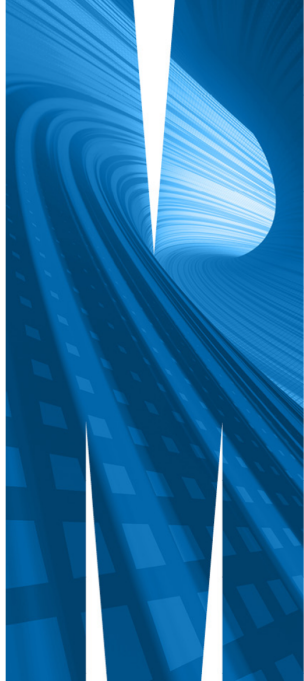


# **Forecast reconciliation: Geometry, optimization and beyond**

Anastasios Panagiotelis

2 July 2025



# Outline

- 1 Hierarchical Data and Forecast Reconciliation
- 2 Probabilistic Forecasts
- 3 Quantile Forecasting
- 4 Non-Linear Forecasting
- 5 Beyond Hierarchies
- 6 Wrap-up

# Outline

- 1 Hierarchical Data and Forecast Reconciliation
- 2 Probabilistic Forecasts
- 3 Quantile Forecasting
- 4 Non-Linear Forecasting
- 5 Beyond Hierarchies
- 6 Wrap-up

# Hierarchical Time Series

- At its most general, **multivariate** data  $\mathbf{y} \in \mathbb{R}^n$  bound together by some constraints.

# Hierarchical Time Series

- At its most general, **multivariate** data  $\mathbf{y} \in \mathbb{R}^n$  bound together by some constraints.
- Typically these constraints are **linear**, although later I will present new work for non-linear constraints.

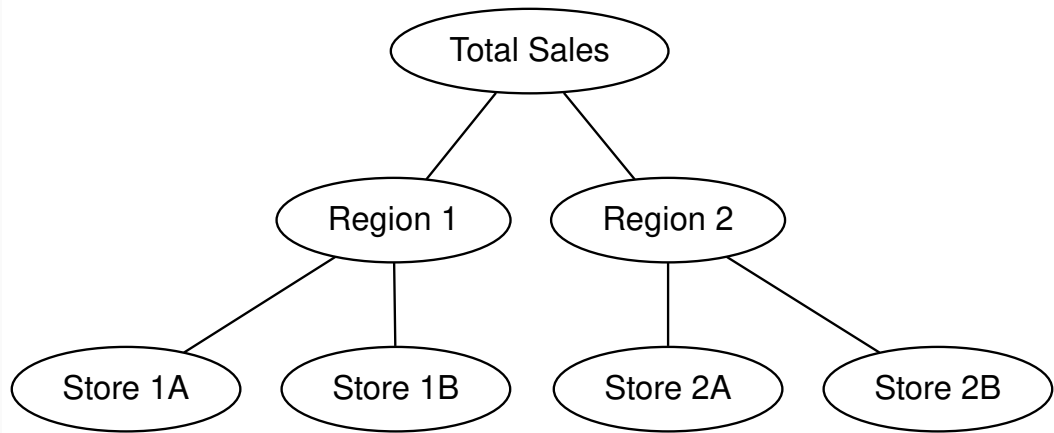
# Hierarchical Time Series

- At its most general, **multivariate** data  $\mathbf{y} \in \mathbb{R}^n$  bound together by some constraints.
- Typically these constraints are **linear**, although later I will present new work for non-linear constraints.
- Most commonly arise due to an **aggregation** structure, hence the name 'hierarchical'.

# Hierarchical Time Series

- At its most general, **multivariate** data  $\mathbf{y} \in \mathbb{R}^n$  bound together by some constraints.
- Typically these constraints are **linear**, although later I will present new work for non-linear constraints.
- Most commonly arise due to an **aggregation** structure, hence the name 'hierarchical'.
- Need not be hierarchical, alternative structures are grouped (or crossed) aggregation, or temporal aggregation.

# Hierarchy





# One representation

- For the simple hierarchy shown earlier:

$$\begin{pmatrix} y_{\text{Tot}} \\ y_1 \\ y_2 \\ y_{1A} \\ y_{1B} \\ y_{2A} \\ y_{2B} \end{pmatrix} = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \times \begin{pmatrix} y_{1A} \\ y_{1B} \\ y_{2A} \\ y_{2B} \end{pmatrix}$$

**y**      =      **S**      ×      **b**

# Forecasting

- Observations will always **cohere** to the constraints.

# Forecasting

- Observations will always **cohere** to the constraints.
- Forecasts generally will not.

# Forecasting

- Observations will always **cohere** to the constraints.
- Forecasts generally will not.
  - ▶ Different forecasts are made by different agents.

# Forecasting

- Observations will always **cohere** to the constraints.
- Forecasts generally will not.
  - ▶ Different forecasts are made by different agents.
  - ▶ Hard to construct a method that guarantees coherence.

# Forecasting

- Observations will always **cohere** to the constraints.
- Forecasts generally will not.
  - ▶ Different forecasts are made by different agents.
  - ▶ Hard to construct a method that guarantees coherence.
- This talk is about **two-stage** processes whereby incoherent **base** forecasts are adjusted to be coherent.

# Forecasting

- Observations will always **cohere** to the constraints.
- Forecasts generally will not.
  - ▶ Different forecasts are made by different agents.
  - ▶ Hard to construct a method that guarantees coherence.
- This talk is about **two-stage** processes whereby incoherent **base** forecasts are adjusted to be coherent.
- Note there is also work on **end-to-end** forecasting (e.g. Rangapuram et al. 2021).

# Forecast reconciliation

- Start with a vector of incoherent forecasts  $\hat{\mathbf{y}}$ .



# Forecast reconciliation

- Start with a vector of incoherent forecasts  $\hat{\mathbf{y}}$ .
- Recall the characterization earlier with  $\mathbf{S}$

# Forecast reconciliation

- Start with a vector of incoherent forecasts  $\hat{\mathbf{y}}$ .
- Recall the characterization earlier with  $\mathbf{S}$
- 'Regress'  $\hat{\mathbf{y}}$  on  $\mathbf{S}$ .

# Forecast reconciliation

- Start with a vector of incoherent forecasts  $\hat{\mathbf{y}}$ .
- Recall the characterization earlier with  $\mathbf{S}$
- 'Regress'  $\hat{\mathbf{y}}$  on  $\mathbf{S}$ .
- Use prediction as reconciled forecasts  $\tilde{\mathbf{y}}$ , i.e.

# Forecast reconciliation

- Start with a vector of incoherent forecasts  $\hat{\mathbf{y}}$ .
- Recall the characterization earlier with  $\mathbf{S}$
- 'Regress'  $\hat{\mathbf{y}}$  on  $\mathbf{S}$ .
- Use prediction as reconciled forecasts  $\tilde{\mathbf{y}}$ , i.e.

$$\tilde{\mathbf{y}} = \mathbf{S}(\mathbf{S}'\mathbf{S})^{-1}\mathbf{S}'\hat{\mathbf{y}}$$

- This is called OLS reconciliation.

# An optimization lens

- Ensure that  $\hat{\mathbf{y}}$  and  $\tilde{\mathbf{y}}$  are 'close'.

$$\tilde{\mathbf{y}} = \underset{\mathbf{y}}{\operatorname{argmin}} ||\mathbf{y} - \hat{\mathbf{y}}||_2$$

# An optimization lens

- Ensure that  $\hat{\mathbf{y}}$  and  $\tilde{\mathbf{y}}$  are 'close'.

$$\tilde{\mathbf{y}} = \underset{\mathbf{y}}{\operatorname{argmin}} ||\mathbf{y} - \hat{\mathbf{y}}||_2$$

- Under certain assumptions, this yields the same solution as the regression interpretation.

# An optimization lens

- Ensure that  $\hat{\mathbf{y}}$  and  $\tilde{\mathbf{y}}$  are 'close'.

$$\tilde{\mathbf{y}} = \underset{\mathbf{y}}{\operatorname{argmin}} \|\mathbf{y} - \hat{\mathbf{y}}\|_2$$

- Under certain assumptions, this yields the same solution as the regression interpretation.
- Also has a game theoretic interpretation (see van Erven and Cugliari 2015)

# Generalizations

- Where OLS works, it makes sense to consider GLS

$$\tilde{\mathbf{y}} = \mathbf{S}(\mathbf{S}'\mathbf{W}\mathbf{S})^{-1}\mathbf{S}'\mathbf{W}\hat{\mathbf{y}}$$



# Generalizations

- Where OLS works, it makes sense to consider GLS

$$\tilde{\mathbf{y}} = \mathbf{S}(\mathbf{S}'\mathbf{W}\mathbf{S})^{-1}\mathbf{S}'\mathbf{W}\hat{\mathbf{y}}$$

- Setting  $\mathbf{W}^{-1}$  to the covariance matrix of  $\mathbf{y} - \hat{\mathbf{y}}$  optimizes expected squared error loss.

# Generalizations

- Where OLS works, it makes sense to consider GLS

$$\tilde{\mathbf{y}} = \mathbf{S}(\mathbf{S}'\mathbf{W}\mathbf{S})^{-1}\mathbf{S}'\mathbf{W}\hat{\mathbf{y}}$$

- Setting  $\mathbf{W}^{-1}$  to the covariance matrix of  $\mathbf{y} - \hat{\mathbf{y}}$  optimizes expected squared error loss.
- This is the well-known **MinT method** of Wickramasuriya, Athanasopoulos, and Hyndman (2019).

# A geometric view

- Another alternative is a geometric view.

# A geometric view

- Another alternative is a geometric view.
- The constraints define a linear subspace  $\mathcal{S} \subset \mathbb{R}^n$ .

# A geometric view

- Another alternative is a geometric view.
- The constraints define a linear subspace  $\mathcal{S} \subset \mathbb{R}^n$ .
  - ▶ Sometimes  $\mathfrak{s}$  notation is used.

# A geometric view

- Another alternative is a geometric view.
- The constraints define a linear subspace  $\mathcal{S} \subset \mathbb{R}^n$ .
  - ▶ Sometimes  $\mathfrak{s}$  notation is used.
- This can easily be extended to other settings

# A geometric view

- Another alternative is a geometric view.
- The constraints define a linear subspace  $\mathcal{S} \subset \mathbb{R}^n$ .
  - ▶ Sometimes  $\mathfrak{s}$  notation is used.
- This can easily be extended to other settings
  - ▶ For discrete data  $\mathcal{S}$  is a set of points (Zhang, Panagiotelis, and Kang 2023).

# A geometric view

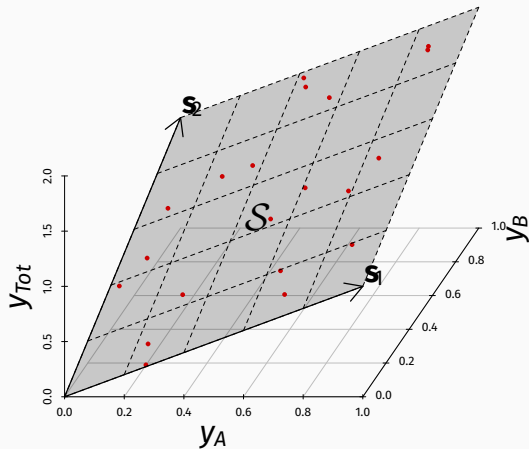
- Another alternative is a geometric view.
- The constraints define a linear subspace  $\mathcal{S} \subset \mathbb{R}^n$ .
  - ▶ Sometimes  $\mathfrak{s}$  notation is used.
- This can easily be extended to other settings
  - ▶ For discrete data  $\mathcal{S}$  is a set of points (Zhang, Panagiotelis, and Kang 2023).
  - ▶ Non linear constraints will be covered later.



# A geometric view

- Another alternative is a geometric view.
- The constraints define a linear subspace  $\mathcal{S} \subset \mathbb{R}^n$ .
  - ▶ Sometimes  $\mathfrak{s}$  notation is used.
- This can easily be extended to other settings
  - ▶ For discrete data  $\mathcal{S}$  is a set of points (Zhang, Panagiotelis, and Kang 2023).
  - ▶ Non linear constraints will be covered later.
- The simplest three-variable hierarchy  $y_{\text{Tot}} = y_A + y_B$  for real-valued data is depicted on the next slide.

# Coherent subspace



# Geometry of reconciliation

- Reconciliation is a **map** from  $\psi : \mathbb{R}^n \rightarrow \mathcal{S}$ .

# Geometry of reconciliation

- Reconciliation is a **map** from  $\psi : \mathbb{R}^n \rightarrow \mathcal{S}$ .
- These maps can be projections e.g. MinT

# Geometry of reconciliation

- Reconciliation is a **map** from  $\psi : \mathbb{R}^n \rightarrow \mathcal{S}$ .
- These maps can be projections e.g. MinT
- They can also be more general.

# Geometry of reconciliation

- Reconciliation is a **map** from  $\psi : \mathbb{R}^n \rightarrow \mathcal{S}$ .
- These maps can be projections e.g. MinT
- They can also be more general.
- For example, later we use  $\psi$  of the form

$$\tilde{\mathbf{y}} = \mathbf{S} (\mathbf{d} + \mathbf{G}\hat{\mathbf{y}})$$

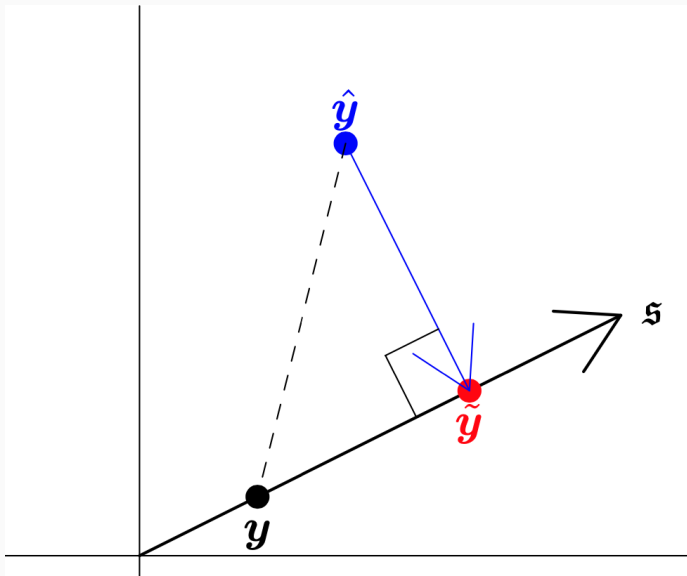
# Geometry of reconciliation

- Reconciliation is a **map** from  $\psi : \mathbb{R}^n \rightarrow \mathcal{S}$ .
- These maps can be projections e.g. MinT
- They can also be more general.
- For example, later we use  $\psi$  of the form

$$\tilde{\mathbf{y}} = \mathbf{S} (\mathbf{d} + \mathbf{G}\hat{\mathbf{y}})$$

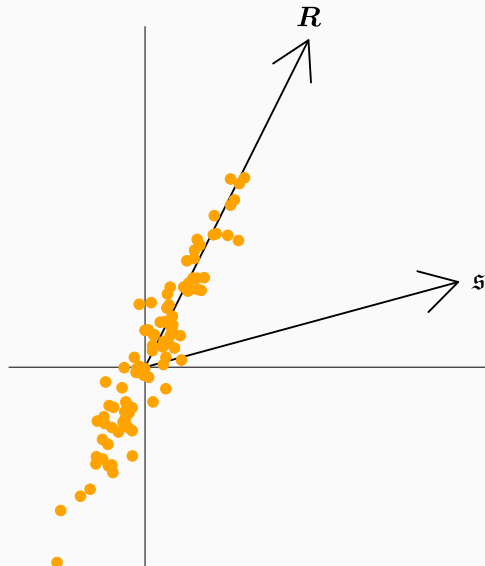
- On the next few slides we will depict projections.

# OLS Reconciliation

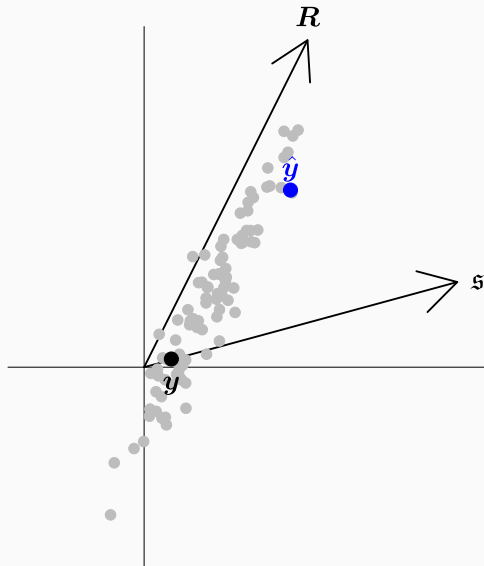




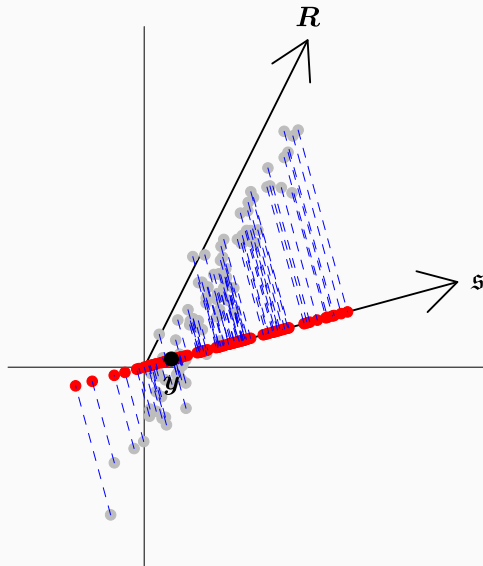
# Why MinT?



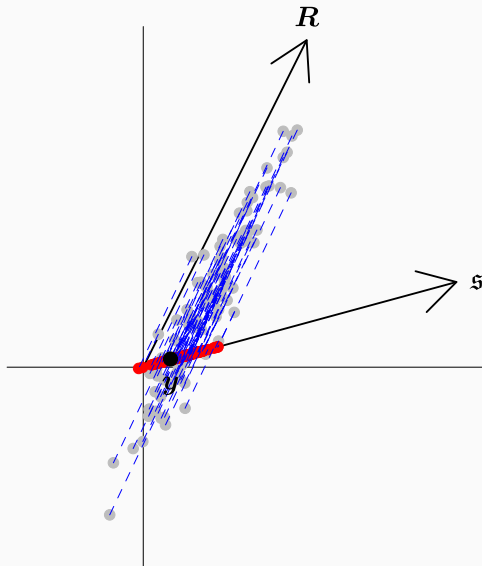
# Why MinT?



# Why MinT?



# Why MinT?



# Outline

- 1 Hierarchical Data and Forecast Reconciliation
- 2 Probabilistic Forecasts
- 3 Quantile Forecasting
- 4 Non-Linear Forecasting
- 5 Beyond Hierarchies
- 6 Wrap-up

# Problem

- The regression interpretation does not naturally lend itself to be extended to probabilistic forecasting.

# Problem

- The regression interpretation does not naturally lend itself to be extended to probabilistic forecasting.
- Alternative approaches define the reconciled distribution using *copulas* (Ben Taieb, Taylor, and Hyndman 2021) or by *conditioning* (see Corani et al. 2021).

# Problem

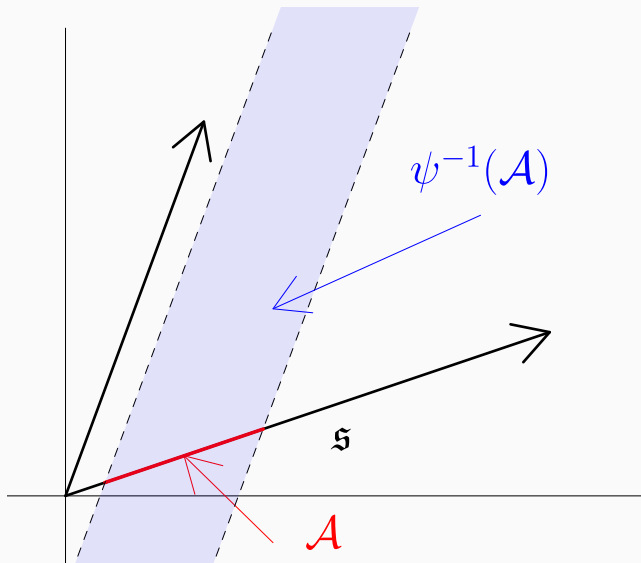
- The regression interpretation does not naturally lend itself to be extended to probabilistic forecasting.
- Alternative approaches define the reconciled distribution using *copulas* (Ben Taieb, Taylor, and Hyndman 2021) or by *conditioning* (see Corani et al. 2021).
- Some notions of reconciling draws from probabilistic distributions (Jeon, Panagiotelis, and Petropoulos 2019).



# Problem

- The regression interpretation does not naturally lend itself to be extended to probabilistic forecasting.
- Alternative approaches define the reconciled distribution using *copulas* (Ben Taieb, Taylor, and Hyndman 2021) or by *conditioning* (see Corani et al. 2021).
- Some notions of reconciling draws from probabilistic distributions (Jeon, Panagiotelis, and Petropoulos 2019).
- Later formalized reconciliation as a **pushforward** (Panagiotelis et al. 2023).

# Probabilistic Reconciliation



# Probabilistic Reconciliation

- Let  $\hat{\mu}$  be a measure on the usual  $\sigma$ -algebra defined on  $\mathbb{R}^n$ .

# Probabilistic Reconciliation

- Let  $\hat{\mu}$  be a measure on the usual  $\sigma$ -algebra defined on  $\mathbb{R}^n$ .
- Let  $\mathcal{A}$  be some region entirely within  $\mathcal{S}$  and  $\psi : \mathbb{R}^n \rightarrow \mathcal{S}$ .

# Probabilistic Reconciliation

- Let  $\hat{\mu}$  be a measure on the usual  $\sigma$ -algebra defined on  $\mathbb{R}^n$ .
- Let  $\mathcal{A}$  be some region entirely within  $\mathcal{S}$  and  $\psi : \mathbb{R}^n \rightarrow \mathcal{S}$ .
- Let  $\mathcal{B}$  be the pre-image of  $\mathcal{A}$ , denoted  $\psi^{-1}(\mathcal{A})$

# Probabilistic Reconciliation

- Let  $\hat{\mu}$  be a measure on the usual  $\sigma$ -algebra defined on  $\mathbb{R}^n$ .
- Let  $\mathcal{A}$  be some region entirely within  $\mathcal{S}$  and  $\psi : \mathbb{R}^n \rightarrow \mathcal{S}$ .
- Let  $\mathcal{B}$  be the pre-image of  $\mathcal{A}$ , denoted  $\psi^{-1}(\mathcal{A})$ 
  - ▶  $\forall \mathbf{x} \in \mathcal{B}, \psi(\mathbf{x}) \in \mathcal{A}$ .

# Probabilistic Reconciliation

- Let  $\hat{\mu}$  be a measure on the usual  $\sigma$ -algebra defined on  $\mathbb{R}^n$ .
- Let  $\mathcal{A}$  be some region entirely within  $\mathcal{S}$  and  $\psi : \mathbb{R}^n \rightarrow \mathcal{S}$ .
- Let  $\mathcal{B}$  be the pre-image of  $\mathcal{A}$ , denoted  $\psi^{-1}(\mathcal{A})$ 
  - ▶  $\forall \mathbf{x} \in \mathcal{B}, \psi(\mathbf{x}) \in \mathcal{A}$ .
- The reconciled measure  $\tilde{\mu}$  is defined as

# Probabilistic Reconciliation

- Let  $\hat{\mu}$  be a measure on the usual  $\sigma$ -algebra defined on  $\mathbb{R}^n$ .
- Let  $\mathcal{A}$  be some region entirely within  $\mathcal{S}$  and  $\psi : \mathbb{R}^n \rightarrow \mathcal{S}$ .
- Let  $\mathcal{B}$  be the pre-image of  $\mathcal{A}$ , denoted  $\psi^{-1}(\mathcal{A})$ 
  - ▶  $\forall \mathbf{x} \in \mathcal{B}, \psi(\mathbf{x}) \in \mathcal{A}$ .
- The reconciled measure  $\tilde{\mu}$  is defined as

$$\tilde{\mu}(\mathcal{A}) = \hat{\mu}(\mathcal{B})$$

- $\tilde{\mu}$  is the **pushforward** of  $\hat{\mu}$  by  $\psi$ , denoted as  $\psi\#\hat{\mu}$



# Optimality

- This merely defines a way of getting a reconciled distribution from some incoherent base distribution.

# Optimality

- This merely defines a way of getting a reconciled distribution from some incoherent base distribution.
- What is the optimal  $\psi$ ?

# Optimality

- This merely defines a way of getting a reconciled distribution from some incoherent base distribution.
- What is the optimal  $\psi$ ?
- In the point forecasting world the  $\psi$  given by MinT is optimal for squared loss.

# Optimality

- This merely defines a way of getting a reconciled distribution from some incoherent base distribution.
- What is the optimal  $\psi$ ?
- In the point forecasting world the  $\psi$  given by MinT is optimal for squared loss.
- What does optimality even mean for distributional forecasts?

# Scoring rules

- A scoring rule  $S : \mathcal{P} \times \mathbb{R}^n \rightarrow \mathbb{R}$  takes a distributional forecast from a family  $\mathcal{P}$  and an observation and assigns a *score* that measures forecast quality.

# Scoring rules

- A scoring rule  $S : \mathcal{P} \times \mathbb{R}^n \rightarrow \mathbb{R}$  takes a distributional forecast from a family  $\mathcal{P}$  and an observation and assigns a *score* that measures forecast quality.
- In the context of probabilistic forecasting, it has been proven that choosing  $\psi$  to be the same projection as MinT is optimal for log score when probabilistic forecasts are Gaussian (Wickramasuriya 2023)

# Scoring rules

- A scoring rule  $S : \mathcal{P} \times \mathbb{R}^n \rightarrow \mathbb{R}$  takes a distributional forecast from a family  $\mathcal{P}$  and an observation and assigns a *score* that measures forecast quality.
- In the context of probabilistic forecasting, it has been proven that choosing  $\psi$  to be the same projection as MinT is optimal for log score when probabilistic forecasts are Gaussian (Wickramasuriya 2023)
- In other cases we can optimize using a data driven approach.

# Score Optimization

- Obtain a sequence of base forecasts  $\hat{\mu}_t$  and corresponding realizations  $\mathbf{y}_t$ .



# Score Optimization

- Obtain a sequence of base forecasts  $\hat{\mu}_t$  and corresponding realizations  $\mathbf{y}_t$ .
- Parameterize  $\psi$  by some parameters  $\theta$  (denoted  $\psi_\theta$ )

# Score Optimization

- Obtain a sequence of base forecasts  $\hat{\mu}_t$  and corresponding realizations  $\mathbf{y}_t$ .
- Parameterize  $\psi$  by some parameters  $\theta$  (denoted  $\psi_\theta$ )
- For example,  $\psi(\mathbf{y}) = \mathbf{S}(\mathbf{d} + \mathbf{G}\mathbf{y})$ ,  $\theta = (\mathbf{d}', \text{vec}(\mathbf{G})')'$

# Score Optimization

- Obtain a sequence of base forecasts  $\hat{\mu}_t$  and corresponding realizations  $\mathbf{y}_t$ .
- Parameterize  $\psi$  by some parameters  $\theta$  (denoted  $\psi_\theta$ )
- For example,  $\psi(\mathbf{y}) = \mathbf{S}(\mathbf{d} + \mathbf{G}\mathbf{y})$ ,  $\theta = (\mathbf{d}', \text{vec}(\mathbf{G})')'$
- Optimize the following

# Score Optimization

- Obtain a sequence of base forecasts  $\hat{\mu}_t$  and corresponding realizations  $\mathbf{y}_t$ .
- Parameterize  $\psi$  by some parameters  $\theta$  (denoted  $\psi_\theta$ )
- For example,  $\psi(\mathbf{y}) = \mathbf{S}(\mathbf{d} + \mathbf{G}\mathbf{y})$ ,  $\theta = (\mathbf{d}', \text{vec}(\mathbf{G})')'$
- Optimize the following

$$\underset{\theta}{\operatorname{argmin}} \sum_t S(\psi_\theta \# \hat{\mu}_t, \mathbf{y}_t)$$

# Some practical points

- Scoring rules that have been used include log score, energy score and variogram score.

## Some practical points

- Scoring rules that have been used include log score, energy score and variogram score.
- Paired forecasts and observations can be obtained using rolling or expanding window schemes.

## Some practical points

- Scoring rules that have been used include log score, energy score and variogram score.
- Paired forecasts and observations can be obtained using rolling or expanding window schemes.
- Different optimal values of  $\theta$  can be obtained for different forecast horizons.

## Some practical points

- Scoring rules that have been used include log score, energy score and variogram score.
- Paired forecasts and observations can be obtained using rolling or expanding window schemes.
- Different optimal values of  $\theta$  can be obtained for different forecast horizons.
- Often draw a sample from  $\hat{\mu}$  rather than work with the distribution itself.



## Some practical points

- Scoring rules that have been used include log score, energy score and variogram score.
- Paired forecasts and observations can be obtained using rolling or expanding window schemes.
- Different optimal values of  $\theta$  can be obtained for different forecast horizons.
- Often draw a sample from  $\hat{\mu}$  rather than work with the distribution itself.
- Optimization by first order methods (e.g. SGD).

# Energy Generation Example

- Consider a moderate sized hierarchy (approx 20 variables) of electricity generation from different sources.

# Energy Generation Example

- Consider a moderate sized hierarchy (approx 20 variables) of electricity generation from different sources.
- Consider four different base forecasts

# Energy Generation Example

- Consider a moderate sized hierarchy (approx 20 variables) of electricity generation from different sources.
- Consider four different base forecasts
  - ▶ Assume Gaussianity or bootstrap

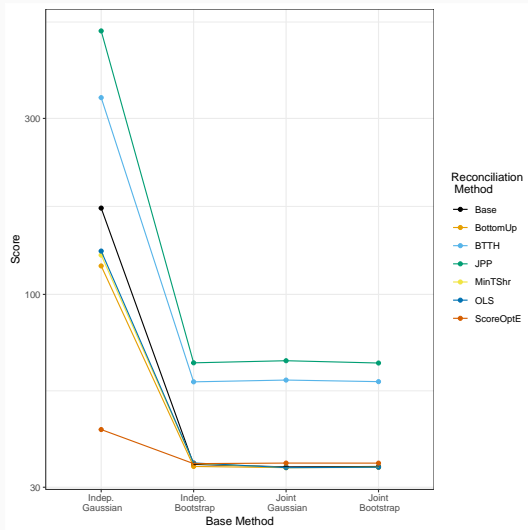
# Energy Generation Example

- Consider a moderate sized hierarchy (approx 20 variables) of electricity generation from different sources.
- Consider four different base forecasts
  - ▶ Assume Gaussianity or bootstrap
  - ▶ Assume independence or dependence

# Energy Generation Example

- Consider a moderate sized hierarchy (approx 20 variables) of electricity generation from different sources.
- Consider four different base forecasts
  - ▶ Assume Gaussianity or bootstrap
  - ▶ Assume independence or dependence
- Reconcile using projections (OLS, MinT) and also by optimising Energy score.

# Energy Generation Example



# Some thoughts

- Interplay between reconciliation and model misspecification.



# Some thoughts

- Interplay between reconciliation and model misspecification.
- Reconciliation via score optimization most effective when base models heavily misspecified.

# Some thoughts

- Interplay between reconciliation and model misspecification.
- Reconciliation via score optimization most effective when base models heavily misspecified.
- For reasonably well specified models

# Some thoughts

- Interplay between reconciliation and model misspecification.
- Reconciliation via score optimization most effective when base models heavily misspecified.
- For reasonably well specified models
  - ▶ Projections are robust

# Some thoughts

- Interplay between reconciliation and model misspecification.
- Reconciliation via score optimization most effective when base models heavily misspecified.
- For reasonably well specified models
  - ▶ Projections are robust
  - ▶ Gains over base forecasts are not as big

# Some thoughts

- Interplay between reconciliation and model misspecification.
- Reconciliation via score optimization most effective when base models heavily misspecified.
- For reasonably well specified models
  - ▶ Projections are robust
  - ▶ Gains over base forecasts are not as big
- Generalization of point forecast reconciliation to probabilistic setting gives forecaster a 'second chance'.

# Outline

- 1 Hierarchical Data and Forecast Reconciliation
- 2 Probabilistic Forecasts
- 3 Quantile Forecasting
- 4 Non-Linear Forecasting
- 5 Beyond Hierarchies
- 6 Wrap-up

# Pinball loss

- Many forecasting problems involve optimizing pinball loss.

$$L_{\alpha}(y, q) = \alpha(y_i - q)I(y_i \geq q) + (1 - \alpha)(q - y_i)I(y_i < q)$$

# Pinball loss

- Many forecasting problems involve optimizing pinball loss.

$$L_{\alpha}(y, q) = \alpha(y_i - q)I(y_i \geq q) + (1 - \alpha)(q - y_i)I(y_i < q)$$

- Here,  $I(.)$  equals 1 when the statement in parentheses is true, 0 otherwise.



# Pinball loss

- Many forecasting problems involve optimizing pinball loss.

$$L_{\alpha}(y, q) = \alpha(y_i - q)I(y_i \geq q) + (1 - \alpha)(q - y_i)I(y_i < q)$$

- Here,  $I(.)$  equals 1 when the statement in parentheses is true, 0 otherwise.
- Quantiles minimize expected pinball loss  $E_Y [L_{\alpha}(y, q)]$

# In reconciliation

- To target quantiles we optimize.

$$\underset{G}{\operatorname{argmin}} \sum_{i=1}^n \sum_{t \in \mathcal{T}_{\text{train}}} L_{\alpha}(y_{i,t}, \tilde{q}_{i,t})$$

# In reconciliation

- To target quantiles we optimize.

$$\underset{G}{\operatorname{argmin}} \sum_{i=1}^n \sum_{t \in \mathcal{T}_{\text{train}}} L_{\alpha}(y_{i,t}, \tilde{q}_{i,t})$$

- Subject to the constraints

$$\tilde{q}_{i,t} = \underset{q}{\operatorname{argmin}} E_{\tilde{y}_{i,t}} [L_{\alpha}(\tilde{y}_{i,t}, q)]$$

# Optimization

- This is an example of **bi-level optimization**.

# Optimization

- This is an example of **bi-level optimization**.
- It is further complicated by the fact that pinball loss is not smooth.

# Optimization

- This is an example of **bi-level optimization**.
- It is further complicated by the fact that pinball loss is not smooth.
- It is also complicated by the need to approximate expectations with sample equivalents

# Smooth pinball loss

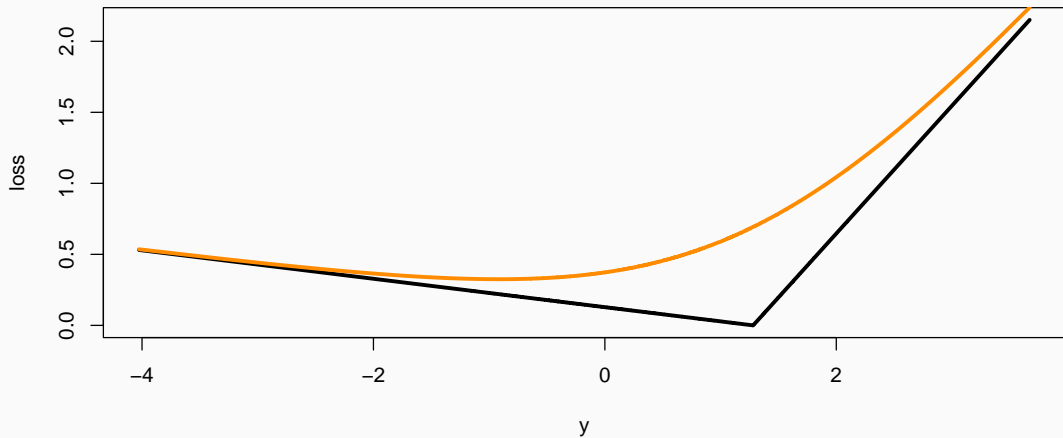
The following function approximates the pinball loss and converges to pinball loss as  $\beta \rightarrow \infty$

$$L_{\alpha}^{\beta}(y, q) = \frac{1}{\beta} \log \left( e^{\beta \alpha (y - q)} + e^{\beta (1 - \alpha) (q - y)} \right)$$

Unlike the pinball function it is smooth, meaning we can use first order methods (like Stochastic Gradient Descent).

# Smoothed pinball loss ( $\beta = 1$ )

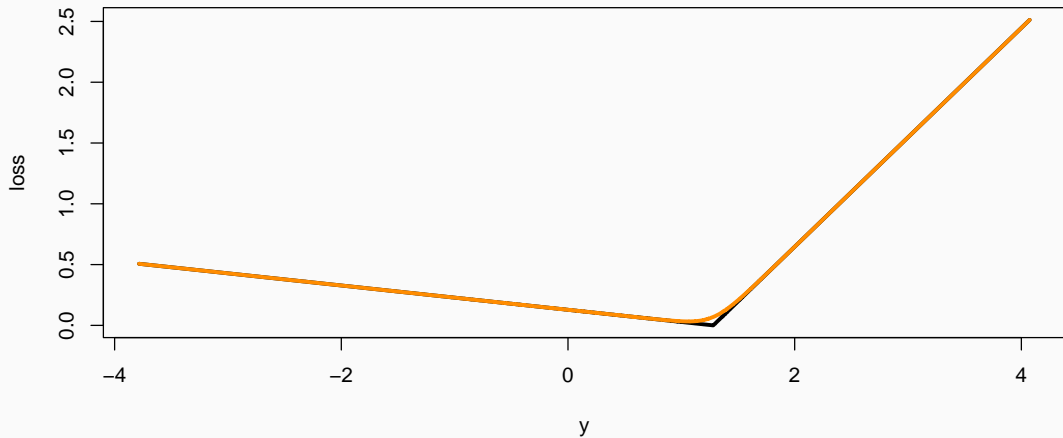
Pinball loss  $\alpha=0.9$ ,  $q=1.2816$





# Smoothed pinball loss ( $\beta = 10$ )

Pinball loss  $\alpha=0.9$ ,  $q=1.2816$



# Smoothed pinball loss ( $\beta = 100$ )

Pinball loss  $\alpha=0.9$ ,  $q=1.2816$



# What we want to solve

$$\underset{G}{\operatorname{argmin}} \sum_{i=1}^n \sum_{t \in \mathcal{T}_{\text{train}}} L_{\alpha}(y_{i,t}, \tilde{q}_{i,t})$$

- Subject to the constraints

$$\tilde{q}_{i,t} = \underset{q}{\operatorname{argmin}} E_{\tilde{Y}_{i,t}} [L_{\alpha}(\tilde{y}_{i,t}, q)]$$

# What we can solve

$$\underset{G}{\operatorname{argmin}} \sum_{i=1}^n \sum_{t \in \mathcal{T}_{\text{train}}} L_{\alpha}^{\beta}(\mathbf{y}_{i,t}, \tilde{\mathbf{q}}_{i,t})$$

subject to

$$\tilde{\mathbf{q}}_{i,t|t-1} = \underset{q}{\operatorname{argmin}} \sum_j L_{\alpha}^{\beta}(\tilde{\mathbf{y}}_{i,t}^{(j)}, q)$$

where  $\tilde{\mathbf{y}}_{i,t}^{(j)} = \psi_{\theta}(\hat{\mathbf{y}}_{i,t}^{(j)})$  and  $\hat{\mathbf{y}}_{i,t}^{(j)} \sim \hat{\mu}_t$  for  $j = 1, \dots, J$

# What has been proven

Let  $f(\mathbf{G})$  be the problem we want to solve,  $f^\beta(\mathbf{G})$  be the smooth approximation of pinball loss, and  $f^{(J)}$  the approximation from using  $J$  draws. We prove

$$\sup_{\theta \in \Theta} |f^\beta(\theta) - f(\theta)| \rightarrow 0 \text{ as } \beta \rightarrow \infty$$
$$\sup_{\theta \in \Theta} |f^{(J)}(\theta) - f^\beta(\theta)| \rightarrow 0 \text{ as } J \rightarrow \infty.$$

This implies the minimizer of the approximate problem converges to the minimizer of the ‘true’ problem.

# Convergence of SGD

- Optimization via SGD, taking care to pass gradient through argmin in lower level.

# Convergence of SGD

- Optimization via SGD, taking care to pass gradient through argmin in lower level.
- Note that the approximation of the expectation in the constraint means that the gradient is a biased estimate

# Convergence of SGD

- Optimization via SGD, taking care to pass gradient through argmin in lower level.
- Note that the approximation of the expectation in the constraint means that the gradient is a biased estimate
- However the variant of SGD we use will converge if



# Convergence of SGD

- Optimization via SGD, taking care to pass gradient through argmin in lower level.
- Note that the approximation of the expectation in the constraint means that the gradient is a biased estimate
- However the variant of SGD we use will converge if
  - ▶ Bounded second moment of the stochastic gradient

# Convergence of SGD

- Optimization via SGD, taking care to pass gradient through argmin in lower level.
- Note that the approximation of the expectation in the constraint means that the gradient is a biased estimate
- However the variant of SGD we use will converge if
  - ▶ Bounded second moment of the stochastic gradient
  - ▶  $L$ -smoothness

# Convergence of SGD

- Optimization via SGD, taking care to pass gradient through argmin in lower level.
- Note that the approximation of the expectation in the constraint means that the gradient is a biased estimate
- However the variant of SGD we use will converge if
  - ▶ Bounded second moment of the stochastic gradient
  - ▶  $L$ -smoothness
- Both are proven to hold for the functions we consider.

# Convergence of SGD

- Optimization via SGD, taking care to pass gradient through argmin in lower level.
- Note that the approximation of the expectation in the constraint means that the gradient is a biased estimate
- However the variant of SGD we use will converge if
  - ▶ Bounded second moment of the stochastic gradient
  - ▶  $L$ -smoothness
- Both are proven to hold for the functions we consider.
- Important to check convergence of SGD.

# Empirical study

- Use Australian tourism data.

# Empirical study

- Use Australian tourism data.
- Grouped hierarchy of states and purpose of travel.

# Empirical study

- Use Australian tourism data.
- Grouped hierarchy of states and purpose of travel.
- Dimension of **S** is  $40 \times 28$ .

# Empirical study

- Use Australian tourism data.
- Grouped hierarchy of states and purpose of travel.
- Dimension of **S** is  $40 \times 28$ .
- Seasonal ARIMA used for base forecasts. Distributional forecasts assume Gaussian errors and skew t errors.



# Empirical study

- Use Australian tourism data.
- Grouped hierarchy of states and purpose of travel.
- Dimension of **S** is  $40 \times 28$ .
- Seasonal ARIMA used for base forecasts. Distributional forecasts assume Gaussian errors and skew t errors.
- Train on 10 years (120 observations), evaluation on 7 years (84 observations).

## Pinball Loss - Out of Sample (Normal errors)

Method	Quantile Level			
	0.05	0.2	0.8	0.95
Base	32*	85*	101	46
OLS	32*	84*	104	51
WLS	<b>31*</b>	<b>82*</b>	112	65
MinT	<b>31*</b>	<b>82*</b>	111	65
QOpt	35*	85*	100*	<b>41*</b>

**Bold** denotes best performing method, asterisk(\*) denotes inclusion in model confidence set (Hansen et. al., 2011).

# Outline

- 1 Hierarchical Data and Forecast Reconciliation
- 2 Probabilistic Forecasts
- 3 Quantile Forecasting
- 4 Non-Linear Forecasting**
- 5 Beyond Hierarchies
- 6 Wrap-up

# The problem

- What if the constraints are non-linear?

# The problem

- What if the constraints are non-linear?
- For example ratios are common quantities of interest.

# The problem

- What if the constraints are non-linear?
- For example ratios are common quantities of interest.
  - ▶ Mortality rates are Deaths divided by Exposure.

# The problem

- What if the constraints are non-linear?
- For example ratios are common quantities of interest.
  - ▶ Mortality rates are Deaths divided by Exposure.
  - ▶ Unemployment rates are number of unemployed divided by labor force.

# The problem

- What if the constraints are non-linear?
- For example ratios are common quantities of interest.
  - ▶ Mortality rates are Deaths divided by Exposure.
  - ▶ Unemployment rates are number of unemployed divided by labor force.
- Both of these examples are also be subject to aggregation.



# Problem formulation

- In general there are  $C$  constraints  $g_1(\mathbf{y}) = 0, \dots, g_C(\mathbf{y}) = 0$ , or more compactly  $g(\mathbf{y}) = \mathbf{0}$ .

# Problem formulation

- In general there are  $C$  constraints  $g_1(\mathbf{y}) = 0, \dots, g_C(\mathbf{y}) = 0$ , or more compactly  $g(\mathbf{y}) = \mathbf{0}$ .
- The level set of points  $\mathbf{y} : g(\mathbf{y}) = \mathbf{0}$  defines a coherent *surface* or *manifold* continue to be denoted as  $\mathcal{S}$ .

# Problem formulation

- In general there are  $C$  constraints  $g_1(\mathbf{y}) = 0, \dots, g_C(\mathbf{y}) = 0$ , or more compactly  $g(\mathbf{y}) = \mathbf{0}$ .
- The level set of points  $\mathbf{y} : g(\mathbf{y}) = \mathbf{0}$  defines a coherent *surface* or *manifold* continue to be denoted as  $\mathcal{S}$ .
- Non-linear reconciliation solves the following problem:

# Problem formulation

- In general there are  $C$  constraints  $g_1(\mathbf{y}) = 0, \dots, g_C(\mathbf{y}) = 0$ , or more compactly  $g(\mathbf{y}) = \mathbf{0}$ .
- The level set of points  $\mathbf{y} : g(\mathbf{y}) = \mathbf{0}$  defines a coherent *surface* or *manifold* continue to be denoted as  $\mathcal{S}$ .
- Non-linear reconciliation solves the following problem:

$$\tilde{\mathbf{y}} = \underset{\mathbf{y}}{\operatorname{argmin}} (\mathbf{y} - \hat{\mathbf{y}})' \mathbf{W} (\mathbf{y} - \hat{\mathbf{y}})$$

- Subject to  $\mathbf{y} \in \mathcal{S}$ .

# Towards Theory

- Note focus is still on point forecasts.

# Towards Theory

- Note focus is still on point forecasts.
- First consider case of convex constraints.

# Towards Theory

- Note focus is still on point forecasts.
- First consider case of convex constraints.
  - ▶ Reconciliation guaranteed to improve base forecast, but only in hypograph.

# Towards Theory

- Note focus is still on point forecasts.
- First consider case of convex constraints.
  - ▶ Reconciliation guaranteed to improve base forecast, but only in hypograph.
- For more general constraints



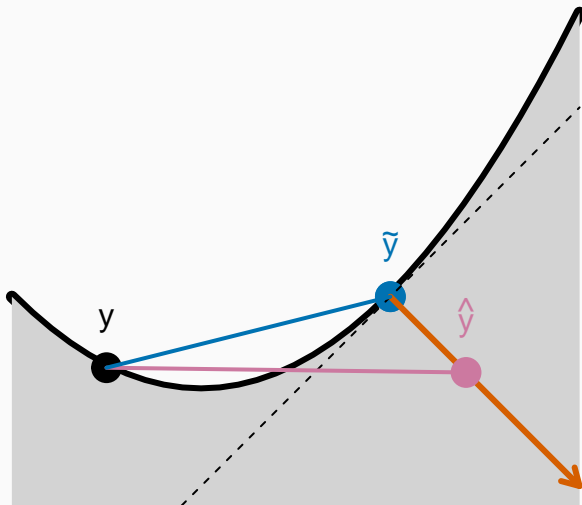
# Towards Theory

- Note focus is still on point forecasts.
- First consider case of convex constraints.
  - ▶ Reconciliation guaranteed to improve base forecast, but only in hypograph.
- For more general constraints
  - ▶ Find closest point on the coherent manifold equidistant from the base and reconciled forecast.

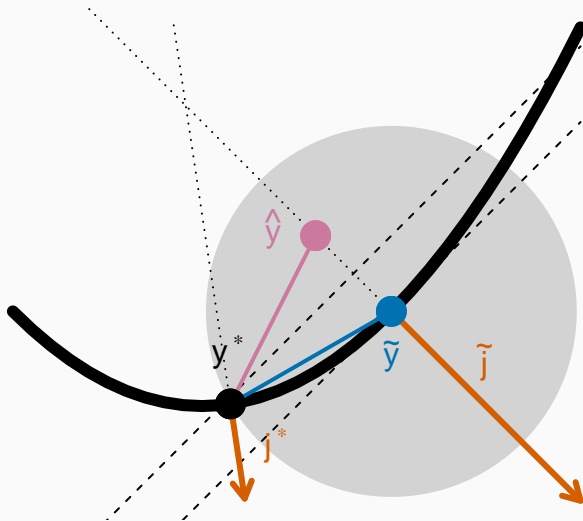
# Towards Theory

- Note focus is still on point forecasts.
- First consider case of convex constraints.
  - ▶ Reconciliation guaranteed to improve base forecast, but only in hypograph.
- For more general constraints
  - ▶ Find closest point on the coherent manifold equidistant from the base and reconciled forecast.
  - ▶ This defines a ball in which reconciliation always outperforms base forecasts.

# Convex Function: Hypograph



# Any function



# Radius of the Ball

- The radius of the ball on the previous slide is given by

$$r = \sqrt{\kappa' \mathbf{J}^{*'} \mathbf{J}^* \kappa + \mu \kappa' \mathbf{J}^{*'} \tilde{\mathbf{J}} \lambda + \frac{\mu^2}{4} \lambda \tilde{\mathbf{J}} \tilde{\mathbf{J}} \lambda}$$

# Radius of the Ball

- The radius of the ball on the previous slide is given by

$$r = \sqrt{\kappa' \mathbf{J}^{*'} \mathbf{J}^* \kappa + \mu \kappa' \mathbf{J}^{*'} \tilde{\mathbf{J}} \lambda + \frac{\mu^2}{4} \lambda \tilde{\mathbf{J}}' \tilde{\mathbf{J}} \lambda}$$

- $\mathbf{J}^*$  and  $\tilde{\mathbf{J}}$  are gradients of the constraint evaluated at  $\mathbf{y}^*$  and  $\tilde{\mathbf{y}}$  respectively.

# Radius of the Ball

- The radius of the ball on the previous slide is given by

$$r = \sqrt{\kappa' \mathbf{J}^{*'} \mathbf{J}^* \kappa + \mu \kappa' \mathbf{J}^{*'} \tilde{\mathbf{J}} \lambda + \frac{\mu^2}{4} \lambda \tilde{\mathbf{J}}' \tilde{\mathbf{J}} \lambda}$$

- $\mathbf{J}^*$  and  $\tilde{\mathbf{J}}$  are gradients of the constraint evaluated at  $\mathbf{y}^*$  and  $\tilde{\mathbf{y}}$  respectively.
- $\lambda$  and  $\kappa$  are Lagrange multipliers associated with certain optimization problems.

# How is this useful?

- This theory tells us that non-linear forecast reconciliation is more likely to succeed when



# How is this useful?

- This theory tells us that non-linear forecast reconciliation is more likely to succeed when
  - ▶ Base forecast is far from coherent manifold.

# How is this useful?

- This theory tells us that non-linear forecast reconciliation is more likely to succeed when
  - ▶ Base forecast is far from coherent manifold.
  - ▶ The constraint function has lower curvature.

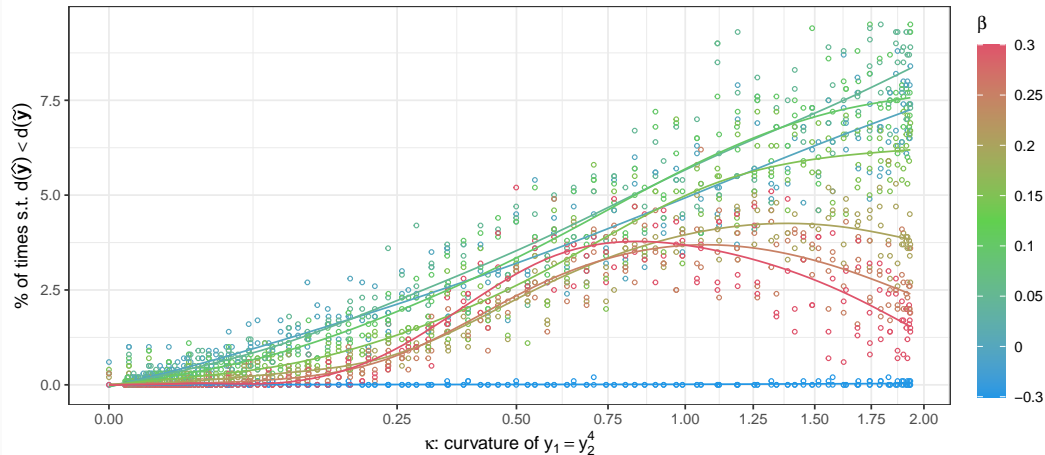
# How is this useful?

- This theory tells us that non-linear forecast reconciliation is more likely to succeed when
  - ▶ Base forecast is far from coherent manifold.
  - ▶ The constraint function has lower curvature.
  - ▶ When reconciled forecast is in a high probability region of the true DGP.

# How is this useful?

- This theory tells us that non-linear forecast reconciliation is more likely to succeed when
  - ▶ Base forecast is far from coherent manifold.
  - ▶ The constraint function has lower curvature.
  - ▶ When reconciled forecast is in a high probability region of the true DGP.
  - ▶ When some constraints are convex and the base forecast is more likely to lie in the hypographs of these constraints.

# Simulation results



# Mortality Data

- Annual (1969-2019) data on

# Mortality Data

- Annual (1969-2019) data on
  - ▶ Exposure ( $E$ )

# Mortality Data

- Annual (1969-2019) data on
  - ▶ Exposure ( $E$ )
  - ▶ Deaths ( $D$ )



# Mortality Data

- Annual (1969-2019) data on
  - ▶ Exposure ( $E$ )
  - ▶ Deaths ( $D$ )
  - ▶ Mortality rates ( $M$ )

# Mortality Data

- Annual (1969-2019) data on
  - ▶ Exposure ( $E$ )
  - ▶ Deaths ( $D$ )
  - ▶ Mortality rates ( $M$ )
- For US as a whole and 9 census regions.

# Mortality Data

- Annual (1969-2019) data on
  - ▶ Exposure ( $E$ )
  - ▶ Deaths ( $D$ )
  - ▶ Mortality rates ( $M$ )
- For US as a whole and 9 census regions.
- $E$  and  $D$  respect aggregation constraints.

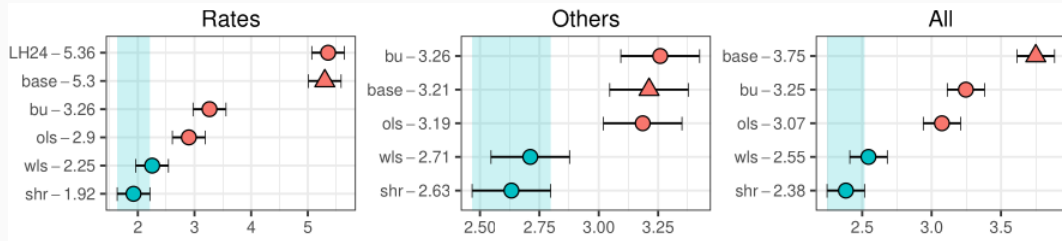
# Mortality Data

- Annual (1969-2019) data on
  - ▶ Exposure ( $E$ )
  - ▶ Deaths ( $D$ )
  - ▶ Mortality rates ( $M$ )
- For US as a whole and 9 census regions.
- $E$  and  $D$  respect aggregation constraints.
- $M$  need not respect hierarchical constraints.

# Mortality Data

- Annual (1969-2019) data on
  - ▶ Exposure ( $E$ )
  - ▶ Deaths ( $D$ )
  - ▶ Mortality rates ( $M$ )
- For US as a whole and 9 census regions.
- $E$  and  $D$  respect aggregation constraints.
- $M$  need not respect hierarchical constraints.
- However  $M = D/E$  for each region.

# Mortality Data



# Outline

- 1 Hierarchical Data and Forecast Reconciliation
- 2 Probabilistic Forecasts
- 3 Quantile Forecasting
- 4 Non-Linear Forecasting
- 5 Beyond Hierarchies**
- 6 Wrap-up

# Setup

- Suppose we are interested in multivariate forecasting but do not have linear (or non-linear) constraints.



# Setup

- Suppose we are interested in multivariate forecasting but do not have linear (or non-linear) constraints.
- Is there anything interesting about forecast reconciliation.

# Setup

- Suppose we are interested in multivariate forecasting but do not have linear (or non-linear) constraints.
- Is there anything interesting about forecast reconciliation.
- Surprisingly... Yes!

# Setup

- Suppose we are interested in multivariate forecasting but do not have linear (or non-linear) constraints.
- Is there anything interesting about forecast reconciliation.
- Surprisingly... Yes!
- New work on Forecast Linear Augmented Projects (FLAP)

# What is FLAP?

- Suppose the target is to forecast  $\mathbf{y}_t \in \mathbb{R}^m$ .

# What is FLAP?

- Suppose the target is to forecast  $\mathbf{y}_t \in \mathbb{R}^m$ .
- We construct new synthetic series  $\mathbf{c}_t \in \mathbb{R}^p$  where  $\mathbf{c}_t = \Phi \mathbf{y}_t$ .

# What is FLAP?

- Suppose the target is to forecast  $\mathbf{y}_t \in \mathbb{R}^m$ .
- We construct new synthetic series  $\mathbf{c}_t \in \mathbb{R}^p$  where  $\mathbf{c}_t = \Phi \mathbf{y}_t$ .
  - ▶ The choice of  $\Phi$  is arbitrary.

# What is FLAP?

- Suppose the target is to forecast  $\mathbf{y}_t \in \mathbb{R}^m$ .
- We construct new synthetic series  $\mathbf{c}_t \in \mathbb{R}^p$  where  $\mathbf{c}_t = \Phi \mathbf{y}_t$ .
  - ▶ The choice of  $\Phi$  is arbitrary.
- The augmented vector  $(\mathbf{c}'_t, \mathbf{y}'_t)'$  coheres to known linear constraints.

# What is FLAP?

- Suppose the target is to forecast  $\mathbf{y}_t \in \mathbb{R}^m$ .
- We construct new synthetic series  $\mathbf{c}_t \in \mathbb{R}^p$  where  $\mathbf{c}_t = \Phi \mathbf{y}_t$ .
  - ▶ The choice of  $\Phi$  is arbitrary.
- The augmented vector  $(\mathbf{c}'_t, \mathbf{y}'_t)'$  coheres to known linear constraints.
- Forecast all components of  $(\mathbf{c}'_t, \mathbf{y}'_t)'$



# What is FLAP?

- Suppose the target is to forecast  $\mathbf{y}_t \in \mathbb{R}^m$ .
- We construct new synthetic series  $\mathbf{c}_t \in \mathbb{R}^p$  where  $\mathbf{c}_t = \Phi \mathbf{y}_t$ .
  - ▶ The choice of  $\Phi$  is arbitrary.
- The augmented vector  $(\mathbf{c}'_t, \mathbf{y}'_t)'$  coheres to known linear constraints.
- Forecast all components of  $(\mathbf{c}'_t, \mathbf{y}'_t)'$
- Reconcile using MinT.

# The key idea

- It is known from the properties of MinT that we will reduce forecast variance for  $(\mathbf{c}'_t, \mathbf{y}'_t)'$ .

# The key idea

- It is known from the properties of MinT that we will reduce forecast variance for  $(\mathbf{c}'_t, \mathbf{y}'_t)'$ .
- We have proven that the same is true when only looking at  $\mathbf{y}_t$ .

# The key idea

- It is known from the properties of MinT that we will reduce forecast variance for  $(\mathbf{c}'_t, \mathbf{y}'_t)'$ .
- We have proven that the same is true when only looking at  $\mathbf{y}_t$ .
- It is this result that allows the benefits of reconciliation to be applied to problems where there are no constraints at all!

# The key idea

- It is known from the properties of MinT that we will reduce forecast variance for  $(\mathbf{c}'_t, \mathbf{y}'_t)'$ .
- We have proven that the same is true when only looking at  $\mathbf{y}_t$ .
- It is this result that allows the benefits of reconciliation to be applied to problems where there are no constraints at all!
- We also prove that the forecast variance is non-increasing as more synthetic components are added.

# No free lunch

- Originally FLAP stood for 'Free Lunch' augmented projection.

# No free lunch

- Originally FLAP stood for 'Free Lunch' augmented projection.
- All proofs assume error covariance matrix used in MinT in **known**. In practice it is estimated.

# No free lunch

- Originally FLAP stood for 'Free Lunch' augmented projection.
- All proofs assume error covariance matrix used in MinT is **known**. In practice it is estimated.
- The quality of covariance matrix estimates deteriorate with higher dimension.



# No free lunch

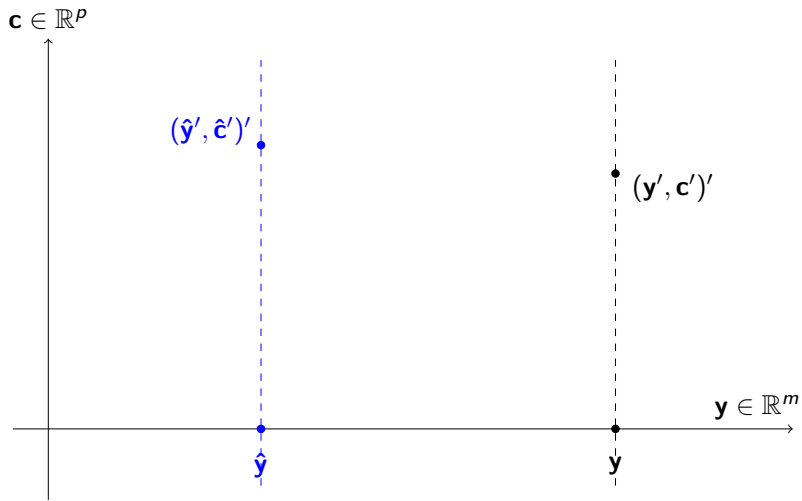
- Originally FLAP stood for 'Free Lunch' augmented projection.
- All proofs assume error covariance matrix used in MinT is **known**. In practice it is estimated.
- The quality of covariance matrix estimates deteriorate with higher dimension.
- However for finite dimension, the benefit of FLAP outweighs errors in estimating covariance matrix.

# Geometry of FLAP

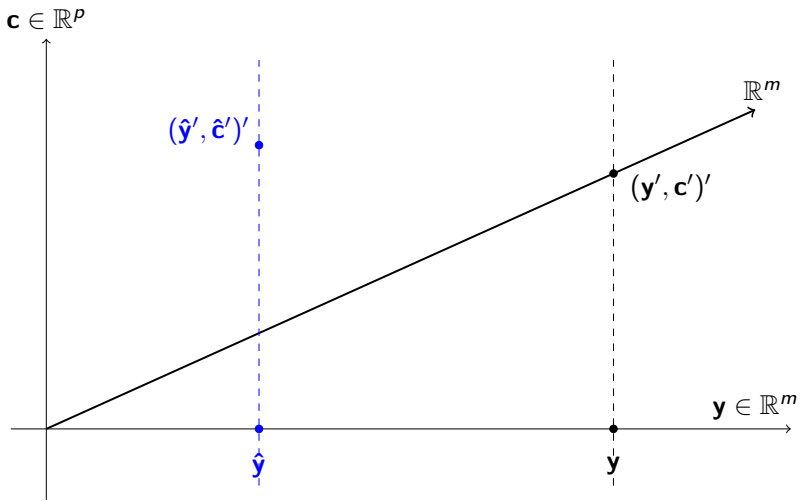
# Geometry of FLAP



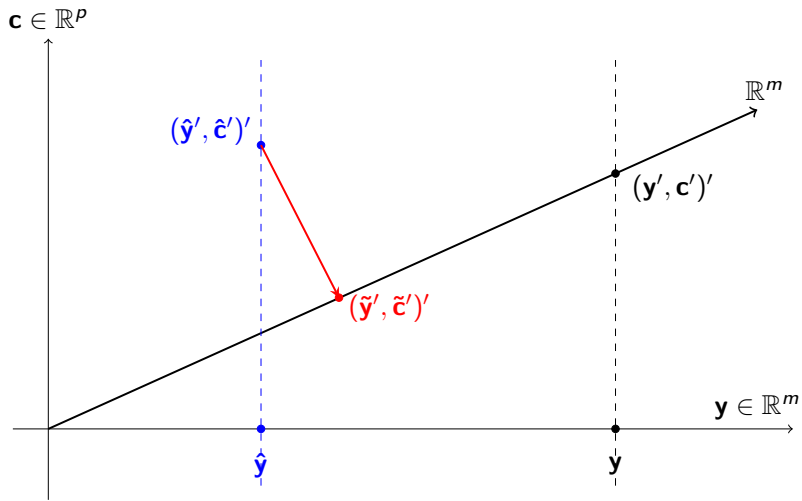
# Geometry of FLAP



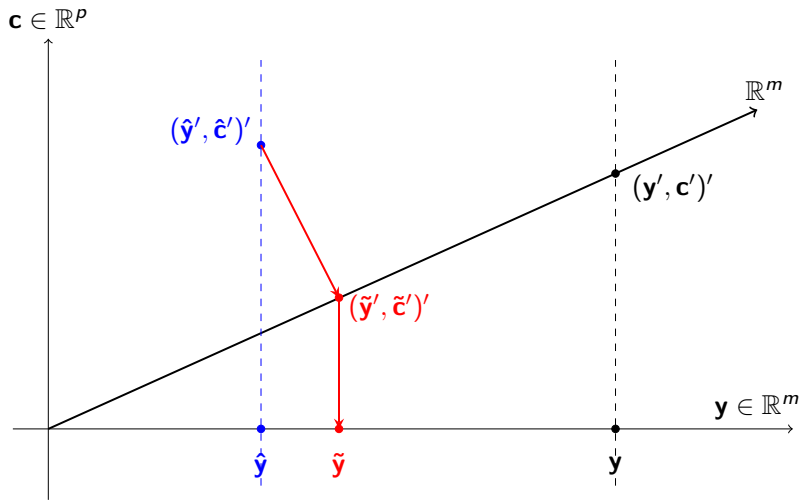
# Geometry of FLAP



# Geometry of FLAP



# Geometry of FLAP



- Monthly data of macroeconomic variables (McCracken and Ng, 2016).

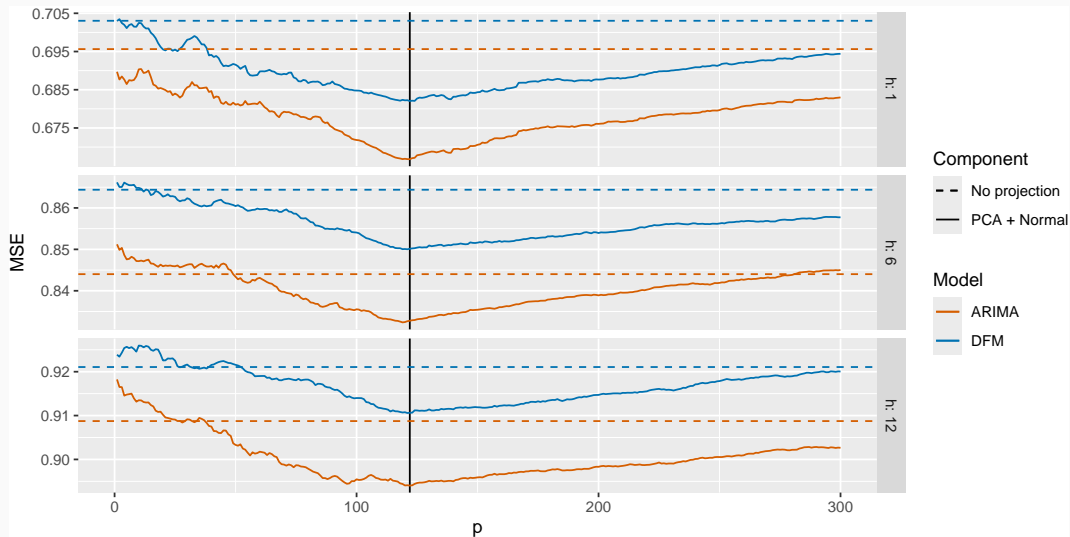


- Monthly data of macroeconomic variables (McCracken and Ng, 2016).
- Data from Jan 1959 – Sep 2023. 777 observations on 122 series.

- Monthly data of macroeconomic variables (McCracken and Ng, 2016).
- Data from Jan 1959 – Sep 2023. 777 observations on 122 series.
- Same cleaning process as per McCracken and Ng (2016).

- Monthly data of macroeconomic variables (McCracken and Ng, 2016).
- Data from Jan 1959 – Sep 2023. 777 observations on 122 series.
- Same cleaning process as per McCracken and Ng (2016).
- All series scaled to have mean 0 and variance 1.

- Monthly data of macroeconomic variables (McCracken and Ng, 2016).
- Data from Jan 1959 – Sep 2023. 777 observations on 122 series.
- Same cleaning process as per McCracken and Ng (2016).
- All series scaled to have mean 0 and variance 1.
- Expanding time series cross-validation with initial size of 25 years and forecast horizon 12 months.



# Working Paper and R Package

YF Yang, G Athanasopoulos, RJ Hyndman, and A Panagiotelis (2024). "Forecast Linear Augmented Projection (FLAP): A free lunch to reduce forecast error variance".

*Department of Econometrics and Business Statistics, Monash University, Working Paper Series 13/24.*

You can install the stable version from CRAN

```
## CRAN.R-project.org/package=flap  
install.packages("flap")
```

or the development version from Github

```
## github.com/FinYang/flap  
# install.packages("remotes")  
remotes::install_github("FinYang/flap")
```

# Outline

- 1 Hierarchical Data and Forecast Reconciliation
- 2 Probabilistic Forecasts
- 3 Quantile Forecasting
- 4 Non-Linear Forecasting
- 5 Beyond Hierarchies
- 6 Wrap-up

# Final thoughts

- Sometimes understanding the same problem in a different way opens new doors in research.



# Final thoughts

- Sometimes understanding the same problem in a different way opens new doors in research.
- Theory, methodology and application all matter. The connections and feedback loops between them are important.

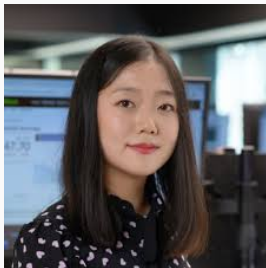
# Final thoughts

- Sometimes understanding the same problem in a different way opens new doors in research.
- Theory, methodology and application all matter. The connections and feedback loops between them are important.
  - ▶ Do not neglect any of these!

# Final thoughts

- Sometimes understanding the same problem in a different way opens new doors in research.
- Theory, methodology and application all matter. The connections and feedback loops between them are important.
  - ▶ Do not neglect any of these!
- Work with the right people

# The right people



# Forecast reconciliation

- Forecast reconciliation is a practical and interesting problem with many open questions.

# Forecast reconciliation

- Forecast reconciliation is a practical and interesting problem with many open questions.
  - ▶ How can we guarantee that reconciled probabilistic forecasts are correctly calibrated?

# Forecast reconciliation

- Forecast reconciliation is a practical and interesting problem with many open questions.
  - ▶ How can we guarantee that reconciled probabilistic forecasts are correctly calibrated?
  - ▶ Multi-objective optimization aspects of reconciliation problem.

# Forecast reconciliation

- Forecast reconciliation is a practical and interesting problem with many open questions.
  - ▶ How can we guarantee that reconciled probabilistic forecasts are correctly calibrated?
  - ▶ Multi-objective optimization aspects of reconciliation problem.
  - ▶ Other loss functions?



# Forecast reconciliation

- Forecast reconciliation is a practical and interesting problem with many open questions.
  - ▶ How can we guarantee that reconciled probabilistic forecasts are correctly calibrated?
  - ▶ Multi-objective optimization aspects of reconciliation problem.
  - ▶ Other loss functions?
- Jump on the bandwagon!



Postdoc opportunity



Link to slides