

# Cuantificación de la incertidumbre en la predicción espacio-temporal con bayesian deep learning: Aplicación a El Niño

Autor: Jose González Abad  
Director: Jorge Baño Medina



# ÍNDICE

- Introducción
  - El Niño
  - Cuantificación de la Incertidumbre
  - Predicción de El Niño
- Deep Learning aplicado al Forecasting
  - Red Neuronal Densa
  - Red Neuronal Recurrente
  - Red Neuronal Convolutacional
  - Deep Learning
  - Bayesian Deep Learning
- Datos y Metodología
- Métodos
  - Modelos de referencia
  - Modelos desarrollados
- Resultados
- Conclusiones

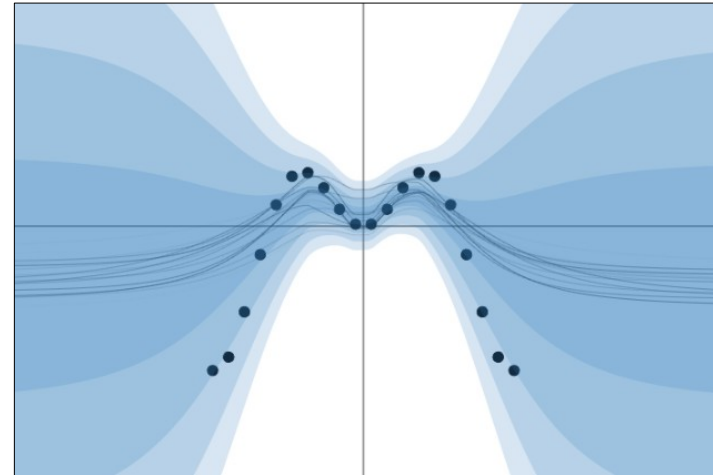
# Introducción

# INTRODUCCIÓN

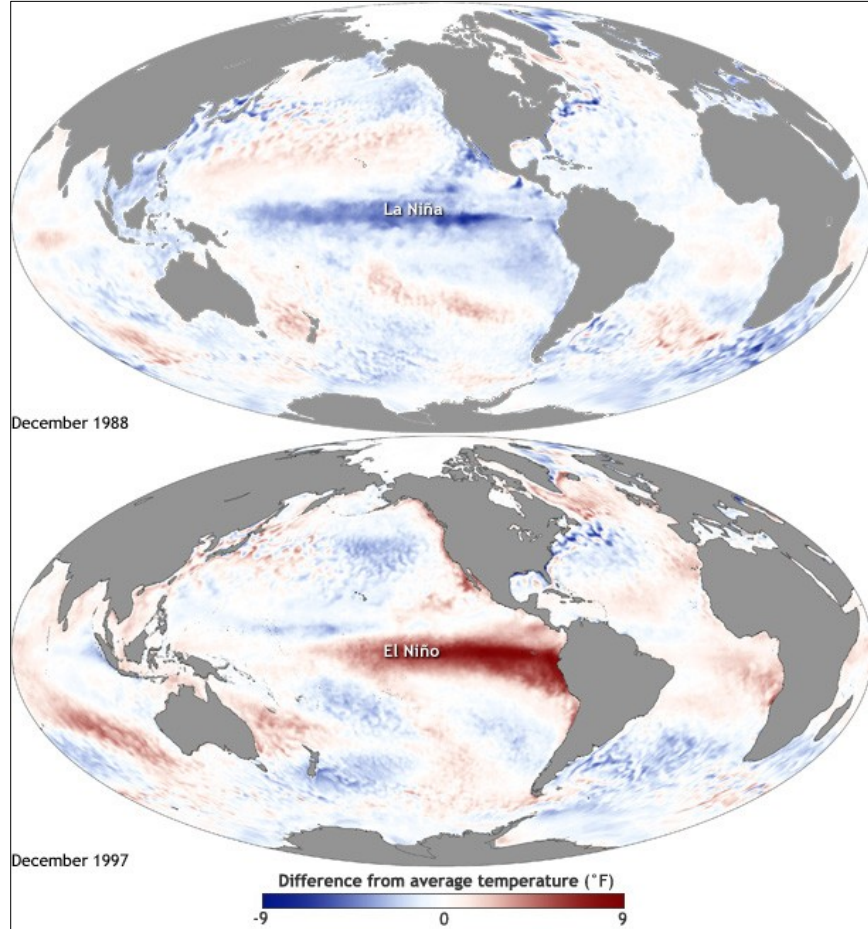
## El Niño



## Cuantificación de la incertidumbre

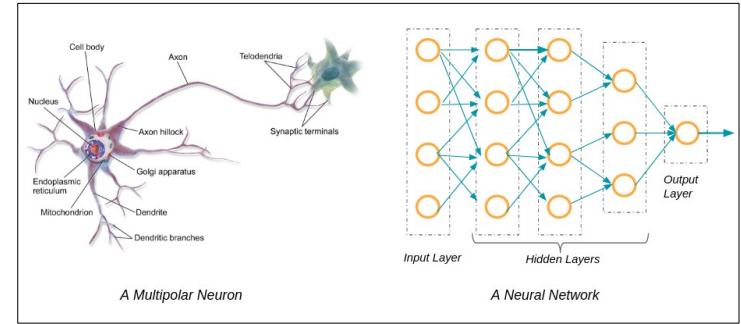
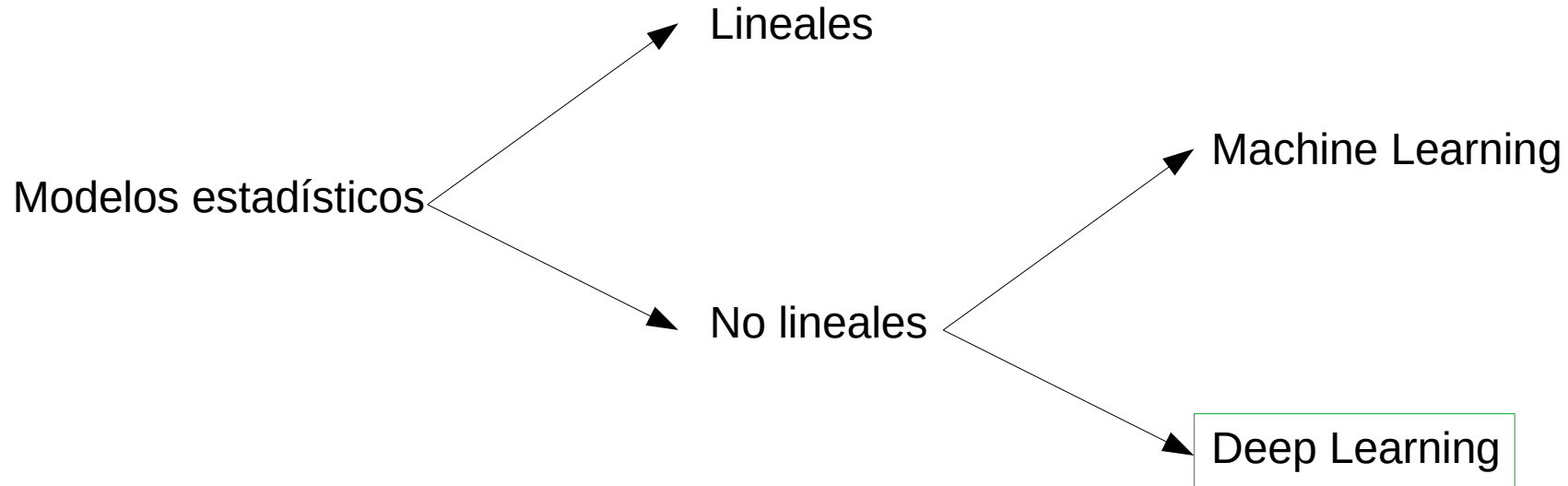


# El Niño



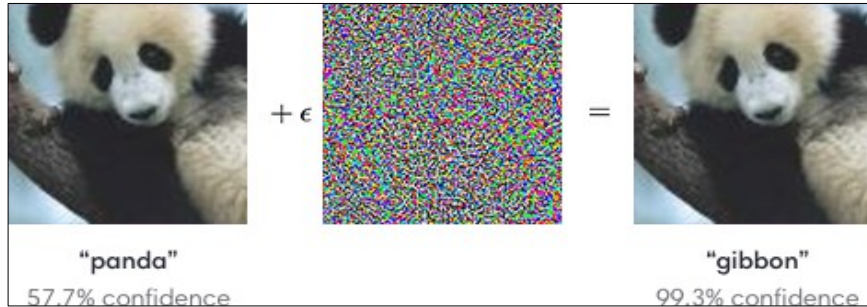
# El Niño

Modelos físico-matemáticos

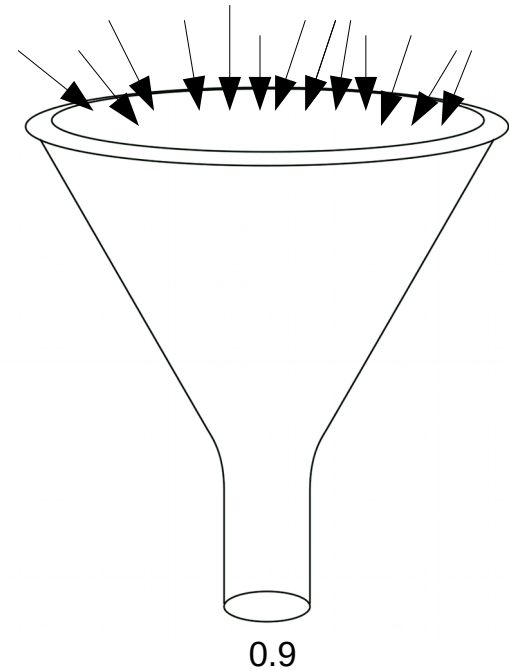


# Cuantificación de la incertidumbre

¿Cómo de seguros están  
nuestros modelos?

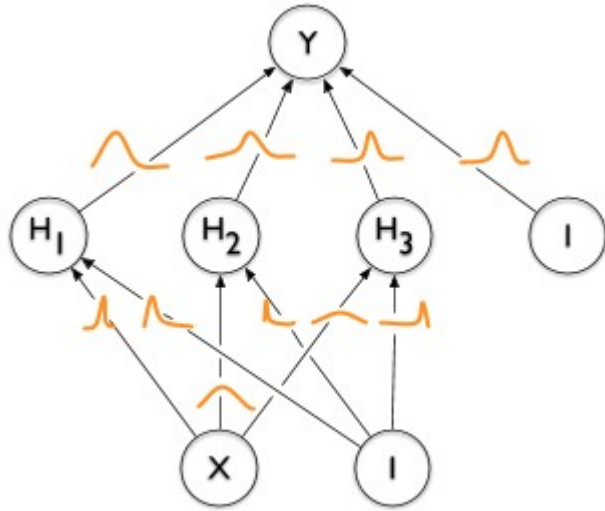


Intervalo de confianza

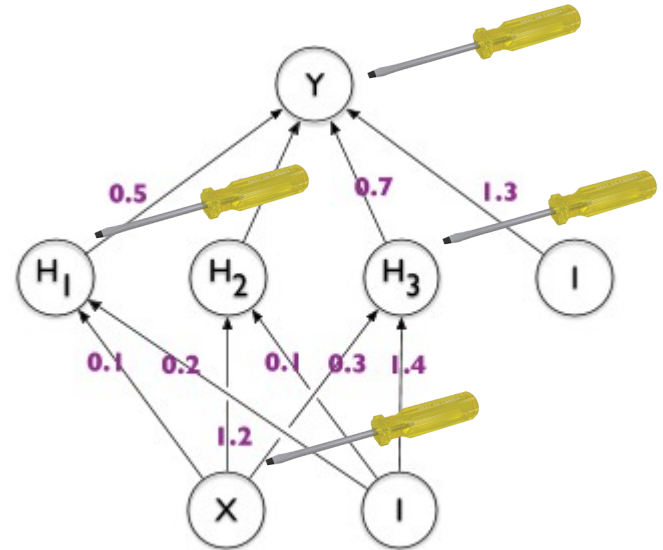


# Cuantificación de la incertidumbre

Redes neuronales bayesianas



Bayesian deep learning





# Predicción de El Niño

## Bayesian Recurrent Neural Network Models for Forecasting and Quantifying Uncertainty in Spatial-Temporal Data

Patrick L. McDermott\* Christopher K. Wikle

Department of Statistics

University of Missouri

February 8, 2018

### Abstract

Recurrent neural networks (RNNs) are nonlinear dynamical models commonly used in the machine learning and dynamical systems literature to represent complex dynamical or sequential relationships between variables. More recently, as deep learning models have become more common, RNNs have been used to forecast increasingly complicated systems. Dynamical spatio-temporal processes represent a class of complex systems that can potentially benefit from these types of models. Although the RNN literature is expansive and highly developed, uncertainty quantification is often ignored. Even when considered, the uncertainty is generally quantified without the use of a rigorous framework, such as a fully Bayesian setting. Here we attempt to quantify uncertainty in a more formal framework while maintaining the forecast accuracy that

\*Department of Statistics, University of Missouri, 146 Middlebush Hall, Columbia, MO 65211 USA;  
E-mail: plmyt7@mail.missouri.edu (corresponding author)

VS



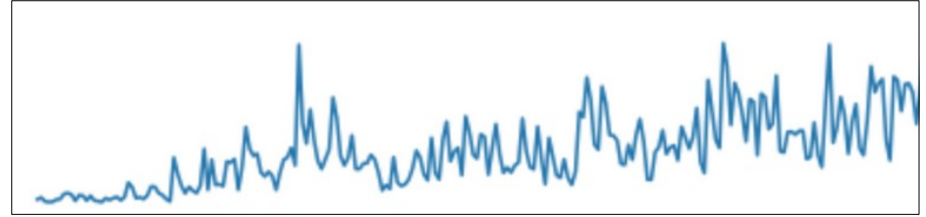
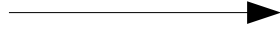
+

Bayesian deep  
learning

# Deep Learning aplicado al Forecasting

# Deep Learning aplicado al Forecasting

El Niño presenta una marcada componente temporal



Enfoque determinista  
Metodología Box-Jenkins

} **Univariante**

VAR } **Multivariante**

} **Lineal**



Con intervalos de confianza

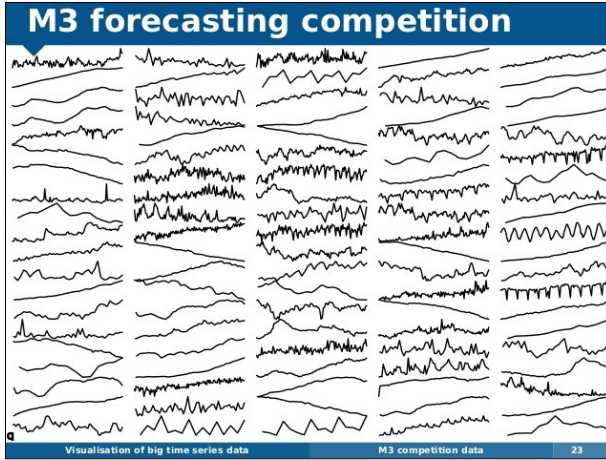
Machine Learning  
Redes Neuronales

} **No lineal**

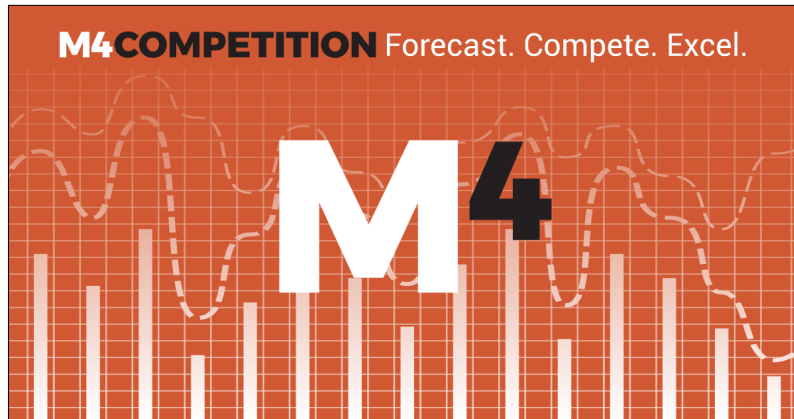


Sin intervalos de confianza

# Deep Learning aplicado al Forecasting

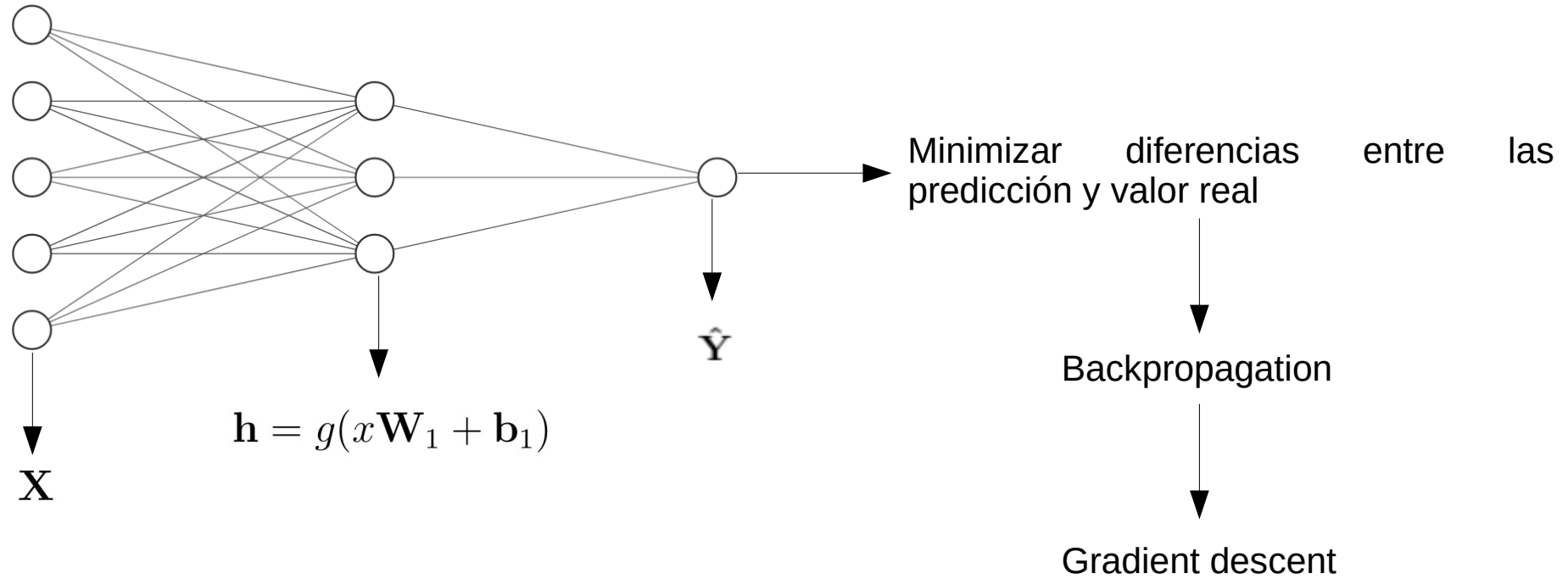


Método Theta



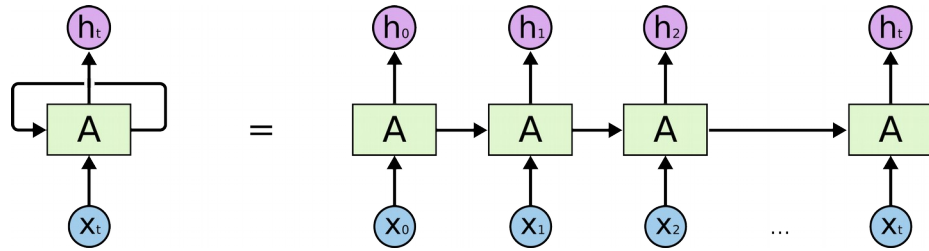
Red neuronal recurrente +  
suavizados exponenciales

# Red neuronal densa



# Red neuronal recurrente

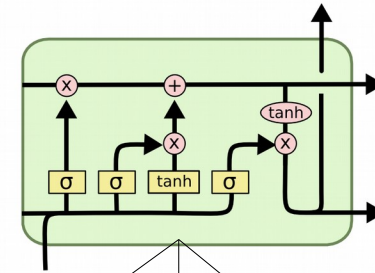
Red neuronal recurrente



$$h_t = g(\mathbf{W}\mathbf{x}_t + \mathbf{U}\mathbf{h}_{t-1})$$

Vanishing Gradient

LSTM

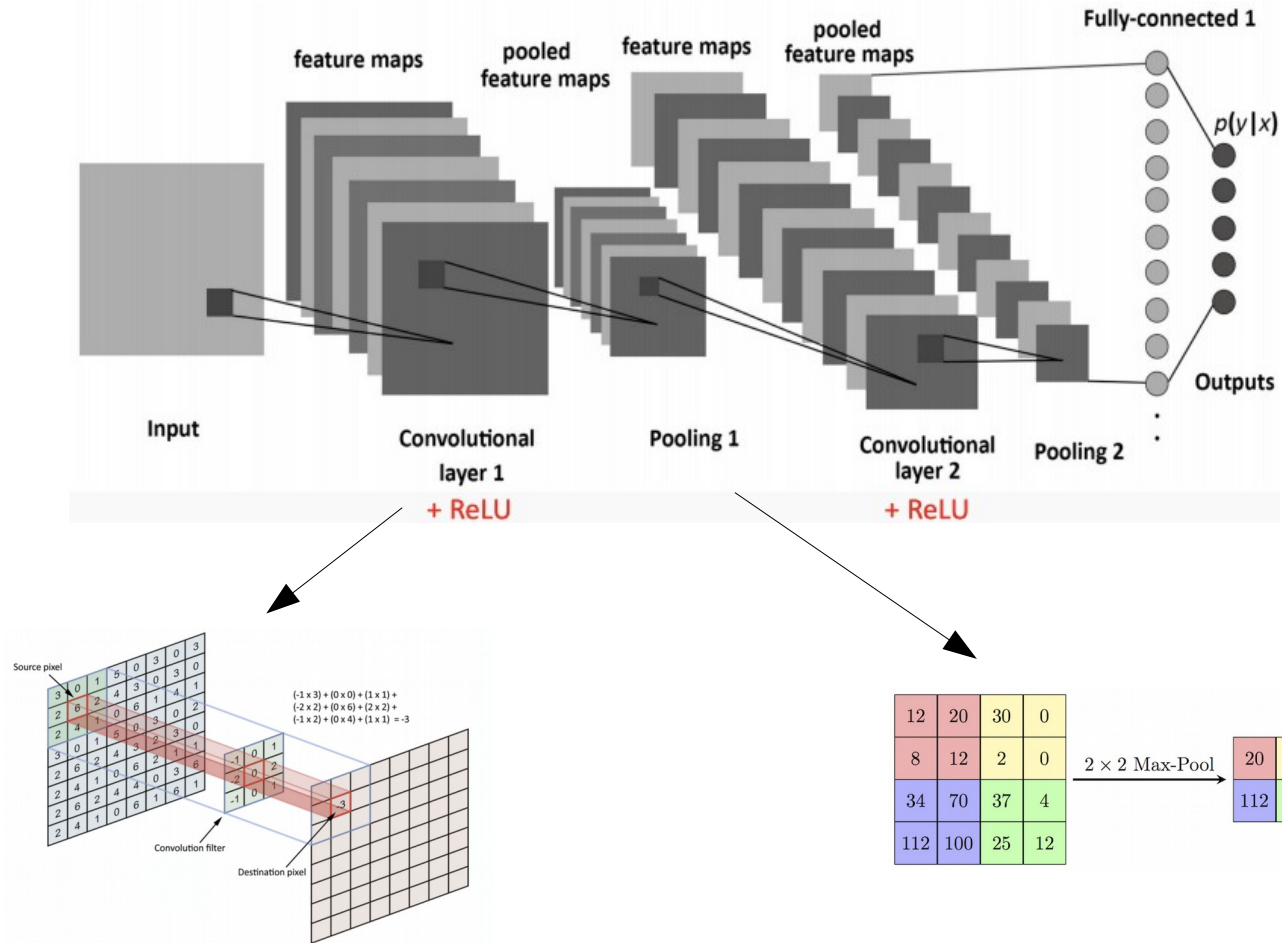


$$f_t = \sigma(\mathbf{W}_f[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_f)$$

$$\begin{aligned} \mathbf{o}_t &:= \sigma(\mathbf{W}_o[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_o) \\ \mathbf{h}_t &:= \mathbf{o}_t \tanh(\mathbf{c}_t) \end{aligned}$$

$$\begin{aligned} \mathbf{i}_t &:= \sigma(\mathbf{W}_i[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_i) \\ \mathbf{c}'_t &:= \tanh(\mathbf{W}_C[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_C) \end{aligned}$$

# Red neuronal convolucional

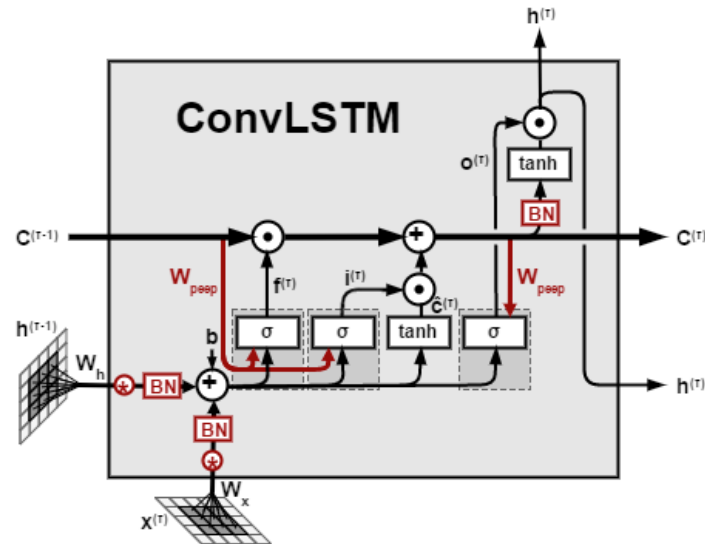


# Red neuronal convolucional

Dinámica temporal —→ LSTM

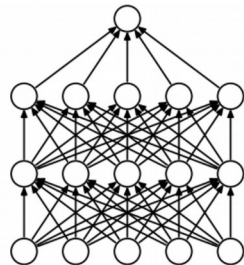
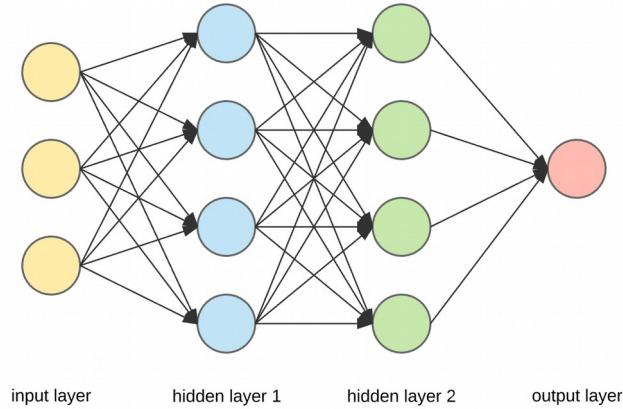
Dinámica espacial —→ Convolucional

**LSTM Convolucional**

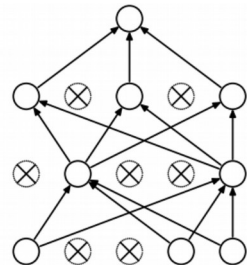




# Deep learning



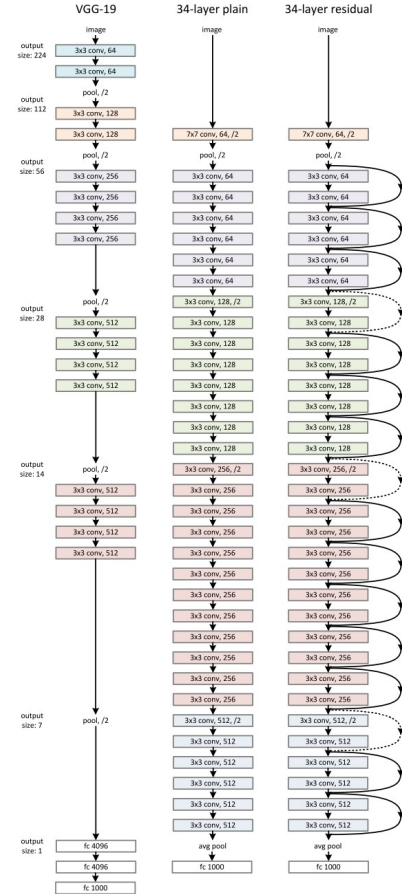
(a) Standard Neural Net



(b) After applying dropout.

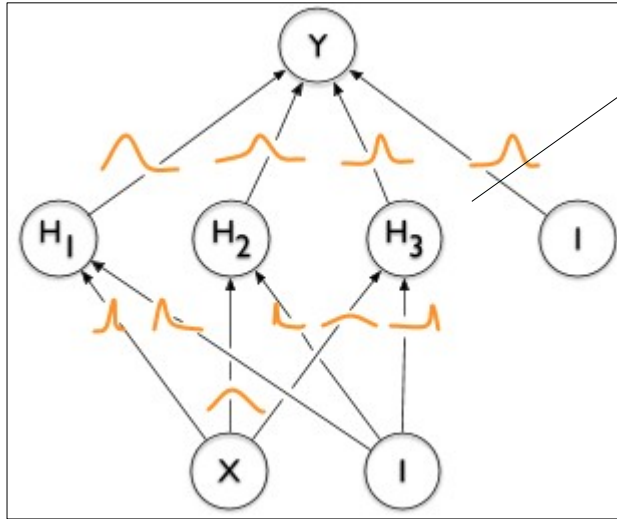
Dropout

- Incremento exponencial en la cantidad de datos
- Desarrollo de hardware específico
- Desarrollo de nuevos algoritmos



# Bayesian deep learning

Red neuronal bayesiana



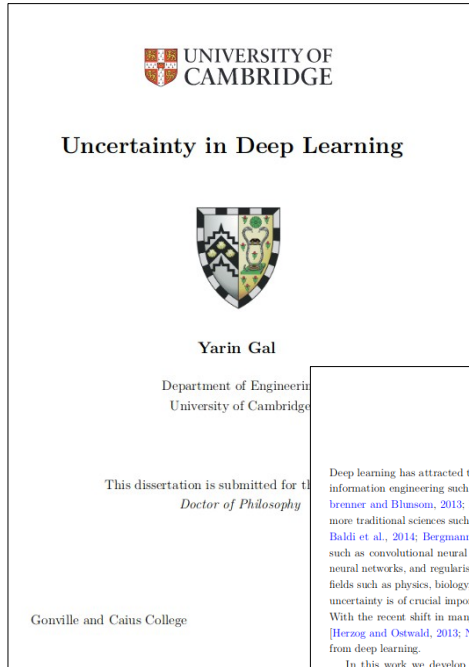
$$P(W|X, Y) = \frac{P(X|Y, W)P(W)}{\int P(X|Y, W)P(W)dW}$$

MCMC

Variational Inference

- Alto coste computacional
- Limitaciones en el diseño de arquitecturas

# Bayesian deep learning

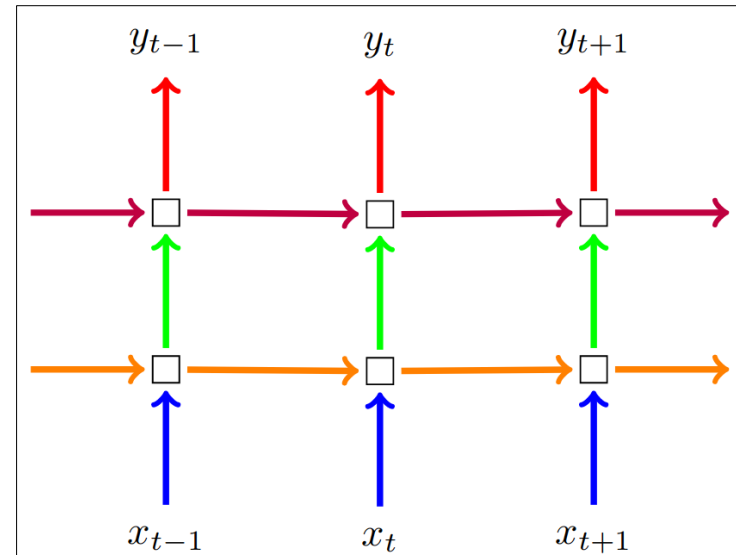


Si se define la función de distribución de aproximación  $Q$  del variational inference de una red neuronal bayesiana a partir de una Bernoulli su optimización es equivalente a la realizada por una red neuronal regularizada con dropout

## Abstract

Deep learning has attracted tremendous attention from researchers in various fields of information engineering such as AI, computer vision, and language processing [Kalchbrenner and Blusson, 2013; Krizhevsky et al., 2012; Mnih et al., 2013], but also from more traditional sciences such as physics, biology, and manufacturing [Anjos et al., 2015; Baldi et al., 2014; Bergmann et al., 2014]. Neural networks, image processing tools such as convolutional neural networks, sequence processing models such as recurrent neural networks, and regularisation tools such as dropout, are used extensively. However, fields such as physics, biology, and manufacturing are ones in which representing model uncertainty is of crucial importance [Ghahramani, 2015; Krzywinski and Altman, 2013]. With the recent shift in many of these fields towards the use of Bayesian uncertainty [Herzog and Ostwald, 2013; Nuzzo, 2014; Trafimow and Marks, 2015], new needs arise from deep learning.

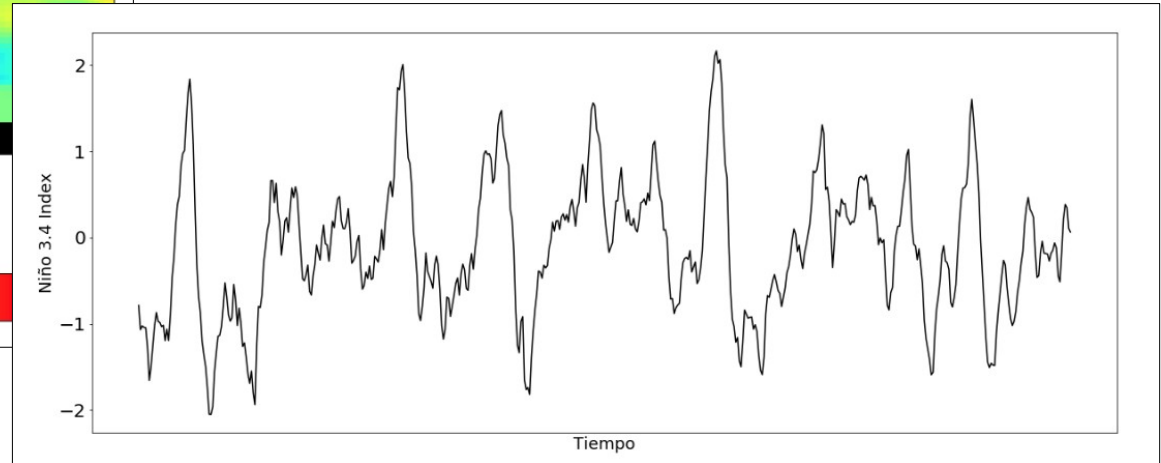
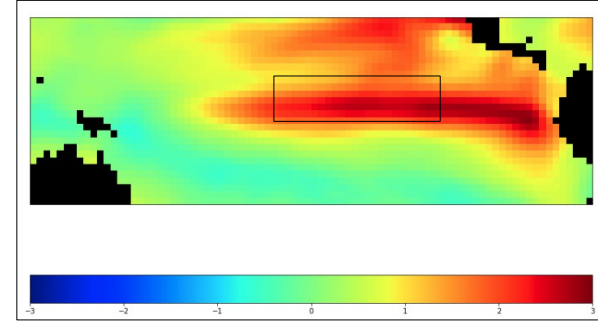
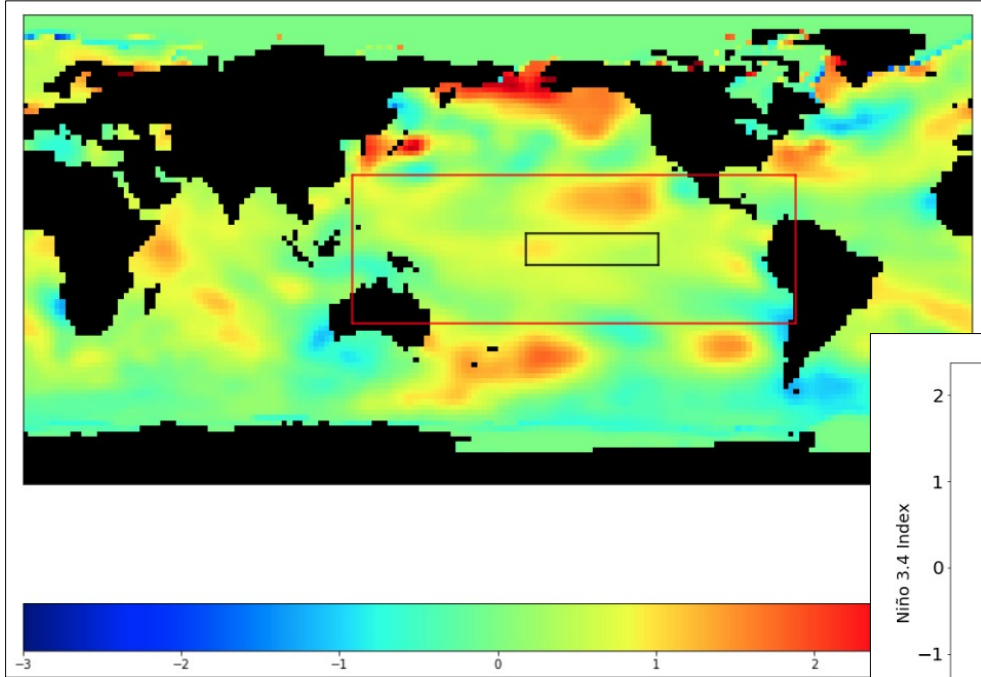
In this work we develop tools to obtain practical uncertainty estimates in deep learning, casting recent deep learning tools as Bayesian models without changing either the models or the optimisation. In the first part of this thesis we develop the theory for such tools, providing applications and illustrative examples. We tie approximate inference in Bayesian models to dropout and other stochastic regularisation techniques, and assess the approximations empirically. We give example applications arising from this connection between modern deep learning and Bayesian modelling such as active learning of image data and data-efficient deep reinforcement learning. We further demonstrate the tools' practicality through a survey of recent applications making use of the suggested techniques in language applications, medical diagnostics, bioinformatics, image processing, and autonomous driving. In the second part of the thesis we explore the insights stemming from the link between Bayesian modelling and deep learning, and its theoretical implications. We discuss what determines model uncertainty properties, analyse the approximate inference analytically in the linear case, and theoretically examine various priors such as spike and slab priors.



# Datos y Metodología

# Datos y Metodología

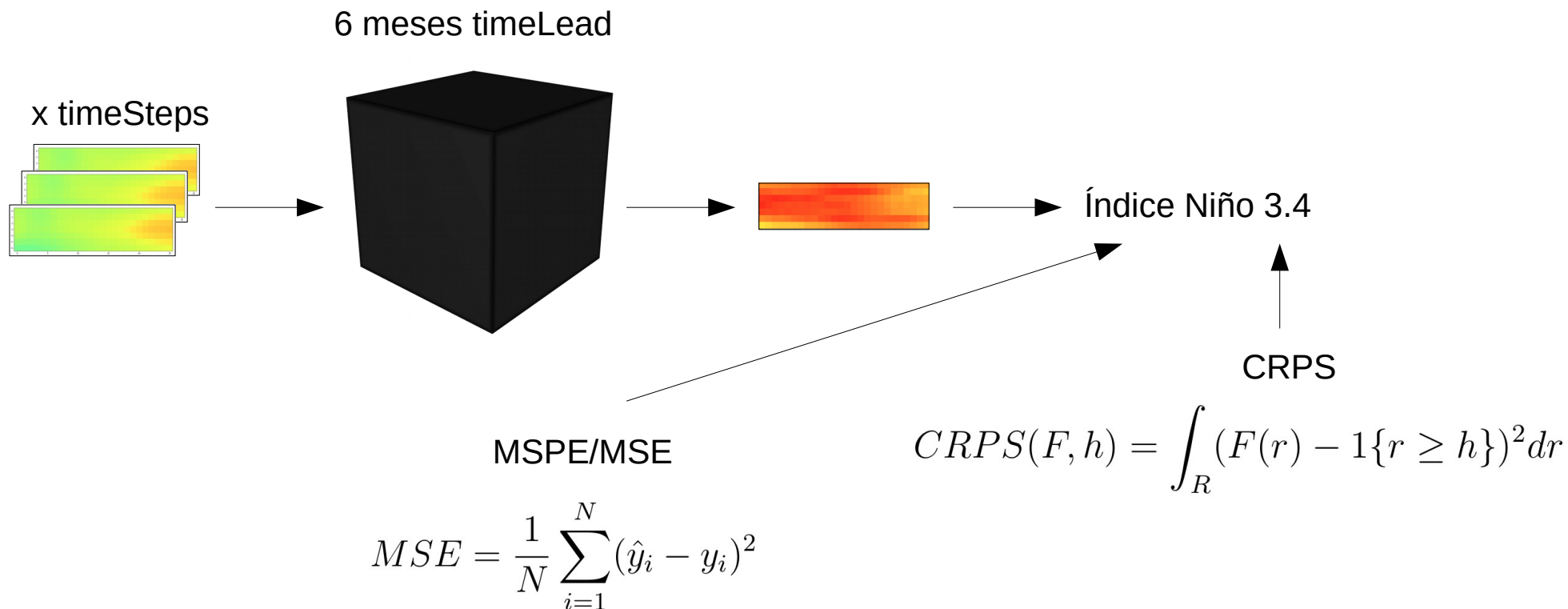
Fuente: NOAA



# Datos y Metodología

Conjunto de entrenamiento: Enero 1970 – Agosto 2014

Conjunto de test: Febrero 2015 – Diciembre 2016 (forma alterna) —► Anomalía ENSO Diciembre 2015



# Métodos

# Modelos de referencia

## Linear DTSM

$$Y_{t,i} = \sum_{j=1}^{n_y} a_{i,j} Y_{t-1,j} + \zeta_{t,i}^{(l)}$$

## GQN

$$Y_{t,i} = \sum_{j=1}^{n_y} a_{i,j} Y_{t-1,j} + \sum_{k=1}^{n_y} \sum_{l=1}^{n_y} b_{i,k,l} Y_{t-1,k} Y_{t-1,l} + \zeta_{t,i}^{(q)}$$

## E-QESN

Reservoir computing

## BAST-RNN

$$Y_t = \boldsymbol{\mu} + \mathbf{V}_1 \mathbf{h}_t + \mathbf{V}_1^2 \mathbf{h}_t^2 + \epsilon_t, \quad \epsilon_t \sim \text{Gau}(0, \mathbf{R}_t)$$

$$\mathbf{h}_t = f\left(\frac{\delta}{|\lambda_w|} \mathbf{W} \mathbf{h}_{t-1} + \mathbf{U} \mathbf{x}_t\right)$$



MCMC



# Modelos desarrollados

Dropout  $p = 0.05, 0.1, 0.2, 0.5, 0.8$

Algoritmo ADAM

Early-stopping sobre un 10% de datos de validación

## LIN

Modelo lineal de Benchmark

Combinación lineal de píxeles

## NN

### Versión sin Dropout

Capa de entrada

+

Capa intermedia con 4 neuronas y ReLU

+

Capa Salida (Lineal)

Batch = 64    Learning Rate = 0.001

### Versión con Dropout

Dropout en capa de entrada

Learning Rate = 0.0001

# Modelos desarrollados

timeSteps = 2, 12, 24

timeSteps = 12

## RNN

### Versión sin Dropout

Capa de entrada  
+  
Capa intermedia con 20 neuronas (ReLU +  
tanh) y regularizador L2  
+  
Capa Salida (Lineal)

Batch = 64    Learning Rate = 0.01

### Versión con Dropout

Dropout de acuerdo al esquema

timeSteps = 2

## LSTM

### Versión sin Dropout

Capa de entrada  
+  
Capa intermedia con 20 neuronas (ReLU +  
tanh)  
+  
Capa Salida (Lineal)

Batch = 64    Learning Rate = 0.01

### Versión con Dropout

Dropout de acuerdo al esquema

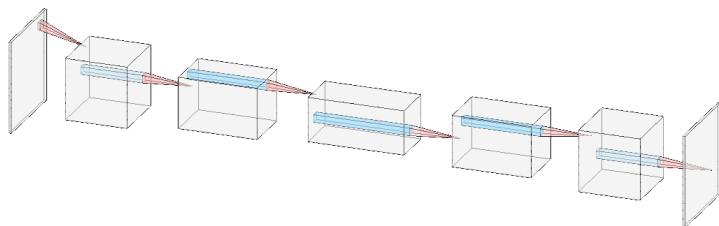
# Modelos desarrollados

TimeSteps = 1, 2, 12, 24

timeSteps = 2

## Conv AE

Versión sin Dropout



Activación ReLU  
Filtros 3x3  
12 24 48 24 12

Batch = 64    Learning Rate = 0.01

timeSteps = 2, 12, 24

timeSteps = 12

## Conv LSTM

Versión sin Dropout

Capa de entrada  
+  
Capa intermedia con 12 filtros  
+  
Capa salida con 1 filtro (Lineal)

Batch = 64    Learning Rate = 0.001

Versión con Dropout

Dropout de acuerdo al esquema

# Resultados

# Resultados

Tabla 5.1: Comparación en términos de MSPE y CRPS de los distintos modelos ajustados. En negrita los de referencia. La  $t$  hace referencia al número de meses utilizados como recurrencia en la capa de entrada. El  $DP$  es el valor del hiperparámetro  $p$  del dropout.

Modelo	Niño 3.4 MSPE	Niño 3.4 CRPS
<b>BAST-RNN</b>	<b>0.193</b>	<b>0.290</b>
NN	0.241	-
LSTM $t=2$	0.259	-
<b>E-QESN</b>	<b>0.261</b>	<b>0.354</b>
RNN $t=12$	0.291	-
LIN	0.357	-
RNN $t=24$	0.4	-
Conv AE $t=2$	0.401	-
RNN $t=2$	0.427	-
LSTM $t=24$	0.438	-
Conv AE $t=1$	0.44	-
LSTM $t=2$ w/ DP=5 %	0.491	0.534
RNN $t=12$ w/ DP=5 %	0.503	0.469
<b>GQN</b>	<b>0.619</b>	<b>0.538</b>
Conv AE $t=12$	0.757	-
NN w/ DP=5 %	0.766	0.689
<b>Lin. DSTM</b>	<b>0.785</b>	<b>0.699</b>
LSTM $t=12$	0.839	-
Conv LSTM $t=12$	0.874	-
Conv AE $t=24$	0.936	-
Conv LSTM $t=2$	1.015	-
Conv LSTM $t=12$ w/ DP=5 %	1.408	0.94
Conv LSTM $t=24$	1.801	-

# Resultados

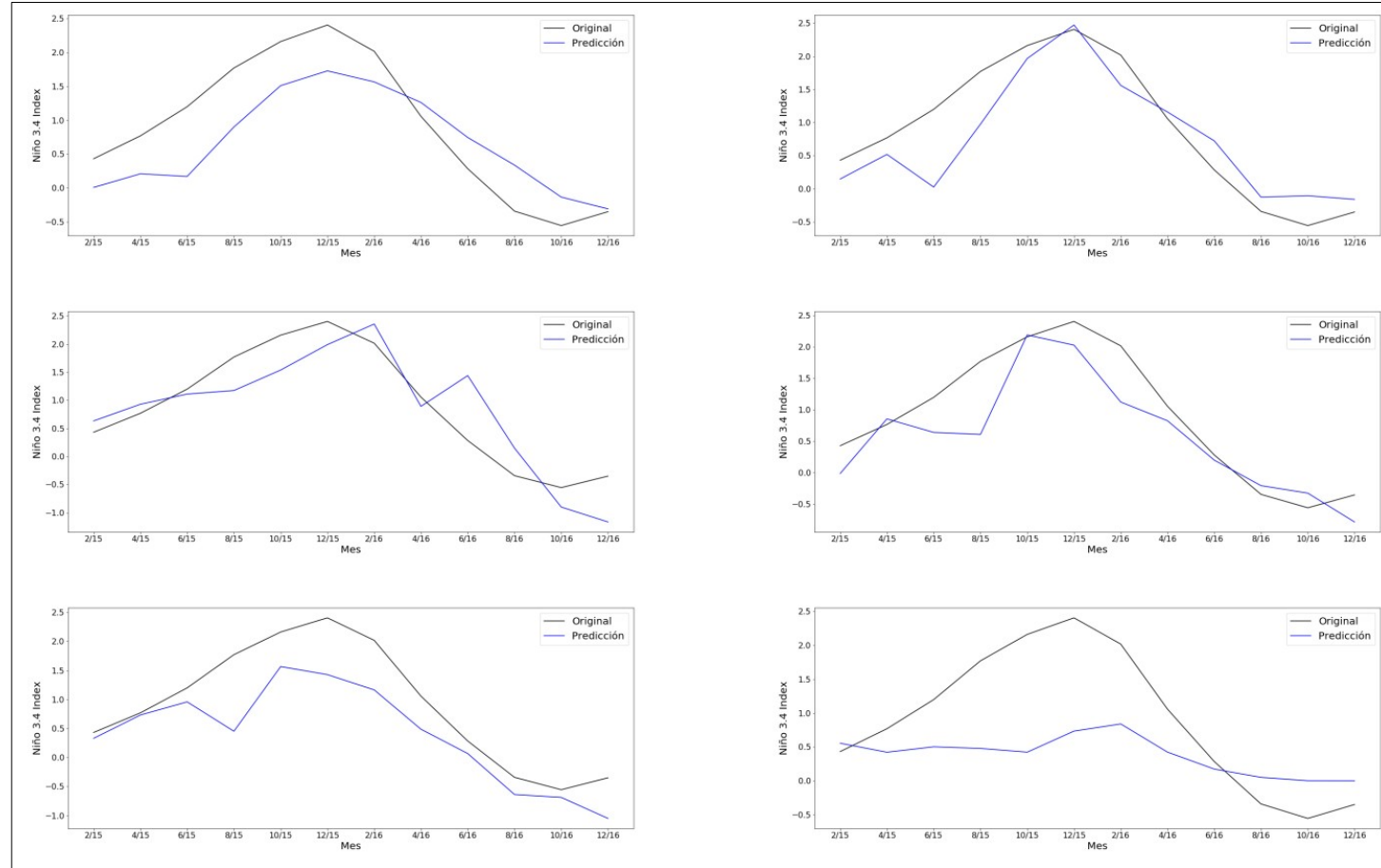
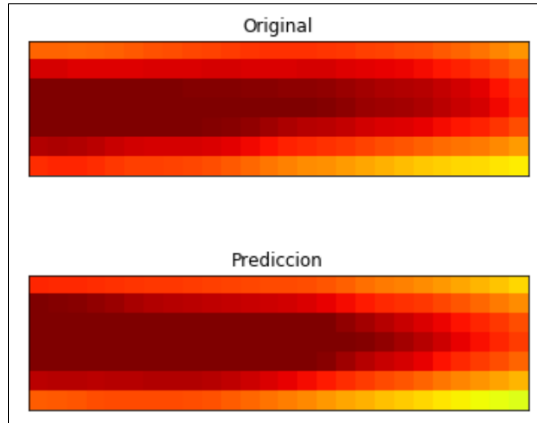
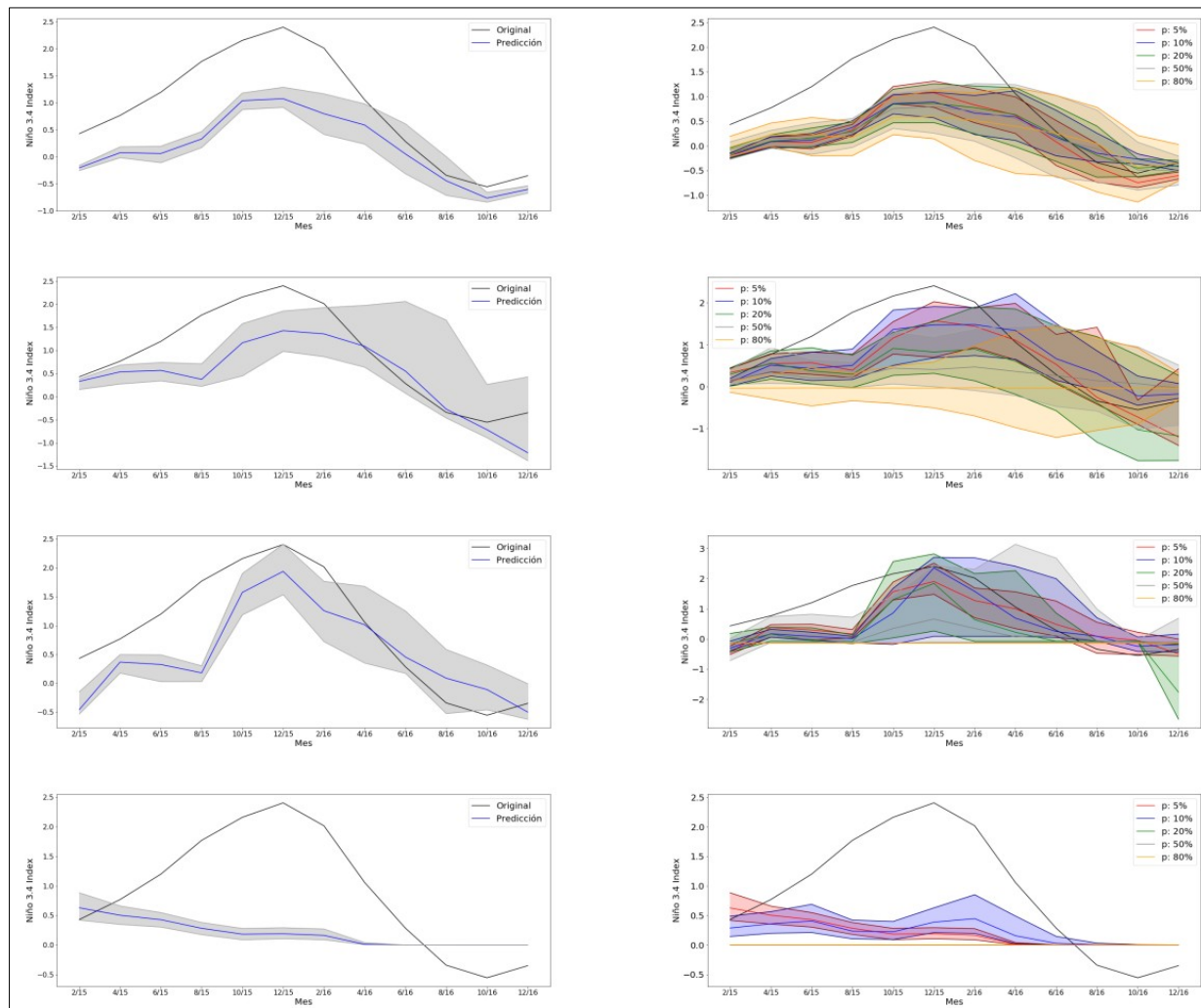


Figura 5.1: Predicciones respecto a los valores originales. En la columna de la izquierda de arriba a abajo se muestran las correspondientes a la red lineal, red recurrente y autoencoder convolucional. En la columna de la derecha de arriba a abajo se muestran las de la red densa, LSTM y LSTM convolucional.

# Resultados

Figura 5.3: En la columna de la izquierda de arriba a abajo se muestran los intervalos de confianza con  $p = 5\%$  para la red densa, red recurrente, LSTM y LSTM convolucional. En la columna de la derecha los intervalos de confianza para distintos valores de  $p$



# Resultados

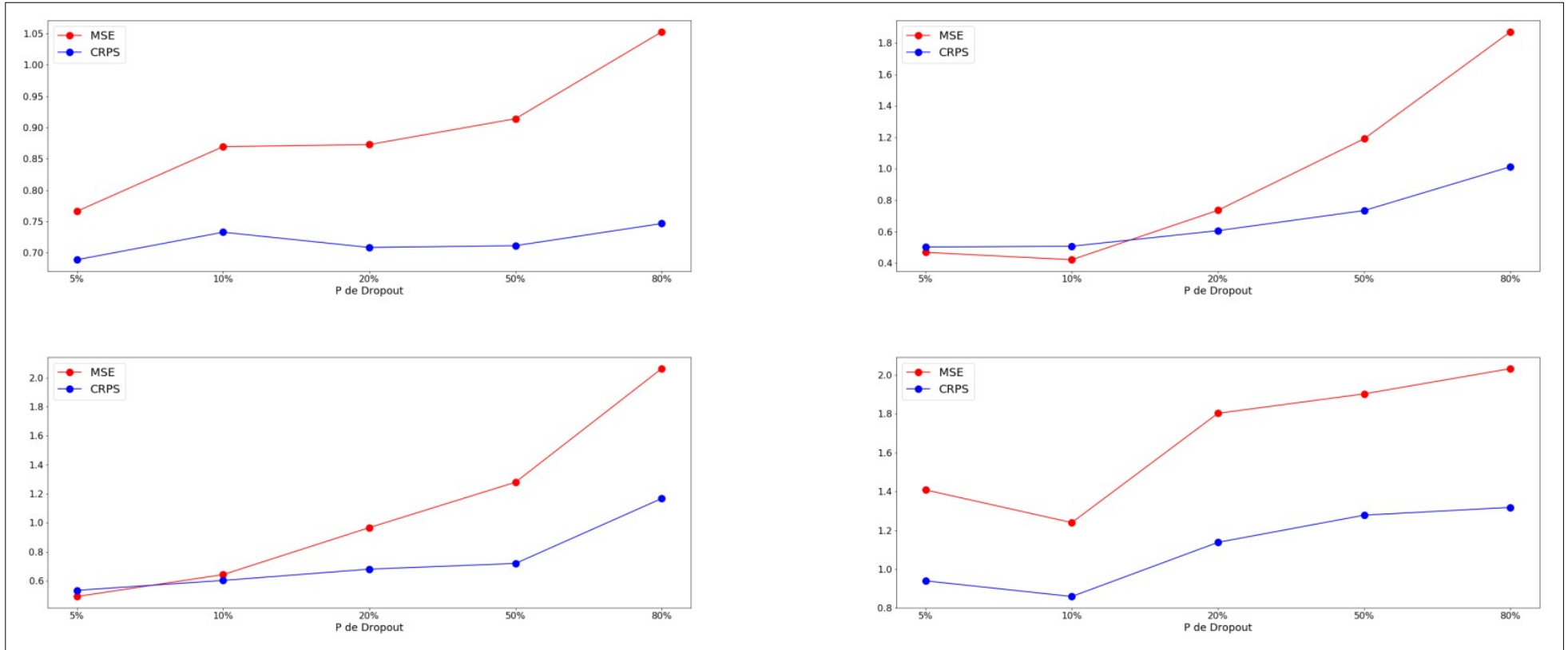


Figura 5.4: Evolución del MSE Y CRPS en función de la  $p$  del dropout. En la columna izquierda de arriba a abajo la correspondiente a la red densa y LSTM. En la columna derecha la de la red recurrente y LSTM convolucional



# Conclusiones

# Conclusiones

- Los modelos desarrollados superan al Benchmark lineal
- Los modelos sencillos son los que mejores resultados obtienen
- La red neuronal densa consigue por primera vez predecir correctamente la anomalía de Diciembre de 2015
- Los intervalos de confianza están muy condicionados al valor de  $p$
- Los mejores intervalos son los de los modelos recurrentes

Trade-off entre la capacidad de aprendizaje y la cuantificación de la incertidumbre

¿Preguntas?