# Sunflower Oil x Envr - r2 (shared)

A useful and not overly complex analyses to get out the door would be an analysis of the wild sunflower seed oil traits and the environmental parameters that correlate with them. A great target journal and submission type would be to the American Journal of Botany as a brief communication (research article). Those types of article are 3000 to 4000 words in length, and have no more than 4 visual items (tables or figures). This article type is concise and gives enough room to explore one or two ideas in an manuscript.

I need to get permission to use there data, find out what environmental analyses were conducted in the original dissertation, and then start building models. It's possible there are more traits in there than oil that have the large latitudinal gradient, and it may be possible to conduct a novel analyses on those data.  However, just sticking to the oil and envr. traits would make for a tight brief communication.

I could get this sent out to review in two months for sure.
Authors: Karl Fetter, Max Barnhart (if he's interested), Ed McAssey, Andy Goeherty, John Burke.

Roles:
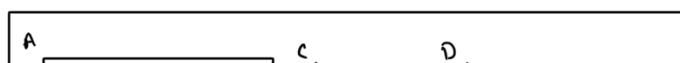KF - fit models and write the MS
MB - collate the data and create figure 1, contextual brainstorming for introduction, edit text
AG - create shapefile of wild sunflower range.
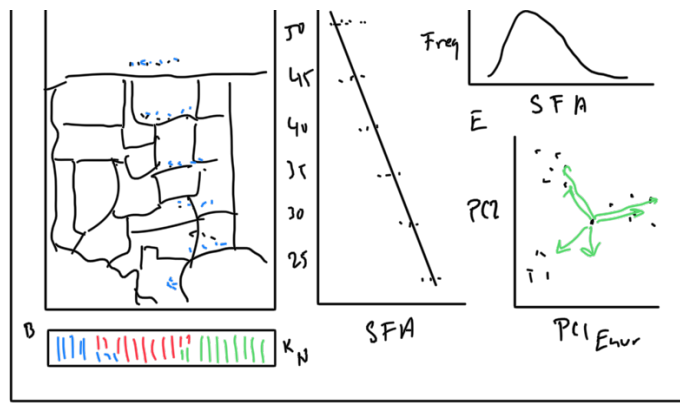JB - keep the lights on, edit MS

Fig. 1

Sample & Envr. Context Figure

Panels

A. Sample map

| A | | C. | D. |

B  Barplot (Adapted from Pub

C  Latitude observation

D  Histogram of trait

E  PCA on Envr.

## Notes:

- Color points According to demo.
- Adapt Structure / Admixture bar plot from original 2016 MS (here for convenience)
- (D) could include hist of one or two other interesting envr params. eg. GSL, Temp, Precip. Params that turn out to be important.  $^D \underset{X_1\ X_2}{\underline{|SFA}}$
- Include Latitude in PCA.

- Create convex hull of wild sunflower accessions from GBIF to serve as a range map for wild Helianthus.
  - Include Range as a shapefile & provide in supplement.
  - then when people use it, they will cite the manuscript.
- Possibly something Andy Doherty can do. He can probably do a really good job quickly.

## Input Data

$Y$ = % Saturated fatty Acids

$X$ = 19 Bioclim  } Abiotic -lat-longs
Growing Season length  } envr.
...

$G$ = Population Effects
Admixture Amt., Kinship matrix?
↳ From where?
Sometimes difficult to get.

Grouping Effects
Demo-code
Pop-code
possibly ind-code? or family-code ← Depends on data str.

## Model.

**Basic Idea:** $Y \sim E + G + \varepsilon$ $\qquad\qquad V_P = V_G + V_E + e$

- Bayesian model
- Or, use a Random Forest (RF) model.

1) $SFA \sim Anc + E + (1 | Deme)$  ← Grouping effect depend
$\qquad\qquad\qquad\qquad\qquad\quad$ ? $\qquad$ on how the data wer
$\qquad\qquad\qquad\qquad\qquad\quad$ Pop $\qquad$ collected.

$\qquad\qquad\downarrow \qquad\qquad\qquad\qquad\qquad\quad\downarrow$ - genetic control of Qmat

$\qquad$ Deme? $\qquad\qquad\qquad\qquad\qquad$ - Random effect of Pop.
$\qquad$ Admixture? $\qquad\qquad\quad$ Depends on data format &
$\qquad$ ... $\qquad\qquad\qquad\qquad\qquad$ my question. Need to think on
$\qquad\qquad\qquad\qquad\qquad\quad$ Then may be an interesting question

1) <u>Basic model:</u>

$\qquad Y \sim Anc. + E + \varepsilon$

if this can be done in Bayesian or <u>RF</u>, do it.
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\downarrow$
$\qquad\qquad\qquad\qquad\qquad\qquad$ • How to include genetic
$\qquad\qquad\qquad\qquad\qquad\qquad\quad$ Control w/ RF models?
$\qquad\qquad\qquad\qquad\qquad\qquad$ - Ask Dominik

<u>Fig 2</u>

Plot conditional effects from model 1



$\Big\}$ etc. • main takeaway figure.
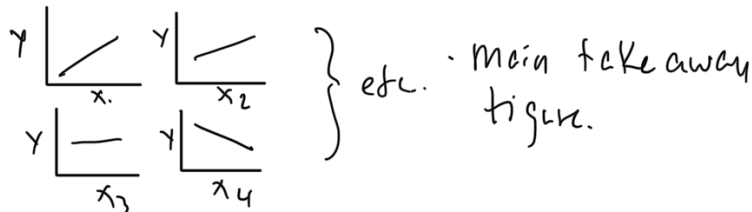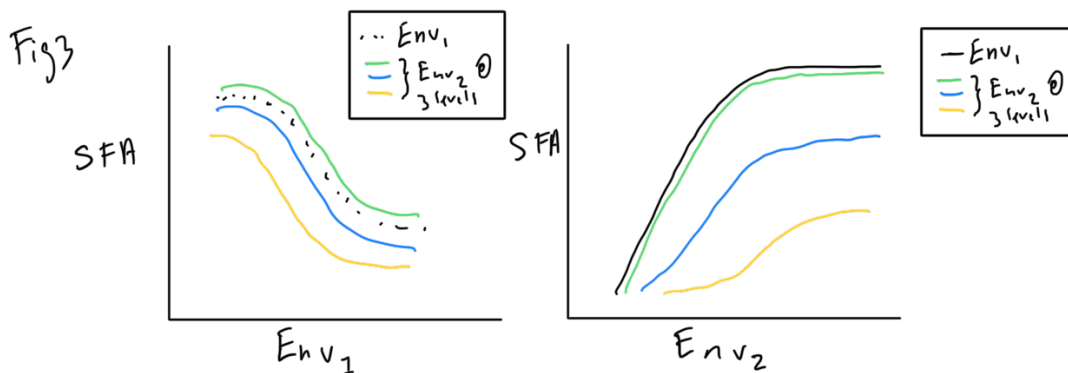
<u>Table 1</u>

model output of Basic model (if using Bayesian)

# Interactions of Climate

- There are certainly interesting interactions. The idea is to explore them in some depth. Do the interactions explain more variance than the single terms?

- I don't think I can fit a full interaction model. Too many parameters. I could fit an interaction model to the subset of params that explain most variance & explore them.

- If you take the Bayesian route. Plot interactions as terms or conditional effects.

Fig 3



- Some visual exploration of interactions

- Would be great if some interactions were negative!

$$y \frac{\boxed{\times}}{x_1} \begin{matrix} x_1 \\ x_2 \end{matrix} - \text{conditical effect}$$

## Bayesian - Only Approach

- Fit All x's & genetic controls

- Take top 3 variables & fit interaction model

## Visuals Outline

Fig 1. Sample & Exploration fig.

Fig 2. Main effects plot

Fig 3. Interactions plot

Table 1. main effects table + interaction effect table $\left(= \begin{array}{c} \text{Table 1 Model output.} \\ \text{mod 1} \\ \overline{\phantom{========}} \\ \overline{\phantom{========}} \\ \text{mod 2} \\ \overline{\phantom{========}} \\ \overline{\phantom{========}} \end{array}\right.$

## RF Approach

- Fit RF model to $X_j$ + generic controls
- Take top 3 vars & fit Bayesian interaction model

## Visuals Outline

Fig 1. Sample & Exploration fig.

Fig 2. RF output (RMSE + Cond. effts)

Fig 3. Interaction models

Table 1. Interaction model output

## Pros / Cons of RF vs Bayesian

- Not sure beyond superficial ideas.

- Read &/or talk to Dominik
- Need to fit the models first, then determine pro/cons

## Validation

Take wild seeds from populations that span the top envr. gradient. If the gradient is easy to replicate in a chamber, do so. This will act as a form of validation for the statistics & it connects to the mechanistic hypothesis that %SFA is under local selection for seed germination.