

Significant Advance in Calculating EOF From a **Very Large** Data set.

Huug van den Dool

Some 1% / 99% comments

- 99% of the Reanalysis users are happy with 1% of the data
- The other 99% of the data is required (if not demanded) by the remaining 1% of the users.
- That 1% of the users thus has a lot of impact on data storage and data transmission issues
- How do you split data wisely into 1%/99%?
One common way: High-Res and Low-Res versions (both time and space).

W.J.A. Kuipers 1970 seminar at Dutch Weather Service (KNMI) about EOF

- This was not about teleconnections, NAO, PNA
- Not about the workings of nature or to 'let the data speak' diagnostics
- It was about reducing a data set to $\sim 0.1-1\%$ of its original size to fit a ~ 1970 computer
- Curiously, while EOFs are maximally efficient for data compression, it is very costly (CPU, Memory) to calculate EOFs, but, keep in mind: you have to do this only once.
- Curiously, data compression is as desirable now as it was in 1970.

- Almost everybody has calculated EOF
- My guess: They **all** did it by first calculating the covariance matrix Q or Q^a .
- Evaluating the elements of the covariance matrix requires $n_s * n_s * n_t$ (Q) or $n_t * n_t * n_s$ (Q^a) multiplications where n_s and n_t are the number of points in space and time.
- Example (CFSR) Reanalysis full resolution data:
 space = $1152 * 576 = 663552$ points, and
 time = $32 * 365 * 24 = 280320$
 ($n_s * n_t = 1.86 * 10^{11}$ is one unit)
- ➔ It would take $5.22 * 10^{16}$ ($1.23 * 10^{17}$) multiplications to fill Q^a (Q) (for one variable). That is before any EOFs are calculated.!! {CPU cost is n_t units}

Help from Bob Kistler:

- Rearranged the do loops
- Applied both methods of multi-tasking: MPI (message passing) and OPENMP (threading)
- George vandenBerghe helped to configure batch jobs for power 6 machine so as to use all 64 logical cpu's on each physical node
- Even so, it is hard in terms of CPU to fill Q (mind you: this is only for a single variable at one level).

The good news:

- We can find EOFs without first filling the covariance matrix.
- This advance addresses CPU problems, not so much memory problems (which are also an obstacle, but 2nd to CPU).
- So we can calculate EOFs from CFSR ?

Lay-out talk

- The new(?) method
- Examples to show that it actually works and yields the expected results
- Why does it work?
- Origin of the method

$$\text{Basics: } f(s, t) = \sum_m \alpha_m(t) e_m(s) \quad (0)$$

- Multiply lhs and rhs of (0) by $\alpha_n(t)$ and sum over all times t (n is a specific mode number). Result:

$$e_m(s) = \sum_t \alpha_m(t) f(s, t) / \sum_t \alpha_m^2(t) \quad (1)$$

- Multiply lhs and rhs side of (0) by $e_n(s)$ and sum over all space s . Result:

$$\alpha_m(t) = \sum_s w(s) e_m(s) f(s, t) / \sum_s w(s) e_m^2(s) \quad (2)$$

where $w(s)$ are spatial weights, not shown below.

$$e_m(s) = \sum_t \alpha_m(t) f(s, t) / \sum_t \alpha_m^2(t) \quad (1)$$

$$\alpha_m(t) = \sum_s e_m(s) f(s, t) / \sum_s e_m^2(s) \quad (2)$$

- The above are orthogonality relationships. If we know $\alpha_m(t)$ and $f(s,t)$, $e_m(s)$ can be calculated trivially. If we know $e_m(s)$ and $f(s,t)$, $\alpha_m(t)$ can be calculated trivially. This is the basis of the iteration.
- Randomly pick (or make) a time series $\alpha^0(t)$, and stick into (1). This yields $e^0(t)$. Stick $e^0(t)$ into (2). This yields $\alpha^1(t)$. This is one iteration. Etc. This generally converges to the first EOF $\alpha_1(t)$, $e_1(s)$. CPU cost 2 units per iteration.
- $f^{\text{reduced}}(s,t) = f(s,t) - \alpha_1(t) e_1(s)$ and repeat. One finds mode#2. Etc.

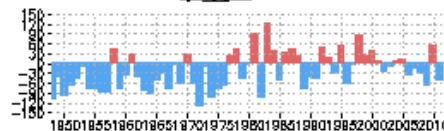
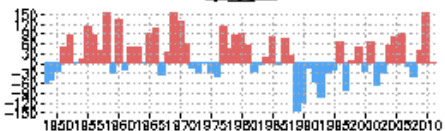
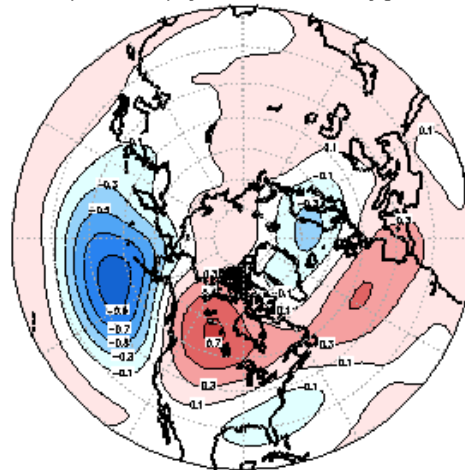
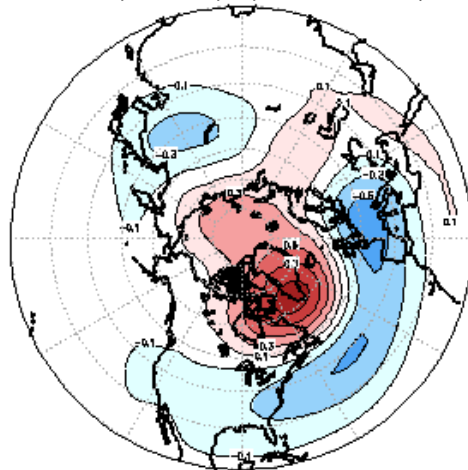
Example to show it works

- This is a lowres example when the cov matrix can be calculated.
- 1948-2011 JFM seasonal Z500 mean at 2.5 degree grid. 64 time levels
- Domain 20N-pole ($144 * 29 = 4176$ gridpoints)

EOF for JFM 1948–2011 HGT 500 mb

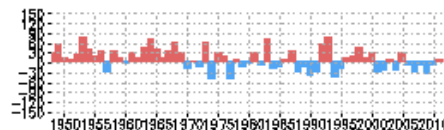
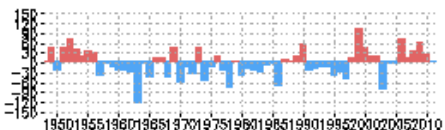
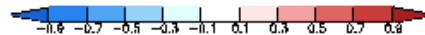
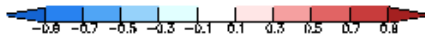
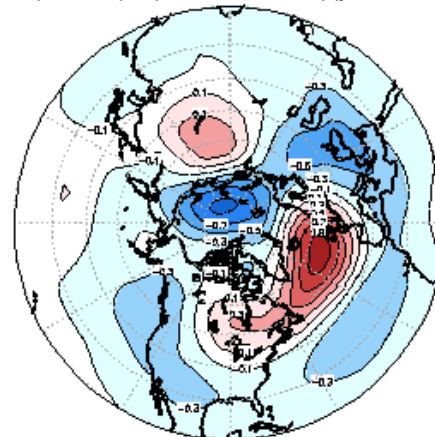
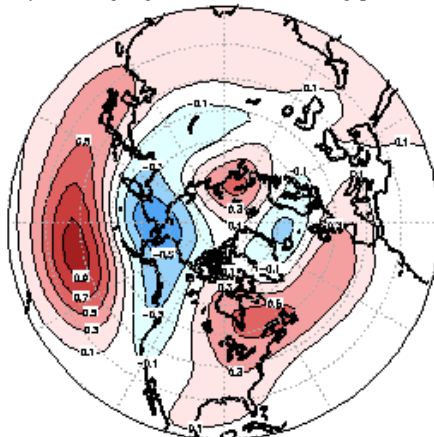
EOF1 (23.0 %EV) (seed=62.5N,55W)

EOF2 (19.4 %EV) (seed=55N,107.5W)(partial 1)



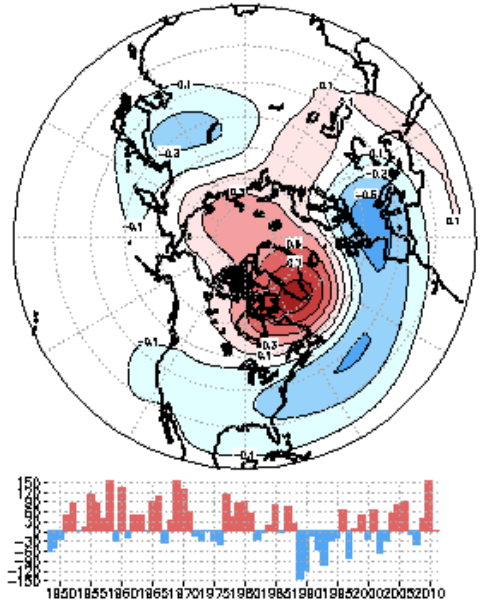
Example using a covariance matrix

EOF3 (9.1 %EV) (seed=37.5N,172.5W)(partial 1&2) EOF4 (7.1 %EV) (seed=57.5N,20W)(partial 1&2&3)

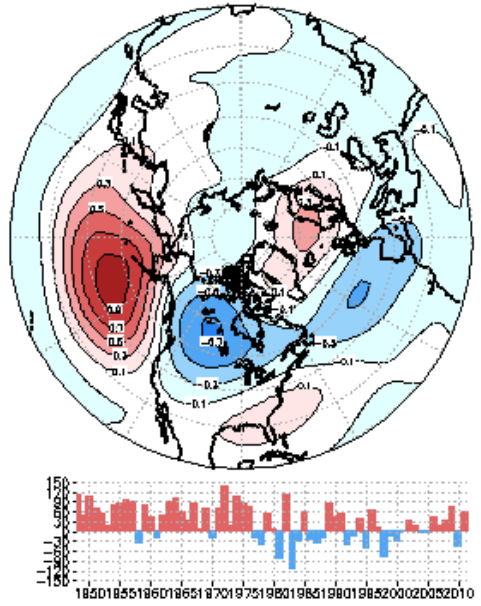


EOF for JFM 1948–2011 HGT 500 mb(iter)

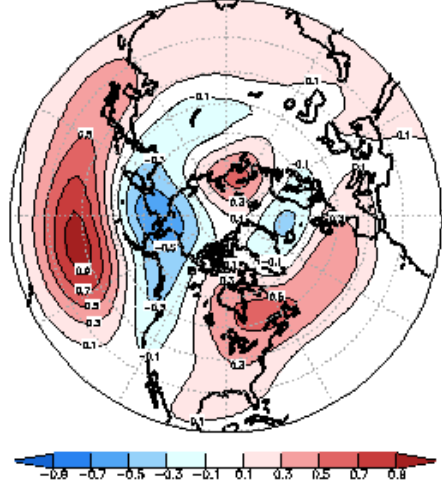
EOF1 (23.0 %EV) (seed=guess)



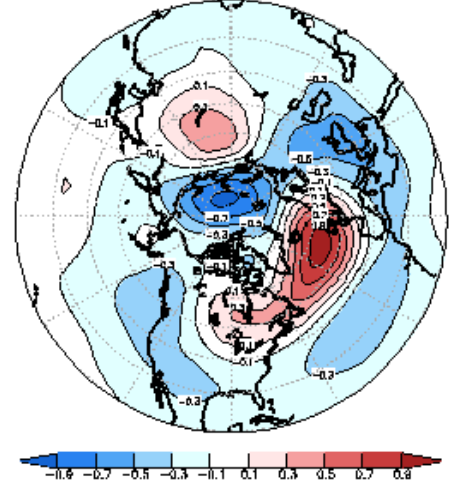
EOF2 (19.4 %EV) (seed=guess)(partial 1)



EOF3 (9.1 %EV) (seed=guess)(partial 1&2)



EOF4 (7.2 %EV) (seed=guess)(partial 1&2&3)

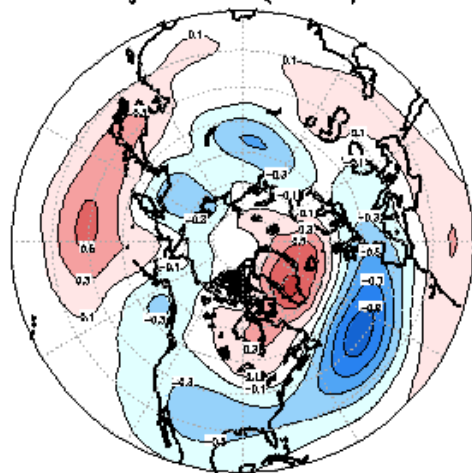


Example using an iterative method (100 iterations)

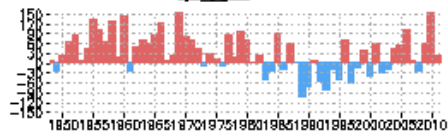
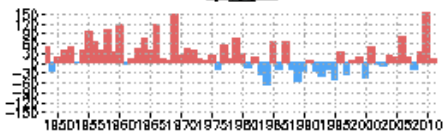
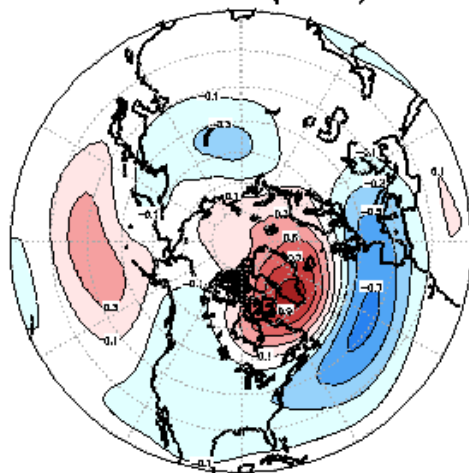
The route to EOF#1 taken by the
iteration from a guess

Path to EOF1 for JFM 1948–2011 Z500mb

guess EOF1 (iter=0.5)

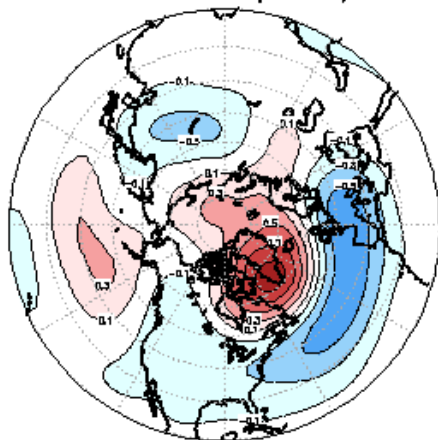


towards EOF1 (iter=1.5)

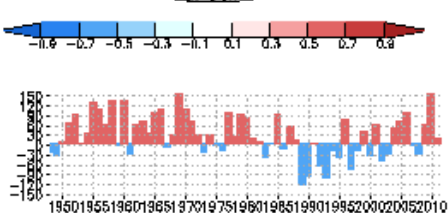
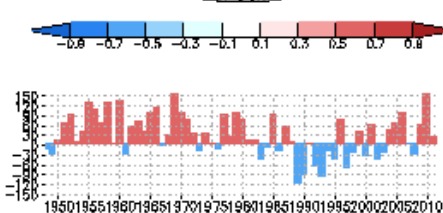
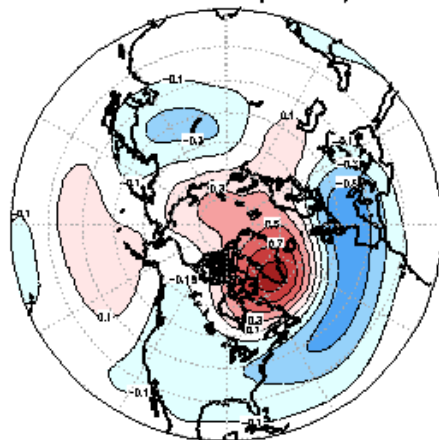


The route to EOF1 from a guess field (pick: (1948+2010)/2)
Iterations 1-4

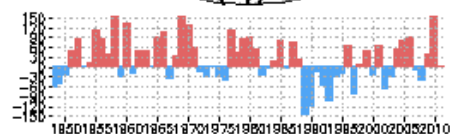
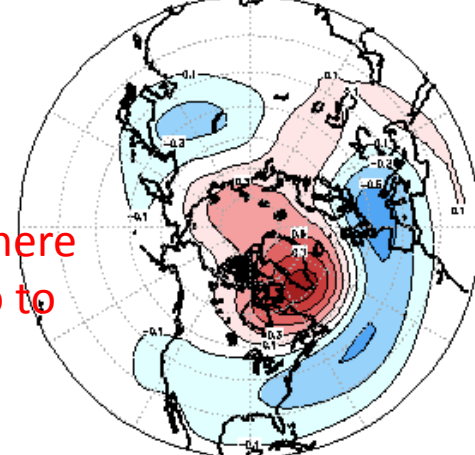
towards EOF1 (iter=2.5)



towards EOF1 (iter=3.5)



EOF1 (23.0 %EV) (seed=guess)

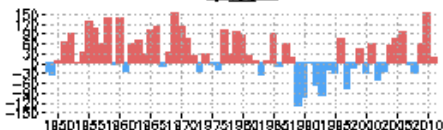
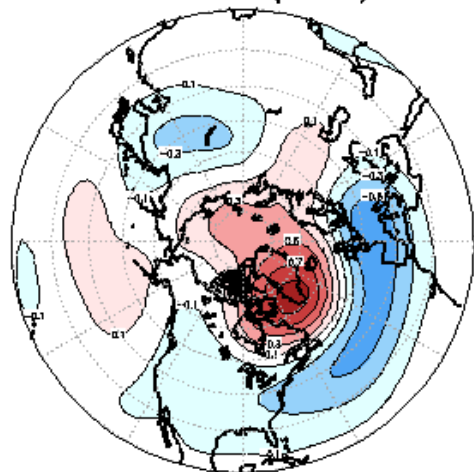


→ Where we go to

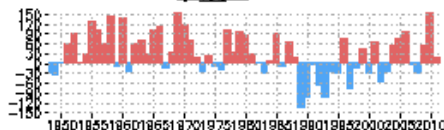
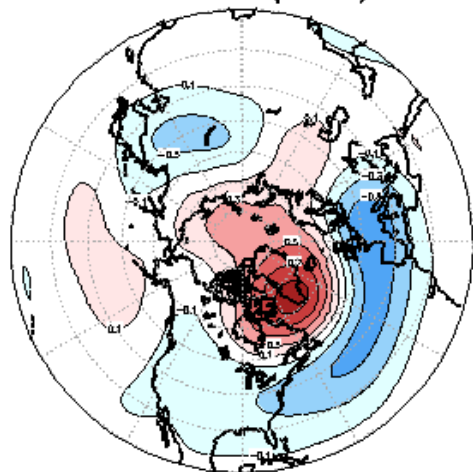
EOF3 (9.1 %EV) (seed=guess)(partial 1&2)

Path to EOF1 for JFM 1948–2011 Z500mb

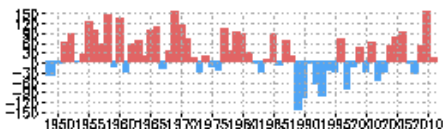
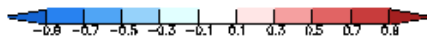
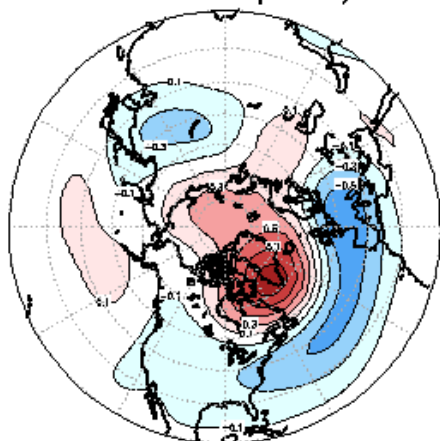
towards EOF1 (iter=4.5)



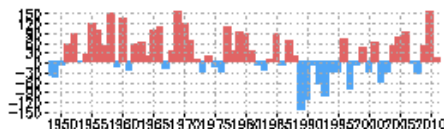
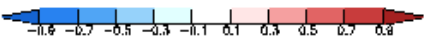
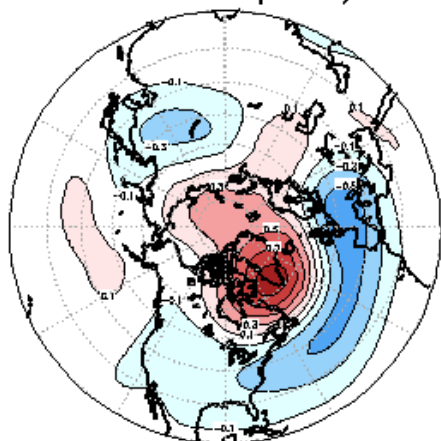
towards EOF1 (iter=5.5)



towards EOF1 (iter=6.5)

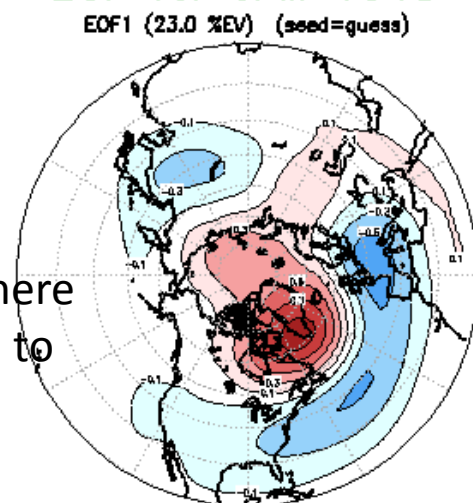


towards EOF1 (iter=7.5)

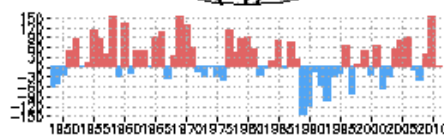


The route to EOF1 from a guess field (pick: (1948+2010)/2)
Iterations 5-8

EOF1 (23.0 %EV) (seed=guess)



→ Where we go to

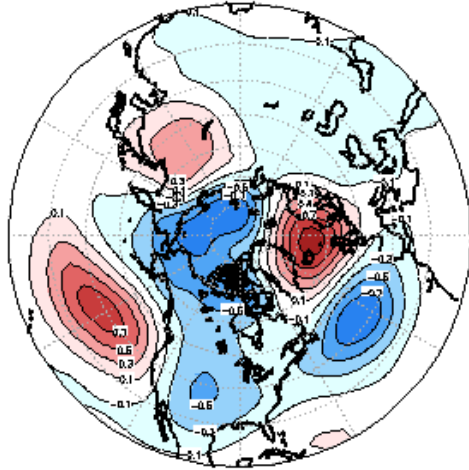


EOF3 (9.1 %EV) (seed=guess)(partial 1&2)

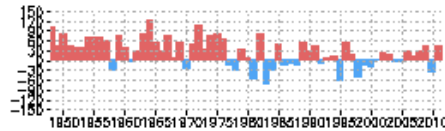
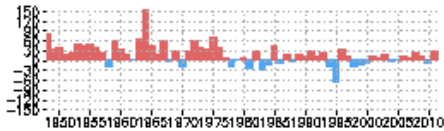
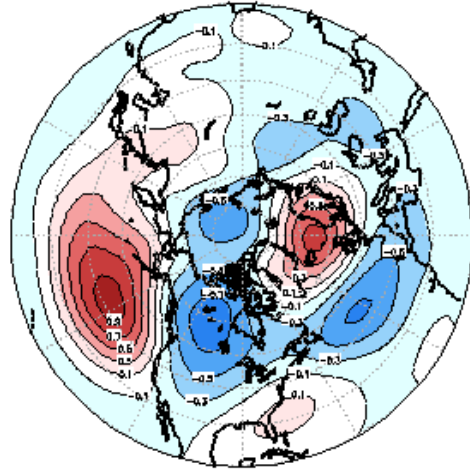
An example of the difficulties with
convergence:

Path to EOF1 for JFM 1948–2011 Z500mb

guess EOF1 (iter=0.5)

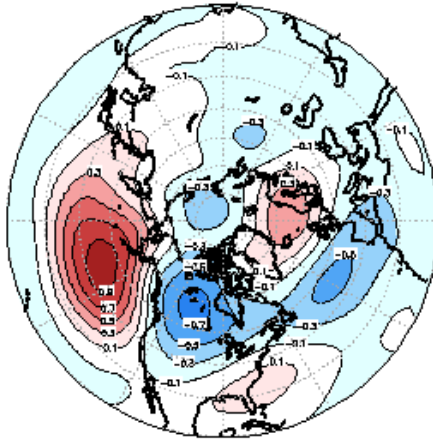


towards EOF1 (iter=1.5)

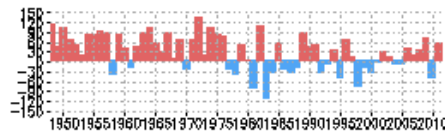
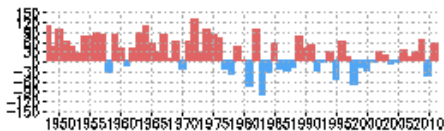
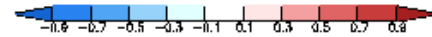
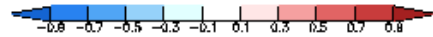
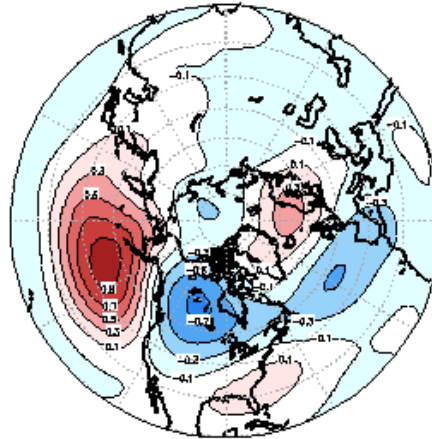


The route to EOF1 from a **very different** guess field (pick: 1964) Iterations 1-4

towards EOF1 (iter=2.5)

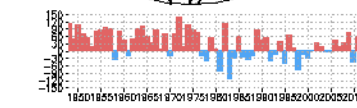
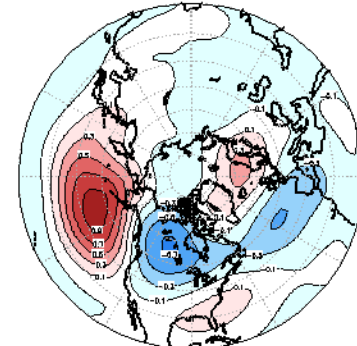
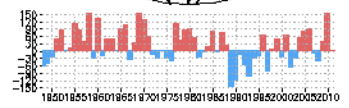
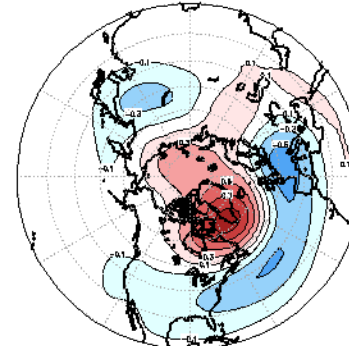


towards EOF1 (iter=3.5)



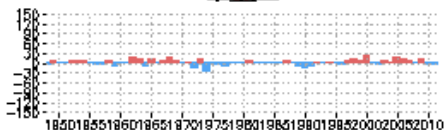
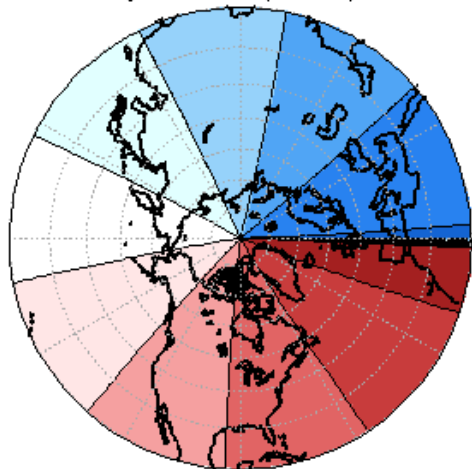
↓ Where we go to

EOF for JFM 1948–2011 HGT 500 mb(iter)
 EOF1 (23.0 %EV) (seed=guess) EOF2 (19.4 %EV) (seed=guess)(partial 1)

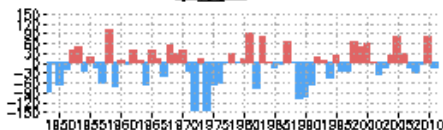
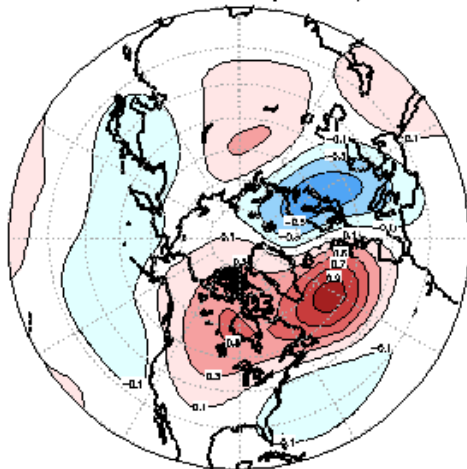


Path to EOF1 for JFM 1948–2011 Z500mb

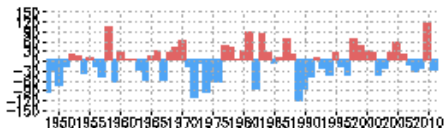
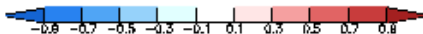
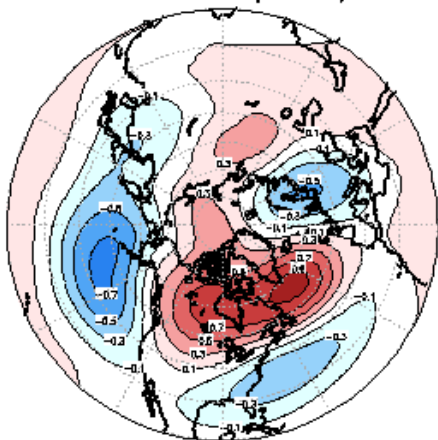
guess EOF1 (iter=0.5)



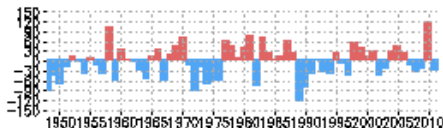
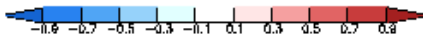
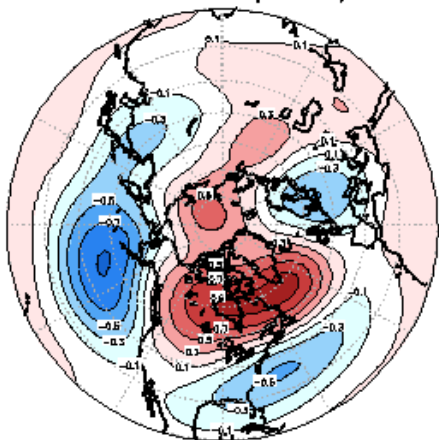
towards EOF1 (iter=1.5)



towards EOF1 (iter=2.5)



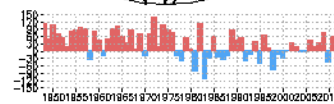
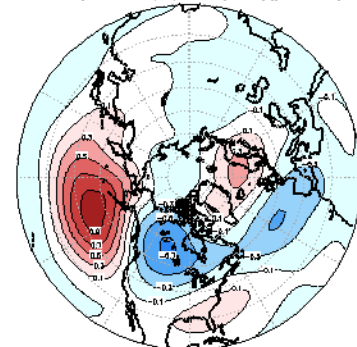
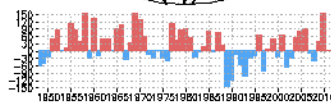
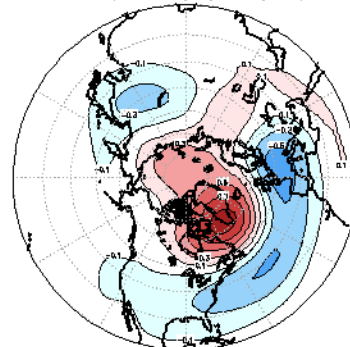
towards EOF1 (iter=3.5)



The route to EOF1 from a **ridiculous** guess field
Iterations 1-4

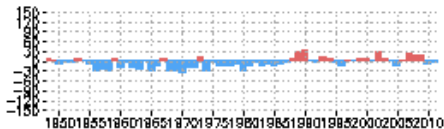
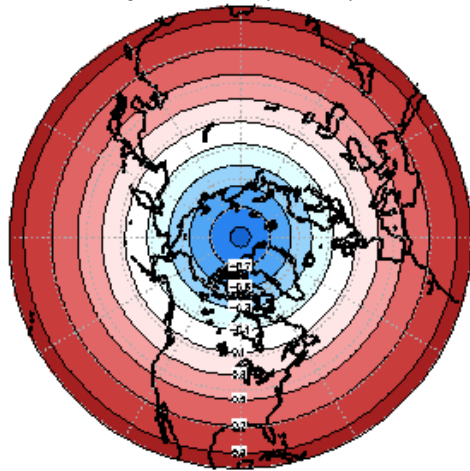
↓ Where
we go to

EOF for JFM 1948–2011 HGT 500 mb(iter)
EOF1 (23.0 %EV) (seed=guess) EOF2 (19.4 %EV) (seed=guess)(partial 1)

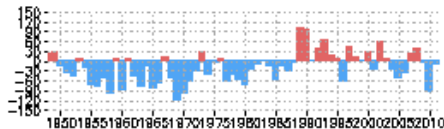
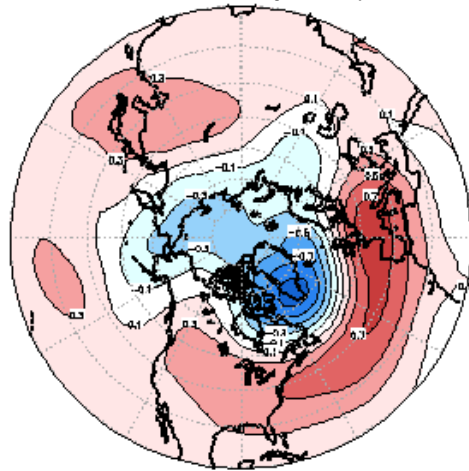


Path to EOF1 for JFM 1948–2011 Z500mb

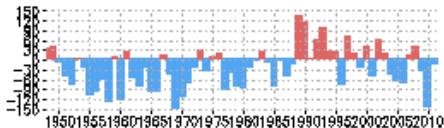
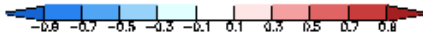
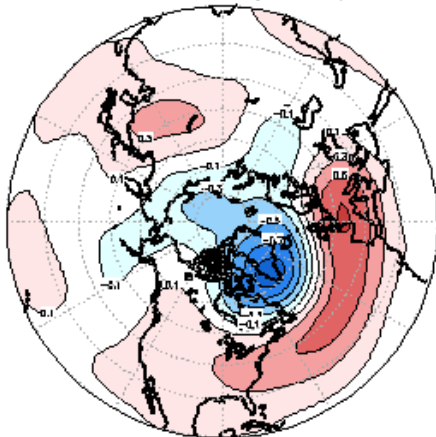
guess EOF1 (iter=0.5)



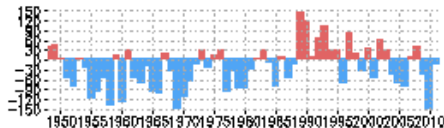
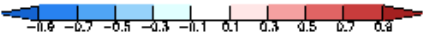
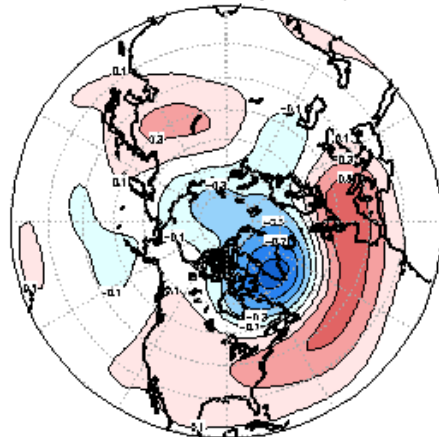
towards EOF1 (iter=1.5)



towards EOF1 (iter=2.5)



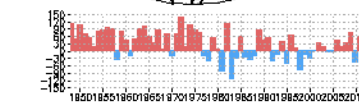
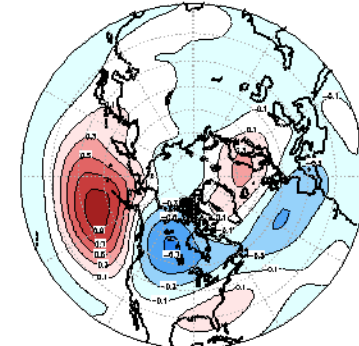
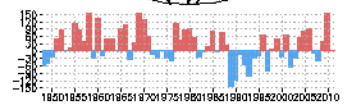
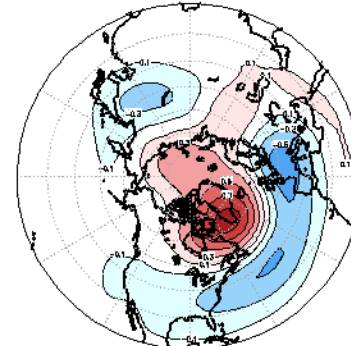
towards EOF1 (iter=3.5)



The route to EOF1 from a **another ridiculous** guess field
Iterations 1-4

↓ Where we go to

EOF for JFM 1948–2011 HGT 500 mb(iter)
EOF1 (23.0 %EV) (seed=guess) EOF2 (19.4 %EV) (seed=guess)(partial 1)



Savings in CPU?

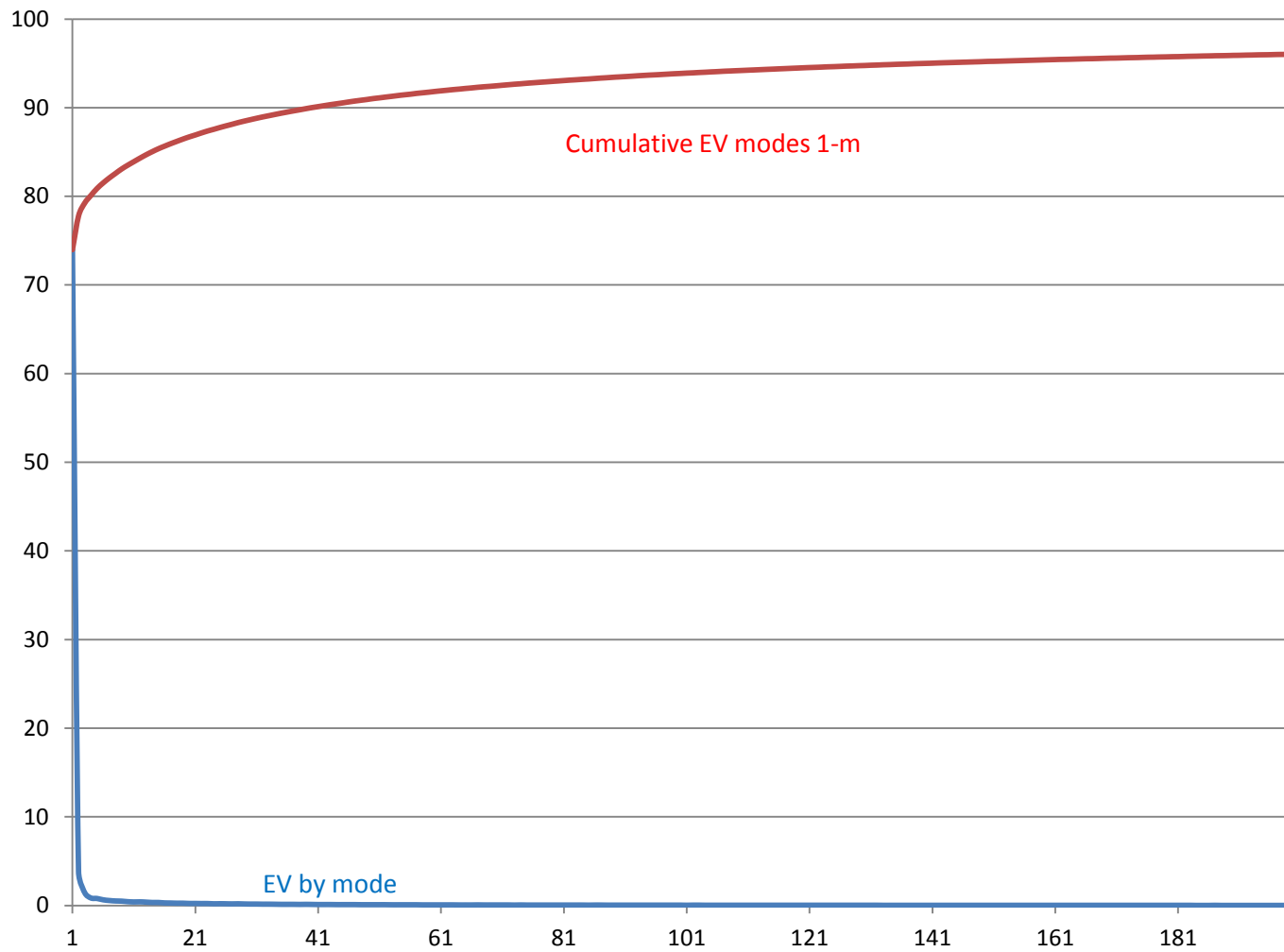
- **None** (to write home about) for a small problem
- Iteration is factor 30 faster for $n_s=n_t=4000$
- When n_s and n_t go higher (10,000 or 100,000) the cov matrix method quickly becomes 'impossible' (depending on computer and smarts of programmer), while iteration is still possible (although not inexpensive).

Other comments

- The calculation is simple (shorter code)
- No issues of space-time reversal apply
- You can start with a guess in space or in time
- Does it always work? (give me an example where it does NOT)
- Double precision issues

Application to highres CFSR daily data

- Something entirely impossible via the cov matrix can actually be done via the iteration method.
- Daily T2m 1979-2010. $nt=11688$, $ns=663552$
- It works well, but, but and but



21	0.23
22	0.22
23	0.22
24	0.19
25	0.2
26	0.19
27	0.18
28	0.19
29	0.17
30	0.16
31	0.16
32	0.15
33	0.15
34	0.14
35	0.13
36	0.13
37	0.13
38	0.12
39	0.13
40	0.12
41	0.11
42	0.11
43	0.11
44	0.1

Application to highres CFSR daily data

- Something entirely impossible via the cov matrix can actually be done via the iteration method.
- Daily T2m 1979-2010. $nt=11688$, $ns=663552$
- It works well, but, but and but
 - Iteration at low spatial res (Cost= $2 \times 30 \times 200$ units)
 - Convergence issues (or are these the real EOFs?)
 - Pulling teeth on T2m
 - Even with 200 modes only 96%EV (residual error 1.15K). Compressibility of data set not good.

As eigenvalues become smaller the spectrum becomes flatter.

One has to iterate more and more to achieve convergence to the next EOF.

The smaller the added EV the longer the iteration. Not a desirable situation.
Not (yet?) a good approach when 'all' modes are needed.

Compressibility

- Regardless of our ability to calculate EOFs, do EOFs compress data sets enough? An example for a presumably very compressible data set: Z500 (daily, 1979-2012, daily climo, iter=100).

	originl	reconstr	res	EV	mode#	
•	93.65	93.28	8.29	99.22	197	
•	93.65	93.28	8.23	99.23	198	
•	93.65	93.29	8.18	99.24	199	
•	93.65	93.29	8.12	99.25	200	(wrt daily climo)
•	5621.5	5621.5	8.14	100.00	200	(absolute)
•	148.05	147.83	8.14	99.70	200	(wrt annual mean)

Why does iteration work?

Understand the Power Method

- In applied math the power method is used to calculate the **largest** eigenvalue (and associated eigenvector) of matrix A .
- Bad press (method gives you only one mode) does not apply to us
- Why&How it works?

Power Method

- A guess $e^0(s)$ can always be written as a linear combination of projections onto the unknown EOFs: $e^0(s) = \sum_m \alpha_m e_m(s)$
- $Q e_m = \lambda_m e_m$ by definition
- Execution of $Q e^0$ (which is done for the power method) thus yields: $\sum_m \lambda_m \alpha_m e_m(s)$.
- Continue, k^{th} iteration: $Q e^k$. One can see the next highest λ emerge.
- Method is 'slow' (i.e. many iterations required, depending on the separation of next highest eigenvalues, 'luck' with initial projection...

- What is the connection between power method and the iteration? Acknowledge Chris Bretherton (Jan 2000).

$$e_m(s) = \sum_t \alpha_m(t) f(s, t) / \sum_t \alpha_m^2(t) \quad (1)$$

$$\alpha_m(t) = \sum_s e_m(s) f(s, t) / \sum_s e_m^2(s) \quad (2)$$

$$e_m(s) = \sum_t \alpha_m(t) f(s, t)$$

$$\alpha_m(t) = \sum_s e_m(s) f(s, t)$$

in essence:

$$\text{Iteration } k+1: \quad \mathbf{e}^{k+1} = \mathbf{f} \mathbf{f}^T \mathbf{e}^k$$

Baldwin, Mark P., David B. Stephenson, Ian T. Jolliffe, 2009: Spatial Weighting and **Iterative Projection Methods for EOFs**. *J. Climate*, **22**, 234–243

“Another method involving iteration between a spatial pattern and a time series was proposed by Clint and Jennings (1970).

Van den Dool et al. (2000) used a similar approach to find the leading EOF beginning from the leading empirical orthogonal teleconnection (EOT) pattern.

Iteration between a time series and spatial pattern to calculate the leading EOF was **discovered** independently by G. Hegerl (2008, personal communication).”

Class Notes Kalnay

Appendix (adapted from Wiki)

- Formally, the singular value decomposition of an $m \times n$ real matrix M is a factorization of the form $M = U \Sigma V^T$,
- where U is an $m \times m$ real matrix, Σ is an $m \times n$ diagonal matrix with nonnegative real numbers on the diagonal, and V^T (the conjugate transpose of V) is an $n \times n$ matrix. The diagonal entries $\Sigma_{i,i}$ of Σ are known as the singular values of M . The m columns of U and the n columns of V are called the left singular vectors and right singular vectors of M , respectively.
- Singular value decomposition and eigendecomposition are closely related. Namely:
 - The left singular vectors of M are eigenvectors of MM^T .
 - The right singular vectors of M are eigenvectors of $M^T M$.
 - The non-zero singular values of Σ are the square roots of the non-zero eigenvalues of $M^T M$ or MM^T .

Toumazou, Vincent, Jean-Francois Cretaux, 2001: Using a Lanczos Eigensolver in the Computation of Empirical Orthogonal Functions. *Mon. Wea. Rev.*, **129**, 1243–1250.

The computation of the singular values can be achieved through three different strategies in terms of linear algebra.

The first one is called the SVD strategy and is based on the singular value decomposition (SVD) algorithm.

The two others (QR and Lanczos) strategies are based on the formulation as an eigenvalue problem.

Acknowledge Michael Tippett at IRI

General conclusions

- Iteration methods allows calculation of a limited # of EOFs, where the covariance matrix method would nearly be impossible (large data sets)
- How limited?
- Using EOF for storage and transmission hindered by fundamental “compressibility” issues. (I would prefer # of modes $\leq 1\%$ of n_t to explain ‘enough’)

Other comments

- The method works for the same reason as the power method (when applied to the cov matrix)
- The calculation is simple (shorter code)
- No issues of space-time reversal apply
- You can start with a guess in space or in time
- Does it always work? (give me an example where it does NOT)
- A convergence criterion? (to save on iterations when not needed)
- What guess to use?
- Double precision issues
- Iterative EOT?
- Swarming a bunch of guess fields simultaneously (Tippett)