- Ronald Richman FIA FASSA CERA ronald.richman@qedact.com
- Associate Director QED Actuaries & Consultants
- 9 September 2020

Time-Series Forecasting of Mortality Rates using Deep Learning



Perla, Francesca and Richman, Ronald and Scognamiglio, Salvatore and Wüthrich, Mario V., Time-Series Forecasting of Mortality Rates using Deep Learning (May 6, 2020). Available at SSRN: <u>https://ssrn.com/abstract=3595426</u> or <u>http://dx.doi.org/10.2139/ssrn.3595426</u>





- **Rationale**
- Lee-Carter Paradigm for Mortality Modelling
- **Deep Learning for Mortality Modelling**
- Time Series Forecasting with Deep Learning



Rationale

- High quality mortality data available for multiple countries in the HMD
- May require forecasts of mortality for multiple countries simultaneously for modelling
 - Insurers operating in multiple jurisdictions need to set assumptions for many populations
- Reasonable to expect that improved forecasts for single countries could be made if all data were utilized to build a model
 - May be of interest even if mortality is only being forecast for a single country



Swiss Female raw log-mortality rates



Swiss Male raw log-mortality rates



- Rationale
- Lee-Carter Paradigm for Mortality Modelling
- **Deep Learning for Mortality Modelling**
- Time Series Forecasting with Deep Learning



Lee-Carter Model

- Forecast mortality rates = key inputs into demographic forecasting, life insurance and pensions models
- Foundational model for mortality forecasting is the Lee-Carter model (Lee and Carter 1992) (LC model) model) suited to old-age mortality (model coefficients of logistic model of qx)
- Mortality over time modeled using:
- i.e. (log) mortality = average rate + rate of change . time index
- Relies on latent variables that must be estimated from data and then multiplied to be fit compared to the t+x effects in the Lee-Carter model.
- and Carter 1992)

Many other approaches; within actuarial literature see Cairns, Blake and Dowd (2006) for an approach (CBD

$\log\left(u_{x,t}\right) = a_x + b_x k_t$

Could use interaction term between the variables Year and Age but this specification would require t.x effects

=> use non-linear/PCA regression to estimate the latent terms (Brouhns, Denuit and Vermunt 2002; Currie 2016; Lee





Producing Forecasts

- Time index k_{t} estimated for years within sample => need to extrapolate k_{t} for out-of-sample forecasts
- Time series models of varying complexity used to forecast k_t
- Two-step process fit model (a_x, b_x, k_t) and extrapolate common to other mortality models, such as CBD model
- Key judgement in LC model: over what period should the LC model be calibrated so that a, & b, appropriate for forecasting period?
- <u>Problem 1: If many forecasts are required, e.g. for multiple populations, then a manual process of selecting</u> calibration periods is required if using the LC model
- Single population extensions Cohort effect (Renshaw and Haberman 2006) Smoothing time series (Currie 2013)





Extending the LC Model

- What about multiple populations?
- Intuition = multi-population mortality forecasting model should produce more robust forecasts technology) Common trends likely captured with more statistical credibility

• => Li and Lee (2005) recommend even if interest is in single series

Augmented Common Factor (Li and Lee 20

Common Age Effect (Kleinow 2015)

<u>Problem 2:</u> Not intended for large scale mortality forecasting - generally applied on smaller sub-set of data => judgment of modeler needed Hard to fit (complex optimization schemes/less known statistical techniques)

• Which specification is better, when, and why?

Common factors (similar socioeconomic circumstances, shared improvements in public health and medical

005)	$\log\left(u_{x,t}\right) = a_x^i + b_x k_t + b_x^i k_t^i$
	$\log\left(u_{x,t}\right) = a_x^i + b_x k_t^i$



Formalizing the LC Paradigm

• Mortality modelling:

- (t, x, i)
- LC paradigm for mortality forecasting for population

 $\log\left(u_{x,t}^{(i)}\right) =$

- $a_x^{(i)} \in \mathbb{R}, \ \boldsymbol{b}_x^{(i)}, \boldsymbol{k}_t^{(i)}$
- Different models produced by choice of q
- q = 1 and parameters estimated for each population => LC model
- q > 1 => multi factor LC model
- estimate $b_x^{(i)} = b_x$ i.e. not depending on population I => CAE model
- q=2 and estimate population-specific and overall time effects => ACF model

$$\mapsto \log \left(u_{x,t}^{(i)} \right)$$

$$i:$$

$$= a_x^{(i)} + \left\langle b_x^{(i)}, k_t^{(i)} \right\rangle$$

$$b_x^{(i)}, k_t^{(i)} \in \mathbb{R}^q$$



LC and Embeddings

• LC model utilizes step functions:

 $\log(u_{x,t}) =$

 $g(x) = \begin{cases} a_1 & \text{for } x = 1, \\ a_2 & \text{for } x = 2, \\ \vdots & \\ a_\omega & \text{for } x = \omega, \end{cases}$

• For a set of categories P with n_p levels, an embedding of dimension q_P is:

 $z_{\mathcal{P}}: \mathcal{P} \to \mathbb{R}^{q_{\mathcal{P}}}, \qquad p \mapsto z_{\mathcal{P}}(p).$

.C model can be interpreted using embeddings (see Richman and Wüthrich (2019)): • L

$$(x,i) \mapsto a_x^{(i)} \in \mathbb{R}.$$

$$= g(x) + h(x)i(t),$$





- Rationale
- Lee-Carter Paradigm for Mortality Modelling
- **Deep Learning for Mortality Modelling**
- Time Series Forecasting with Deep Learning



Deep Learning

to represent abstract concepts

Features in lower layers composed of simpler features constructed at higher layers => complex concepts can be represented automatically

- models, where each layer learns a new representation of the features.
- The principle: Provide raw data to the network and let it figure out what and how to learn.
- level abstractions that would be useful to represent the kind of complex functions needed for AI tasks."



Deep Learning = representation learning technique that automatically constructs hierarchies of complex features

Typical example of deep learning is feed-forward neural networks, which are multi-layered machine learning

Desiderata for AI by Bengio (2009): "Ability to learn with little human input the low-level, intermediate, and high-







Recurrent NN – Temporal data

- Data with temporal structure implies that previous observations should influence the current observation
- **Recurrent network maintains state of hidden** neurons over time

Past representation useful for current prediction i.e. network has a 'memory'

- Key challenge difficult to train due to vanishing gradients
- Several implementations of the recurrent concept which control how network remembers and forgets state

Long Short Term Memory (LSTM) Gated Recurrent Unit (GRU)

RNNs can make predictions at last time step or for each input



x = Input vector S = hidden state (layers) 0 = output Arrows indicate the direction in which data flows.

Folded

Unfolded





Convolutional NN - Images

- **Prior features in images are position invariant i.e. can** recognize at any position within an image Also applies to audio/speech and text/time series data
- **Convolutional network is locally connected and shares** weights => expresses prior of position invariance Far fewer parameters than FCN
- Each neuron (i.e. feature map) in network derived by applying filter to input data

Weights of filter learned when fitting network Multiple filters can be applied

Can also be used for time-series applications



			Dat	a M	latri	X				
0	7	70	9	0	0	0	0	0	0	
0	0	0	3	4	/4/	4	1	0	0	
-θ	:Q	, t,	لي الم	0	0		/74	4	(4)	
0	0	3	0	0	3	4	Ύ.	0	0	Filter
0	1	0	0	1	4	2	1	Ø	0	1 1 1
0	0	0	1	4	2	1	0	0	0	000
0	0	1	4	1	1	0	0	0	0	-1 -1 -1
0	0	4	4	4	4	4	4	0	0	
0	0	0	0	0	0	0	0	0	0	
0	0	0	0	0	0	0	0	0	0	0*1 0*1 0*1 0
										0*0 0*0 0*0 = 0 0
-1	-4	-4	-3	-1	-5	-5	-4			0*-1 0*-1 1*-1 0 0
-3	0	4	-8	5	1	0	0			
0	3	3	-2	-6	-2	2	3			
3	2	-2	-4	0	5	4	1			
0	-4	-5	-1	5	6	3	1			
-4	-7	-7	-5	-5	-9	-7	-4			
1	5	6	6	2	1	0	0			
4	8	12	12	12	12	8	4			
		Fea	atur	e N	lap					





Extending LC – two perspectives

- Lee Carter model = regression model using features derived from data using PCA CAE + ACF = LC-type regression models with features derived at a regional level
- <u>Perspective 1:</u> Use a neural network to model the regression problem and let it decide on the feature set LC is a linear model once parameters are known => use a NN to derive non-linear model to derive more predictive specification
- Perspective 2: use a more general step function formulation to specify the multi-population model LC uses step functions of single dimension for each parameter:

 $\log\left(u_{x,t}\right) = g(x) + h(x)i(t),$



Use NN to derive multi-dimensional vectors for each parameter using an embedding layer

CAE/ACF specify the types of interaction between population-level and regional level parameters => use a NN





Lee-Carter Neural Network

Multi-population mortality forecasting model (Richman and Wüthrich 2018) Input Layer Supervised regression on HMD data (inputs = Year, Country, Age; outputs = m_x) **5 layer deep FCN** Generalizes the LC model in both ways mentioned before Note that no time-series forecasting is done Year enters the model as a numerical variable. Forecasts made by predicting using values for Year beyond the range of input data **Output Layer**







Multi-population results

- Results of comparing the models (LC/ACF/CAE/LCNN)
- Best performing model is <u>deep neural</u>
 <u>network...</u>
- ...produces the best out-of-time forecasts 51 out of 76 times
- For purposes of large scale mortality forecasting, deep neural networks <u>dramatically</u> <u>outperform traditional single and multi-</u> <u>population forecasting models</u>

	Model	Average MSE	Median MSE	Best Perform
1	LC_SVD	5.50	2.48	
2	ACF_SVD_region	3.46	2.50	
3	ACF_SVD_country	7.30	4.77	
1	ACF_BP	6.12	3.00	

	Model	Average MSE	Median MSE	Best Perform
1	LC_SVD	5.50	2.48	
2	CAE_SVD	4.76	2.35	
3	CAE2_SVD	12.01	1.79	
4	CAE2_BP	5.59	3.46	

	Model	Average MSE	Median MSE	Best Perform
L	LC_SVD	5.50	2.48	
2	LC_ACF_region	3.46	2.50	
3	ACF_BP	6.12	3.00	
1	CAE_BP	5.59	3.46	
5	DEEP	2.68	1.38	









Features in last layer of network

- Representation = output of last layer (128 dimensions) with dimension reduced using PCA
- Can be interpreted as relativities of mortality rates estimated for each period
- Output shifted and scaled to produce final results
- Generalization of Brass Logit Transform where base table specified using NN (Brass 1964)

Country • GBRTENW • ITA • USA



 $y_x = logit of mortality at age x$ a,b = regression coefficients $z_x^{Ref} = logit of reference mortality$







Learned embeddings

- Age embeddings extracted from LCNN model
- Five dimensions reduced using PCA
- Age relativities of mortality rates
- In deeper layers of network, combined with other inputs to produce representations specific to:
 - Country Gender Time
- **First dimension of PCA is shape of lifetable**
- Second dimension is shape of child, young and older adult mortality relative to middle age and oldest age mortality



Application to Insurance Data

- Applied the same network structure to data from a reinsurer in Rossouw and Richman (2019)
 - Consists of mortality and morbidity rates from 4 contributing companies over ~15 years Trained on 5 years of data and forecast 4 years
- Compared results of NN to several other models (LASSO/GBM) using Poisson deviance as criterion
- NN beat other models but had bias at portfolio level (see Wüthrich (2019))
- Debiased results: Poisson deviance: 22 836 AvE: 99.7%

Model	Poisson deviance	Actual vs. predicted
glm_trad_ibnr	22 944	99.3%
glm_trad_ebner	22 947	99.5%
glm	22 883	97.0%
glmnet	22 826	98.3%
xgb	22 822	95.9%
dl	22 799	93.8%







LC go Machine Learning: RNNs

- Mortality rates have a time-series structure, e.g. Swiss Female mortality in 1990-2001
- Can NNs <u>exploit time-series structure</u> to forecast mortality directly?
 RNNs seem to be a natural choice due to sequential processing
 LCNN does not rely on time-series structure directly
- Swiss mortality rates forecast using RNNs in Richman and Wüthrich (2019)
 Models trained on data 1950-1999
 Forecasts made for 2000-2016
 Models fit for each gender separately
- To reduce volatility of input data, matrices of rates at ages x-x+4 fed into networks to forecast rates at age x+2



	in-sample		out-of-sample	
	female	male	female	mal
LSTM3 $(T = 10, (\tau_0, \tau_1, \tau_2, \tau_3) = (5, 20, 15, 10))$	2.5222	6.9458	0.3566	1.350
GRU3 $(T = 10, (\tau_0, \tau_1, \tau_2, \tau_3) = (5, 20, 15, 10))$	2.8370	7.0907	0.4788	1.243
LC model with SVD	3.7573	8.8110	0.6045	1.815





Extending the RNN model

NN models flexible enough to incorporate extensions easily, e.g., joint modelling of **both genders**

Gender incorporated explicitly using dummy-coding and implicitly through rates input to the network

- **Results improved versus single gender** models LSTM model beats the GRU model.
- **Direct forecasting using RNNs leads to** unstable results, which are much improved via model averaging (ensembling)
- **Ensemble model captures improvements,** particularly in young adult mortality significantly better than LC.







	in-sample	out-of-	sample
	both genders	female	\mathbf{male}
A3 $(T = 10, (\tau_0, \tau_1, \tau_2, \tau_3) = (5, 20, 15, 10))$	4.7643	0.3402	1.1346
3 $(T = 10, (\tau_0, \tau_1, \tau_2, \tau_3) = (5, 20, 15, 10))$	4.6311	0.4646	1.2571
odel with SVD	6.2841	0.6045	1.8152



Male out-of-sample losses over 100 SGDs



Male observed 2000–2016



Male LC predictions 2000–2016



Male LSTM predictions 2000–2016







Combining LC + RNNs

- Deep Learning Integrated Lee–Carter Model of Nigri, Levantesi, Marino, Scognamiglio & Perla, F. (2019)
- Apply SVD to derive a_x , b_x , k_t of LC model
- Instead of ARIMA, forecast k_t using RNN

Country	Male		Fer	nale
Australia	MAE	RMSE	MAE	RMSE
κ_t ARIMA	24.75	28.04	13.95	15.55
κ_t LSTM	2.57	3.24	3.12	3.83
Denmark	MAE	RMSE	MAE	RMSE
κ_t ARIMA	11.10	16.10	7.99	10.70
κ_t LSTM	2.97	3.84	6.62	8.18
Italy	MAE	RMSE	MAE	RMSE
κ_t ARIMA	55.41	63.46	41.12	45.74
κ_t LSTM	4.63	5.48	8.04	10.69
Spain	MAE	RMSE	MAE	RMSE
κ_t ARIMA	20.24	26.33	26.10	33.95
κ_t LSTM	7.69	8.96	15.61	17.67
the USA	MAE	RMSE	MAE	RMSE
κ_t ARIMA	8.39	9.48	10.81	12.54
κ_t LSTM	2.31	2.86	3.32	4.18
Japan	MAE	RMSE	MAE	RMSE
κ_t ARIMA	9.61	10.50	15.12	17.52
κ_t LSTM	4.71	5.24	7.89	10.23

Australia Female - k(t) 1921 to 1995 - Forecasting: 1996 to 2014 ARIMA (blue) Vs. LSTM (red)











Denmark Female - k(t) 1835 to 1980 - Forecasting: 1981 to 2016 ARIMA (blue) Vs. LSTM (red)



Spain Female - k(t) 1908 to 1994 - Forecasting: 1995 to 2016 ARIMA (blue) Vs. LSTM (red)



Japan Female - k(t) 1947 to 2002 - Forecasting: 2003 to 2016 ARIMA (blue) Vs. LSTM (red)













- Rationale
- Lee-Carter Paradigm for Mortality Modelling
- **Deep Learning for Mortality Modelling**
- **Time Series Forecasting with Deep Learning**



Processing Time Series with DL

- RNN models still specialized for single population can this be expanded to the multi-population case?
- Furthermore, can the excessive volatility of the RNN calibrations be reduced?
- In the LCNN model, we have the following regression function learned by the network:

- Hypothesis: can neural network models designed for directly processing sequential data outperform more general network architectures applied to sequential data?
- i.e. we wish to map directly from observed mortality rates of many populations to a time feature:

 $U^{(i)} \vdash$

Addressed in Perla, Richman, Scognamiglio and Wüthrich (2020)

 $(t, x, i) \mapsto \log\left(u_{x,t}^{(i)}\right)$

$$ightarrow ~ oldsymbol{k}_t^{(i)} ~\in~ \mathbb{R}^q$$
 .





Defining the Model

• For modelling, work with scaled version of mortality rates:

$$y_{t_0+T+1}^{(i)} = \frac{\log(u_{t_0+T+1}^{(i)}) - y_0}{y_1 - y_0} \in [0, 1]^d,$$

where y_1 is the maximum of the observed log rates, and y_0 is the minimum.

For region r and gender g, can extend LC paradigm to:

$$\sigma^{-1}\left(\widehat{y}_{x,t_0+T+1}^{(r,g)}\right) = w_{x,0} + \left\langle W_x^{\mathcal{R}}, z_{\mathcal{R}}(r) \right\rangle + \left\langle W_x^{\mathcal{G}}, z_{\mathcal{G}}(g) \right\rangle + \left\langle W_x^f, z_f(U_{t_0}^{(i)}) \right\rangle$$

$$a_x \qquad b_x \cdot k_t$$

where:

- $z_R(r)$ is a region embedding and W_x^R is a region coefficient
- $z_G(g)$ is a gender embedding and W_x^G is a gender coefficient
- a trend age adjustment coefficient

 $z_f(U_{t0}^i)$ is a representation learned using a neural network with input = U_{t0}^i (matrix of past mortality rates) and W_x^f is





Model Structure

- Similar to LCNN model...
- ... however, time variable replaced with outputs of a NN processing layer
- **Diagram shows a CNN being used to** derive the mortality trend
- Can also use an RNN
- **CNNs** appear to perform better than RNNs
- LCCONV model forecasts mortality rates in a single step...
- ... no time series forecasting required.
- Forecasts derived by adding single year forecasts to matrix of rates, and reapplying the model.





Forecast Results – HMD (1)

- The CNN model (LCCONV) achieves better performance versus the LC model on **75/76 populations in the HMD**
- LCCONV beats the LCNN model in an extra 8 populations and achieves a substantially lower out-of-sample MSE
- **Residual plot shows that model is** substantially better for males, whereas the performance is similar for females
- Using RNNs to process the data, instead of predicting also leads to good performance

moo LCO LC LC LC LC LC LC LC LC DE



del	$test_loss$	ensemble MSE	# populations
CONV	2.27	2.24	75/76
$\operatorname{CONV_tanh}$	2.62	2.58	61/76
CONV_relu	3.26	3.10	57/76
LSTM1	2.86	2.54	69/76
$LSTM1_tanh$	3.32	3.03	58/76
LSTM1_relu	3.33	3.25	52/76
LSTM2	2.43	2.32	74/76
$LSTM2_tanh$	2.36	2.27	75/76
LSTM2_relu	3.44	3.11	56/76
EP	2.83	2.53	67/76



Forecast Results – HMD (2)

- **Populations sorted by size LC model** error increases with size of population
- Similar results for NN models, however, lower than LC model in most instances
- LCCONV lower than LCNN mainly on large male populations







Forecast Results – USMD

- **Does the LCNN model generalize beyond** the HMD?
- Fit model to the USMD (50 states + DC)
- Beats LC model 101/102 times
- **RNN models more competitive on USMD** than on HMD
- Cluster 1 includes states located in the south of US (Arizona, Florida, New Mexico, Texas) / US states geographically distant from the US zone such as Alaska and Hawaii.
- **Cluster 2 groups countries in the southeast zone** ۲ (Louisiana, Georgia, South Carolina, Mississippi and Alabama)
- Cluster 3 central (Colorado, Kansas, Utah and ۲ Missouri) and the northwest (Washington, Oregon, **Montana and**
- North Dakota)
- **Cluster 4-** The last cluster groups countries in the northeast of US (Pennsylvania, Connecticut, Massachusetts, Virginia and West Virginia).

model	test loss	ensemble MSE	# populations
LCCONV	0.50	0.49	101/102
$LCCONV_{tanh}$	0.55	0.54	100/102
LCCONV_relu	0.74	0.73	86/102
LCLSTM1	0.63	0.58	98/102
$LCLSTM1_{tanh}$	0.66	0.53	99/102
LCLSTM1_relu	1.20	1.10	52/102
LCLSTM2	0.50	0.46	102/102
$LCLSTM2_{tanh}$	0.53	0.50	101/102
LCLSTM2_relu	1.09	0.96	62/102





Conclusion

- Deep learning models provide new opportunities for modelling mortality
- Potential gains of moving from traditional specification of models for mortality
- Some models can directly process mortality rates to produce forecasts with increased accuracy over more general models...
- however finding an optimal model architecture can be challenging
- **Future research should address:**
 - reasons for high variability of RNN models applied to forecast mortality rates uncertainty bounds on predictions





References

See https://gist.github.com/RonRichman/655cca0dd79afcd20b33d3131c537414 •

