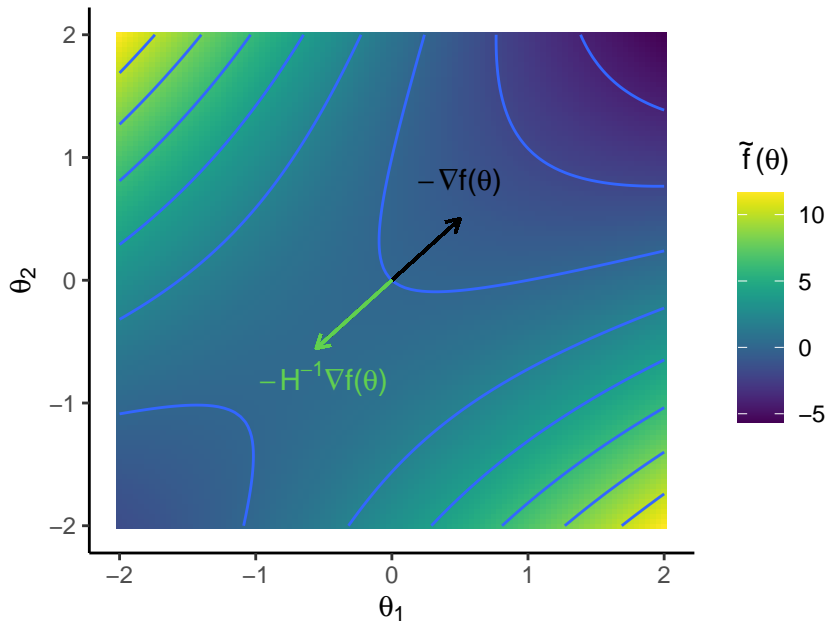
- ▶ Set  $\mathbf{H} \leftarrow \mathbf{H} + c\mathbf{I}$  for large enough  $c > 0$
- ▶ When  $f(\theta)$  is neg. log-likelihood use  $\mathbb{E}(\mathbf{H})$  (**Fisher scoring**)

How does the “eigen-fix” modify the local model, visually?

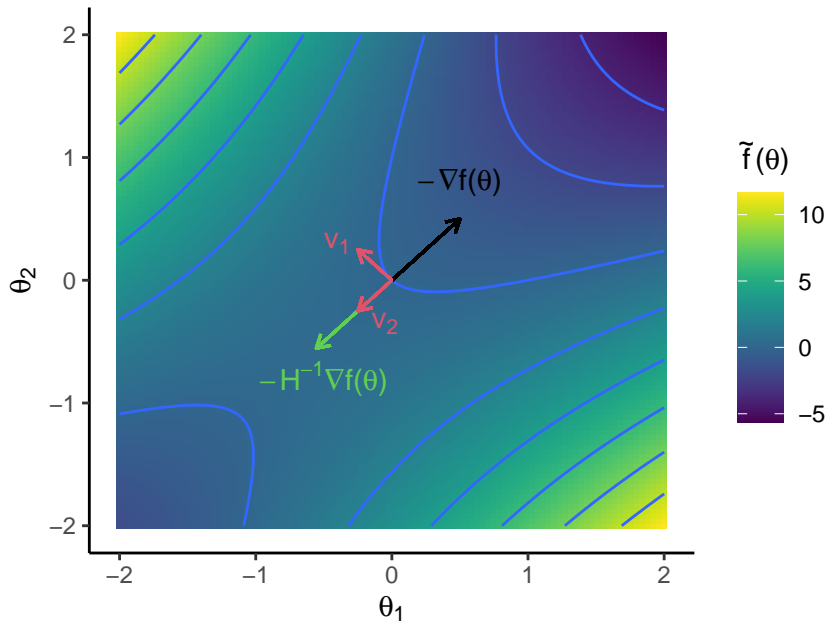
$$\nabla f(\theta) = \begin{bmatrix} -0.5 \\ -0.5 \end{bmatrix}, \quad \nabla^2 f(\theta) = \begin{bmatrix} 1 & -1.9 \\ -1.9 & 1 \end{bmatrix}, \quad \lambda_1 = 2.9, \quad \lambda_2 = -0.9.$$



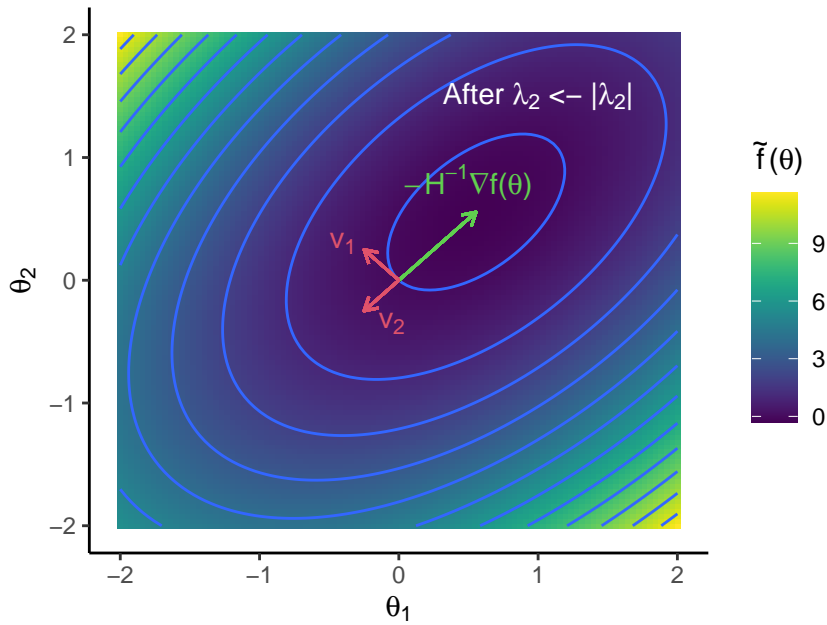
$$\nabla f(\theta) = \begin{bmatrix} -0.5 \\ -0.5 \end{bmatrix}, \quad \nabla^2 f(\theta) = \begin{bmatrix} 1 & -1.9 \\ -1.9 & 1 \end{bmatrix}, \quad \lambda_1 = 2.9, \quad \lambda_2 = -0.9.$$



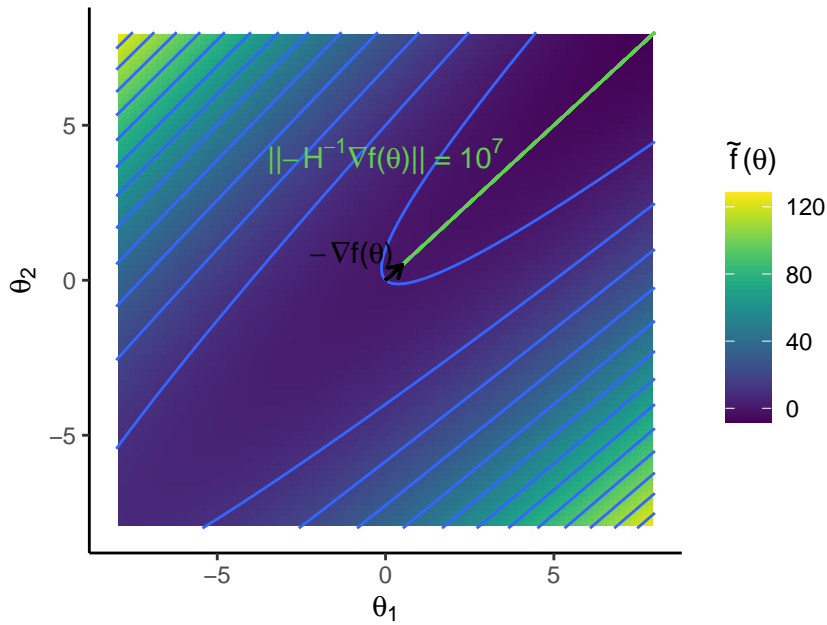
$$\nabla f(\theta) = \begin{bmatrix} -0.5 \\ -0.5 \end{bmatrix}, \quad \nabla^2 f(\theta) = \begin{bmatrix} 1 & -1.9 \\ -1.9 & 1 \end{bmatrix}, \quad \lambda_1 = 2.9, \quad \lambda_2 = -0.9.$$



$$\nabla f(\theta) = \begin{bmatrix} -0.5 \\ -0.5 \end{bmatrix}, \quad \nabla^2 f(\theta) = \begin{bmatrix} 1 & -1.9 \\ -1.9 & 1 \end{bmatrix}, \quad \lambda_1 = 2.9, \quad \lambda_2 = -0.9.$$

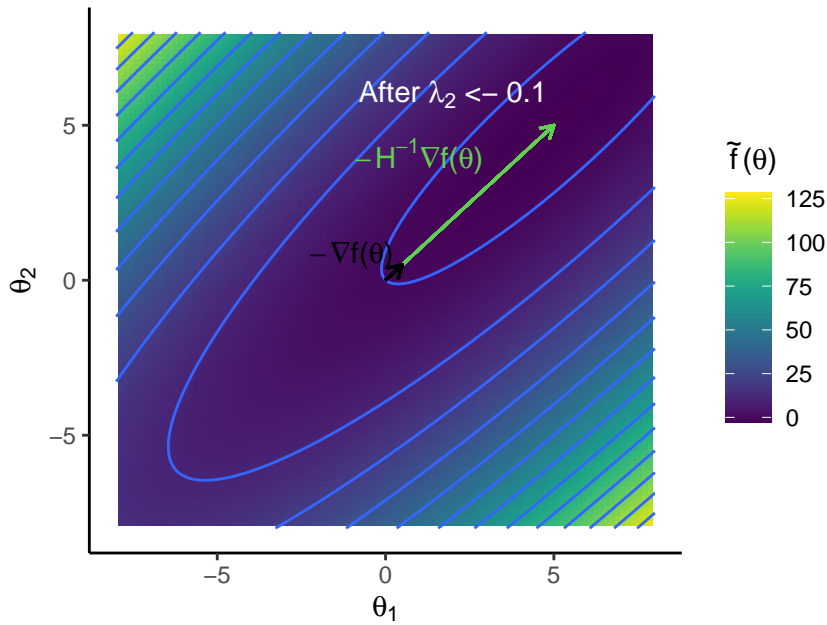


If  $\nabla^2 f(\boldsymbol{\theta}) = \begin{bmatrix} 1 & -(1 - 10^{-7}) \\ -(1 - 10^{-7}) & 1 \end{bmatrix}$  then  $\lambda_2 = 10^{-7}$ .





If  $\nabla^2 f(\theta) = \begin{bmatrix} 1 & -(1 - 10^{-7}) \\ -(1 - 10^{-7}) & 1 \end{bmatrix}$  then  $\lambda_2 = 10^{-7}$ .



We have a “bullet proof” implementation of Newton’s method.

At  $k = 0$  we start at  $\theta^{[0]}$  and iteratively:

1. evaluate  $f(\theta^{[k]})$ ,  $\nabla f(\theta^{[k]})$  and  $\nabla^2 f(\theta^{[k]})$ .
2. Test whether  $\theta^{[k]}$  is a minimum and, if it is, stop here.
3. If  $\nabla^2 f(\theta^{[k]})$  is not p.d., perturb it to make it p.d..
4. Solve

$$\nabla^2 f(\theta^{[k]}) \mathbf{\Delta} = -\nabla f(\theta^{[k]}).$$

to find the search direction  $\mathbf{\Delta}$ .

5. If  $f(\theta^{[k]} + \mathbf{\Delta})$  is not sufficiently lower than  $f(\theta^{[k]})$ , repeatedly halve  $\mathbf{\Delta}$  until it is.
6. Set  $\theta^{[k+1]} \leftarrow \theta^{[k]} + \mathbf{\Delta}$  and back to step 1.

Let’s look at Newton’s performance using

`FLtools::optimisation()`

## Conclusion on SD and Newton's method

Weaknesses of steepest descent:

- ▶ Provides direction but no information on step-length;
- ▶ Tends to zig-zag.

Improved by Newton's method via a better local model, but:

1. Needs to compute the Hessian matrix  $\mathbf{H}$ ;
2. Need to ensure its positive definiteness.

Both need line search to guarantee convergence.

Stochastic SD is widely used, stochastic Newton less so.

# References

On optimisation methods:

- ▶ Nocedal and Wright (2006) Numerical Optimization, 2nd ed.

On problems when fitting chaotic models:

- ▶ Fasiolo, M., Pya, N. and Wood, S.N., 2016. A comparison of inferential methods for highly nonlinear state space models in ecology and epidemiology. *Statistical Science*, pp.96-118.