SOCI 620: QUANTITATIVE METHODS 2

Agenda

Transformations & assessing fit with predictive plots

1. Administrative

- 2. Interpretation with transformed variables
- 3. Prior predictive plots
- 4. Visualizing model predictions
- 5. Hands on:
 - Visualizing model predictions in R

Transformations



MONTREAL BIKE TRAFFIC

287 7 117

24/01/25 16:43 re journ

Daily 2024 ridership at the 27 highesttraffic bike counters

3



INTERPRETING COEFFICIENTS

Predicting ridership using temperature



$$egin{aligned} R_i &\sim \operatorname{Norm}(\mu_i,\sigma)\ \mu_i &= lpha + eta T_i \end{aligned}$$
 $egin{aligned} extsf{Post.mean}\ lpha & extsf{878.28}\ eta & extsf{90.96} \end{aligned}$

If the temperature changes from t_1 to t_2 , how much do we expect ridership to change? $Note: \hat{\alpha} = E(\alpha) \text{ and } \hat{\beta} = E(\beta)$ $E(R_i|T_i = t_2, \mathbf{M}) - E(R_i|T_i = t_1, \mathbf{M}) = (\hat{\alpha} + \hat{\beta}t_2) - (\hat{\alpha} + \hat{\beta}t_1)$ $= \hat{\beta}(t_2 - t_1)$ $= 90.96(t_2 - t_1)$

INTERPRETING COEFFICIENTS

Predicting ridership using temperature



$$egin{aligned} R_i &\sim \operatorname{Norm}(\mu_i,\sigma)\ \mu_i &= lpha + eta T_i \end{aligned}$$
 $egin{aligned} extsf{Post.mean}\ lpha & extsf{878.28}\ eta & extsf{90.96} \end{aligned}$

Units of temperature : Degrees celsius

Units of ridership: Number of riders

Interpretation of β : For every increase of **one degree celsius**, the model predicts an average of **90.96 more riders** per day at each bike counter

STANDARDIZED VARIABLES

Predicting St(ridership) using St(temperature)



$\mathrm{St}(R_i) \sim \mathrm{Norm}(\mu_i, \sigma) \ \mu_i = lpha + eta \mathrm{St}(T_i)$		
	Post. mean	
lpha	0.0	
β	0.644	_
		_

Units of temperature : Standard deviations of temperature

Units of ridership: Standard deviations of ridership

Interpretation of : For every increase of **one standard deviation of temperature**, the model predicts an average of **0.644 more standard deviations of ridership** per day at each bike counter

LOG OF OUTCOME VARIABLE

Predicting Log(ridership) using temperature



$$\log(R_i) \sim \operatorname{Norm}(\mu_i, \sigma) \ \mu_i = lpha + eta T_i$$

	Post. mean	Exp(mean)
α	6.12	454.47
β	0.085	1.089

If the temperature changes from t_1 to t_2 , how much do we expect ridership to change?

$$egin{aligned} rac{\mathrm{E}(R_i|T_i=t_2,\mathbf{M})}{\mathrm{E}(R_i|T_i=t_1,\mathbf{M})} &= \exp\left(\log\left(rac{\mathrm{E}(R_i|T_i=t_2,\mathbf{M})}{\mathrm{E}(R_i|T_i=t_1,\mathbf{M})}
ight)
ight) \ &= \exp\left((\hatlpha+\hateta t_2)-(\hatlpha+\hateta t_1)
ight) \ &= \exp(\hateta(t_2-t_1)) \ &= \exp(0.085(t_2-t_1)) \end{aligned}$$

LOG OF OUTCOME VARIABLE



LOG OF OUTCOME VARIABLE

Predicting Log(ridership) using temperature



$$\log(R_i) \sim \operatorname{Norm}(\mu_i, \sigma)$$
 $\mu_i = lpha + eta T_i$
Post mean Exp(mean)

	Post. mean	Exp(mean)
lpha	6.12	454.47
β	0.085	1.089

Units of temperature : Degrees celsius

Units of ridership: Log ridership

Interpretation of β : For every increase of **one degree celsius**, the model predicts an average increase of **8.9% in ridership** per day at each bike counter

Prior and posterior predictive plots



 $egin{aligned} \log(R_i) &\sim \operatorname{Norm}(\mu_i, \sigma) \ \mu_i &= lpha + eta \operatorname{St}(T_i) \end{aligned}$

 $egin{aligned} lpha &\sim \operatorname{Norm}(10,5) \ eta &\sim \operatorname{Norm}(0,3) \ \sigma &\sim \operatorname{Unif}(0,5) \end{aligned}$

Prior predictive plots allow you to visualize the implications of a set of priors

 $egin{aligned} \log(R_i) &\sim \operatorname{Norm}(\mu_i, \sigma) \ \mu_i &= lpha + eta \operatorname{St}(T_i) \end{aligned}$

 $egin{aligned} lpha &\sim \operatorname{Norm}(10,5) \ eta &\sim \operatorname{Norm}(0,3) \ \sigma &\sim \operatorname{Unif}(0,5) \end{aligned}$



 $\frac{\alpha}{12.45} \frac{\beta}{2.58}$

 $egin{aligned} \log(R_i) &\sim \operatorname{Norm}(\mu_i, \sigma) \ \mu_i &= lpha + eta \operatorname{St}(T_i) \end{aligned}$



lpha	eta
12.45	2.58
8.01	0.51

 $egin{aligned} \log(R_i) &\sim \operatorname{Norm}(\mu_i, \sigma) \ \mu_i &= lpha + eta \operatorname{St}(T_i) \end{aligned}$



α	eta
12.45	2.58
8.01	0.51
12.55	1.42

 $egin{aligned} \log(R_i) &\sim \operatorname{Norm}(\mu_i, \sigma) \ \mu_i &= lpha + eta \operatorname{St}(T_i) \end{aligned}$



α	eta
12.45	2.58
8.01	0.51
12.55	1.42
8.01	-3.34

 $egin{aligned} \log(R_i) &\sim \operatorname{Norm}(\mu_i, \sigma) \ \mu_i &= lpha + eta \operatorname{St}(T_i) \end{aligned}$



α	eta
12.45	2.58
8.01	0.51
12.55	1.42
8.01	-3.34
18.19	3.85

 $egin{aligned} \log(R_i) &\sim \operatorname{Norm}(\mu_i, \sigma) \ \mu_i &= lpha + eta \operatorname{St}(T_i) \end{aligned}$



α	β
12.45	2.58
8.01	0.51
12.55	1.42
8.01	-3.34
18.19	3.85
13.11	3.24
11.01	0.66
15.54	1.36
8.97	-4.10
8.11	-0.42
8.48	3.76
•	•

VISUALIZING PREDICTIONS



To better illustrate uncertainty in posterior estimates, we will use a random subsample of 400 counter-days.

 $egin{aligned} \log(R_i) &\sim \operatorname{Norm}(\mu_i, \sigma) \ \mu_i &= lpha + eta \operatorname{St}(T_i) \end{aligned}$

JAI I7ING PREDIC

 σ



3. Calculate quantiles (say, 10% and 90%) of the posterior draws $\{\mu_n\}$ at each value of t.

VISUALIZING PREDICTIONS

lo	$egin{aligned} & arphi(R_i) &\sim \ & \mu_i = \ & lpha &\sim \ & eta &\sim \ & eta &\sim \ & \sigma &\sim \ & \sigma &\sim \end{aligned}$	$egin{array}{l} \operatorname{Norm}(\mu_i,\sigma)\ lpha+eta \operatorname{St}(T_i) \end{array} \ Norm(10,5)\ \operatorname{Norm}(0,3)\ \operatorname{Unif}(0,5) \end{array}$	Provide the second seco
	Post. median	80% post. interval	Posterior predictive distribution: $\Pr(\log(R_i) T=t)$
α	6.877	(6.823, 6.931)	1. Take a sample of size N from posterior $\Pr(lpha, eta, \sigma D)$.
β	0.940	(0.887, 0.996)	2. For each value of standardized temperature t , calculate N values $\mu_n = \alpha_n + \beta_n t$.
σ	0.841	(0.804, 0.981)	3. For each of the N posterior draws for temperature t , draw $p_n \sim \mathrm{Norm}(\mu_n, \sigma_n)$.
			4. Calculate quantiles (say, 10% and 90%) of the posterior draws $\{p_n\}$ at each value of t .

POSTERIOR MEAN VS. PREDICTED



For any given mean temperature, μ is the value of log ridership that is "expected" by the model on a day with that temperature.

The *posterior distribution* of μ describes our uncertainty about the value of μ after seeing the data.

This distribution takes into account the model coefficients α and β , but *not* the model standard deviation σ .

The 80% posterior interval of μ should get narrower as more data is added.



For any given mean temperature, the *posterior predictive distribution* describes the range of riderships we would expect for any day with that temperature.

The posterior predictive distribution describes our uncertainty about the value of $\log(R_i)$ after seeing the data.

This distribution takes into account all the model parameters: α , β , and σ .

The 80% posterior *predictive* interval should contain about 80% of the data, and should not get appreciably narrower as more data is added.

POSTERIOR MEAN VS. PREDICTED



22

Image credit



Figures by Peter McMahan (<u>source</u> <u>code</u>)



Alfred Hitchcock in Funhouse Mirror by Globe Photos

Photo via Flickr user Midnight Believer